

**ANALISIS KATA KUNCI YANG RELEVAN DENGAN GAS
RUMAH KACA DI INDONESIA MENGGUNAKAN *TEXT*
MINING DAN TEORI LUHN**

LAPORAN TUGAS AKHIR

Oleh:

Kiagus Muhammad Arsyad

105219002



**FAKULTAS SAINS DAN ILMU KOMPUTER
PROGRAM STUDI ILMU KOMPUTER
UNIVERSITAS PERTAMINA
AGUSTUS 2023**

ANALISIS KATA KUNCI YANG RELEVAN DENGAN GAS RUMAH KACA DI INDONESIA MENGGUNAKAN *TEXT* *MINING* DAN TEORI LUHN

LAPORAN TUGAS AKHIR

Oleh:

Kiagus Muhammad Arsyad

105219002



**FAKULTAS SAINS DAN ILMU KOMPUTER
PROGRAM STUDI ILMU KOMPUTER
UNIVERSITAS PERTAMINA
AGUSTUS 2023**



LEMBAR PENGESAHAN

Judul Tugas Akhir : Analisis Kata Kunci yang Relevan dengan Gas
Rumah Kaca di Indonesia Menggunakan *Text Mining* dan Teori Luhn
Nama Mahasiswa : Kiagus Muhammad Arsyad
Nomor Induk Mahasiswa : 105219002
Program Studi : Ilmu Komputer
Fakultas : Sains dan Ilmu Komputer
Tanggal Lulus Sidang Tugas Akhir : 2 Agustus 2023

Jakarta, 31 Juli 2023

MENGESAHKAN

Pembimbing I

Pembimbing II

Dr. Tasmi, S.Si., M.Si.
NIP. 116109

Dr. Ariana Yunita
NIP. 116015

MENGETAHUI,

Ketua Program Studi

Ade Irawan, Ph.D

NIP. 116130

LEMBAR PERNYATAAN

Dengan ini saya menyatakan bahwa Tugas Akhir berjudul Analisis Kata Kunci yang Relevan dengan Gas Rumah Kaca di Indonesia Menggunakan *Text Mining* dan Teori Luhn ini adalah benar-benar merupakan hasil karya saya sendiri dan tidak mengandung materi yang ditulis oleh orang lain kecuali telah dikutip sebagai referensi yang sumbernya telah dituliskan secara jelas sesuai dengan kaidah penulisan karya ilmiah.

Apabila dikemudian hari ditemukan adanya kecurangan dalam karya ini, saya bersedia menerima sanksi dari Universitas Pertamina sesuai dengan peraturan yang berlaku.

Demi pengembangan ilmu pengetahuan, saya menyetujui untuk memberikan kepada Universitas Pertamina hak bebas royalti noneksklusif (*non-exclusive royalty-free right*) atas Tugas Akhir ini beserta perangkat yang ada. Dengan hak bebas royalti noneksklusif ini Universitas Pertamina berhak menyimpan, mengalih media/format-kan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan Tugas Akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Jakarta, 31 Juli 2023

Yang membuat pernyataan,

– Materai –

Kiagus Muhammad Arsyad

ABSTRAK

Kiagus Muhammad Arsyad. 105219002. Analisis Kata Kunci yang Relevan dengan Gas Rumah Kaca di Indonesia Menggunakan *Text Mining* dan Teori Luhn.

Penyebab Gas Rumah Kaca (GRK) menjadi isu global yang mendesak untuk ditangani, termasuk di Indonesia. Penelitian ini bertujuan untuk mengidentifikasi faktor penyebab GRK di Indonesia melalui analisis data abstrak dari studi literatur lima tahun terakhir. Pendekatan yang digunakan melibatkan teknik *text mining*, *clustering*, dan seleksi menggunakan teori Luhn. Dengan menggunakan algoritma K-Means dalam teknik *clustering*, data abstrak dikategorikan menjadi beberapa kluster berdasarkan kesamaan konten. Hasil analisis menunjukkan bahwa aspek yang memiliki relevansi GRK di Indonesia meliputi diantaranya konsumsi listrik, produksi energi, penggunaan energi, dan frasa yang lainnya. Metode seleksi menggunakan teori Luhn juga mampu mengidentifikasi kata-kata kunci yang paling relevan dan mewakili relevansi GRK di Indonesia. Pada penelitian ini juga memberikan pemahaman yang cukup tentang faktor-faktor yang berhubungan pada aspek GRK di Indonesia dan memberikan informasi penting bagi pembuat kebijakan, para ahli, para pembaca untuk mengembangkan strategi mitigasi yang efektif.

Kata kunci: faktor, gas rumah kaca, kata kunci, nlp, teori luhn, *text mining*

ABSTRACT

Kiagus Muhammad Arsyad. 105219002. Analysis of Relevant Keywords Related to Greenhouse Gas in Indonesia Using Text Mining and Luhn's Theory.

The causes of Greenhouse Gas (GHG) have become a pressing global issue to be addressed, including in Indonesia. This study aims to identify the factors causing GHG in Indonesia through the analysis of abstract data from literature studies in the last five years. The approach used involves text mining, clustering, and selection techniques using Luhn's theory. By employing the K-Means algorithm in the clustering technique, the abstract data is categorized into several clusters based on content similarity. The results of the analysis show that aspects relevant to GHG in Indonesia include electricity consumption, energy production, energy use, and other phrases. The selection method using Luhn's theory also successfully identifies the most relevant keywords representing the relevance of GHG in Indonesia. This study provides a comprehensive understanding of the factors related to GHG aspects in Indonesia and offers important information for policymakers, experts, and readers to develop effective mitigation strategies.

Keywords: factors, greenhouse gas, keywords, luhn's theory, nlp, text mining

KATA PENGANTAR

Dengan memanjatkan segala puji syukur kehadirat Allah SWT yang telah melimpahkan rahmat, taufik, dan hidayah-Nya sehingga penulis dapat menyelesaikan skripsi berupa laporan Tugas Akhir (TA) yang berjudul **“Analisis Kata Kunci Faktor Penyebab Gas Rumah Kaca di Indonesia Melalui Pendekatan Kuantitatif dan Kualitatif pada Studi Literatur dengan *Text Mining* dan Teori Luhn”**, sebagai salah satu persyaratan dalam menyelesaikan masa belajar di Program Studi S1 Ilmu Komputer, Fakultas Sains dan Ilmu Komputer Universitas Pertamina. Penulis menyadari bahwa laporan TA ini tidak akan terselesaikan tanpa adanya dukungan, bimbingan, bantuan, dan nasihat dari berbagai pihak selama penyusunan laporan TA ini. Pada kesempatan ini, penulis ingin menyampaikan terima kasih setulus-tulusnya kepada:

1. Bapak Ade Irawan, Ph.D selaku Ketua Program Studi Ilmu Komputer Universitas Pertamina dan juga sebagai dosen penguji pertama yang telah memberikan masukan dan saran pada skripsi ini.
2. Bapak Rangga Ganzar Noegraha, Ph.D selaku dosen penguji kedua yang telah memberikan masukan dan saran pada skripsi ini.
3. Ibu Dr. Tasmi, S.Si, M.Si. dan Ibu Dr. Ariana Yunita selaku dosen pembimbing pertama dan dosen pembimbing kedua pada skripsi ini yang telah memberikan bimbingan, motivasi, dan arahan selama penyusunan laporan tugas akhir dari awal hingga akhir saat ini.
4. Seluruh staf pengajar yang telah memberikan ilmu pengetahuan yang tak ternilai selama penulis menempuh pendidikan di Universitas Pertamina.
5. Seluruh *civitas academica* di lingkungan Universitas Pertamina yang telah menunjang layanan dan dukungan selama proses pembelajaran di Universitas Pertamina.
6. Keluarga, saudara, dan teman-teman penulis yang telah memberikan doa, kasih sayang, motivasi, dan dorongan di setiap langkah perjuangan penulis.

Penulis menyadari bahwa laporan TA ini belum sempurna dan banyak kekurangan, oleh karena itu, segala saran dan kritik yang membangun diharapkan oleh penulis sebagai penyempurnaan penulisan laporan ini di masa yang mendatang. Semoga laporan TA ini dapat memberikan manfaat bagi penulis dan para pembaca.

Jakarta, 31 Juli 2023

Kiagus Muhammad Arsyad

DAFTAR ISI

KATA PENGANTAR	iii
DAFTAR ISI	iv
DAFTAR TABEL	vi
DAFTAR GAMBAR	vii
BAB I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	2
1.3 Batasan Masalah	2
1.4 Tujuan penelitian	2
1.5 Manfaat penelitian	2
BAB II TINJAUAN PUSTAKA	4
2.1 Gas Rumah Kaca (GRK)	4
2.2 <i>Text Mining</i>	5
2.2.1 <i>Natural Language Processing (NLP)</i>	5
2.2.2 <i>Term Frequency-Inverse Document Frequency (TF-IDF)</i>	6
2.2.3 <i>Principal Component Analysis (PCA)</i>	7
2.2.4 K-Means	7
2.3 Teknik Kitchenham	7
2.4 <i>Luhn's Theory</i>	8
BAB III METODE PENELITIAN	11
3.1 Tahapan Penelitian	11
3.2 <i>Data Collection</i>	13
3.3 Tahapan <i>Text Mining</i> dan Interpretasinya	14
BAB IV HASIL DAN PEMBAHASAN	16

4.1	<i>Preprocessing Data</i>	16
4.2	Hasil Pengelompokkan Kata Kunci berdasarkan teknik <i>clustering</i>	17
4.3	Hasil Pemilihan Kata Kunci berdasarkan Teori Luhn	21
4.4	Hasil Kata Kunci yang Didapatkan Berdasarkan Teknik <i>Clustering</i> dan Seleksi Teori Luhn	24
BAB V	KESIMPULAN DAN SARAN	27
5.1	Kesimpulan	27
5.2	Saran	27
	DAFTAR PUSTAKA	29
LAMPIRAN A	Lampiran 1. Tren Total Emisi GRK Indonesia	33
LAMPIRAN B	Lampiran 2. Frasa dari Kata Kunci yang Didapatkan	34

DAFTAR TABEL

Tabel 4.1	Hasil Kata Kunci yang Terpilih Berdasarkan Teknik <i>Clustering</i> dan Implementasi Teori Luhn	25
-----------	---	----

DAFTAR GAMBAR

Gambar 2.1	Diagram <i>Word Frequency</i>	9
Gambar 3.1	Diagram Alir Penelitian	11
Gambar 3.2	Tahapan <i>Data Collection</i> hingga hasil interpretasi	12
Gambar 3.3	Tahapan proses pemilihan data studi literatur	13
Gambar 3.4	Alur tahapan saat dilakukan <i>text mining</i> hingga hasil interpretasi	14
Gambar 4.1	Proses mendapatkan kata-kata dasar dari data studi literatur	16
Gambar 4.2	Visualisasi hasil dari PCA	18
Gambar 4.3	<i>Elbow Plot</i> untuk menentukan nilai klaster dengan K-Means	18
Gambar 4.4	Metrik Evaluasi nilai K dengan <i>Davies-Bouldin index</i> dan <i>library ElbowVisualizer</i>	19
Gambar 4.5	Klasterisasi Pewarnaan dengan K-Means	20
Gambar 4.6	<i>Silhouette Plot</i> untuk K-Means <i>Clustering</i>	20
Gambar 4.7	Representasi visual frekuensi data dengan <i>word cloud</i>	22
Gambar 4.8	Dua puluh kata dasar dengan frekuensi terbesar	22
Gambar 4.9	<i>Preview</i> kata-kata dengan frekuensi terbesar	23
Gambar 4.10	<i>Preview</i> kata-kata dengan frekuensi terendah	23
Gambar 4.11	Hasil Kata Kunci Berdasarkan Hasil Analisis Manual dan Implementasi Teori Luhn	24
Gambar 1.1	Total emisi GRK di Indonesia 2010-2019	33
Gambar 2.1	Hasil frasa yang dapat dibentuk (1)	34
Gambar 2.2	Hasil frasa yang dapat dibentuk (2)	35



BAB I

PENDAHULUAN

1.1 Latar Belakang

Emisi gas rumah kaca (GRK) di Indonesia diperkirakan meningkat pada periode 2021-2030. Informasi tersebut diperoleh dari artikel DataIndonesia.Id yang ditulis oleh Rizaty [1], pada artikel tersebut juga disebutkan bahwa emisi GRK nasional sudah mencapai 259,1 juta ton CO₂ pada tahun 2021. Proyeksi emisi GRK tahun 2030 diprediksi akan meningkat sebesar 29,13% menjadi 334,6 juta ton CO₂.

Selain itu, Indonesia akan terus meningkat diikuti dengan kemajuan teknologi, hal ini menyebabkan peningkatan kebutuhan energi [2]. Pertumbuhan penduduk ini berdampak pada penggunaan bahan bakar fosil, seperti pembakaran kendaraan bermotor dan kegiatan industri, yang menjadi penyumbang salah satu faktor emisi GRK [3]. Dampak dari peningkatan emisi GRK dan konsumsi energi di dunia juga sangat signifikan terhadap lingkungan, seperti kenaikan suhu global, perubahan iklim ekstrem, serta perubahan pola cuaca [4]. Hal ini perlu diwaspadai mengingat dampak dari emisi GRK yang cukup membahayakan.

Di Indonesia, terdapat program yang dinamakan Indonesia's FOLU Net Sink 2030, yang merupakan sebuah inisiatif dari pemerintah Indonesia untuk mencapai keseimbangan antara emisi dan penyerapan karbon pada tahun 2030. Dalam rangka mengatasi isu perubahan iklim, Indonesia berkomitmen dalam memperkuat sektor kehutanan, pertanian, dan penggunaan lahan lainnya sebagai upaya mencapai "*net sink*" yaitu penyerapan karbon yang lebih banyak [5].

Inovasi, upaya, dan ide yang dihasilkan dari penelitian yang membahas tentang emisi gas rumah kaca (GRK) di Indonesia juga diperlukan saat ini dalam memberikan kontribusi yang lebih efektif dalam mencapai target Indonesia's FOLU Net Sink 2030. Dalam rangka mencari solusi terhadap masalah emisi GRK di Indonesia, pemanfaatan teknologi dan kecerdasan buatan, seperti *Natural Language Processing* (NLP) dan *text mining*, dapat membantu mengatasi tantangan tersebut [6, 7]. Teks-teks dari berbagai sumber, seperti jurnal penelitian, artikel, dan laporan yang membahas mengenai GRK dan terdapat hubungannya di Indonesia dapat dijadikan sebagai data. Data teks tersebut dapat diolah dan dianalisis dengan efisien menggunakan proses *text mining* dan NLP dalam mengidentifikasi faktor-faktor penyebab GRK yang relevan. Selain itu, terdapat pendekatan menggunakan teori Luhn [8] sebagai *text summarizer* [9] atau dalam pemilihan kata kunci sesuai dengan konsep yang ingin dicari [10], sehingga dapat meningkatkan kualitas dan validitas hasil dari proses tersebut.

Oleh sebab itu, penulis ingin menerapkan pendekatan gabungan berupa kuantitatif dan kualitatif yang dapat menentukan kata-kata kunci yang relevan mengenai GRK di Indonesia berdasarkan data yang tersedia dengan didukung pencarian studi literatur. Diharapkan dari penelitian ini dapat memberikan pengetahuan baru dan ide terkait solusi yang dapat dilakukan untuk mengurangi emisi GRK. Selain itu juga dapat mendukung upaya kepada pembaca, para pakar, peneliti, instansi maupun peme-

rintah dalam mengoptimalkan kebijakan dan strategi untuk mengatasi perubahan iklim secara lebih efisien dan efektif khususnya dalam aspek GRK di Indonesia.

1.2 Rumusan Masalah

Berdasarkan latar belakang yang telah dijelaskan, rumusan masalah pada penelitian ini yaitu bagaimana cara menentukan kata-kata kunci yang memiliki relevansi faktor penyebab gas rumah kaca di Indonesia berdasarkan data studi literatur dengan *text mining* berupa teknik *clustering* dan teori Luhn?

1.3 Batasan Masalah

Adapun batasan masalah pada penelitian ini ialah menggunakan data berupa studi literatur yang merupakan artikel atau penelitian lima tahun terakhir sejak tahun 2018 yang didapatkan melalui *search engine* Scopus yang mengulas terkait GRK di negara Indonesia.

1.4 Tujuan penelitian

Tujuan dari penelitian ini yaitu untuk menentukan kata-kata kunci yang memiliki relevansi faktor penyebab gas rumah kaca di Indonesia berdasarkan data studi literatur dengan *text mining* berupa teknik *clustering* dan teori Luhn.

1.5 Manfaat penelitian

Penelitian ini diharapkan dapat memberikan manfaat yang signifikan dalam menentukan faktor penyebab emisi Gas Rumah Kaca (GRK) di Indonesia berdasarkan keterbaruan dari salah satu basis data literatur ilmiah yaitu Scopus. Dengan menerapkan pendekatan gabungan yaitu kuantitatif dan kualitatif menggunakan *text mining* dengan NLP serta K-Means dan teori Luhn, penelitian ini diharapkan juga dapat memberikan kemudahan dan efisiensi dalam mengidentifikasi kata kunci yang relevan mengenai GRK secara valid berdasarkan fakta dan data yang terverifikasi. Selain itu, hasil penelitian ini dapat menjadi dasar untuk mengembangkan kebijakan dan strategi yang lebih baik berdasarkan aspek relevansi dalam menghadapi tantangan perubahan iklim dan menjaga keberlanjutan lingkungan bagi masa depan yang lebih baik terutama dalam kasus GRK di Indonesia.



BAB II

TINJAUAN PUSTAKA

2.1 Gas Rumah Kaca (GRK)

Perubahan iklim telah menjadi isu keamanan internasional yang penting dan nontradisional. Emisi gas rumah kaca (GRK) yang berlebihan merupakan salah satu penyebab utama dari meningkatnya permasalahan iklim. Perubahan iklim tidak hanya secara langsung mempengaruhi kesehatan populasi dengan meningkatkan frekuensi dan intensitas gelombang panas, kekeringan, dan curah hujan yang tinggi, tetapi juga secara tidak langsung dengan meningkatkan polusi udara, mempercepat penyebaran vektor penyakit, serta mempengaruhi keamanan pangan dan kesehatan mental. Emisi GRK yang berlebihan menjadi faktor yang memperburuk perubahan iklim dan dampaknya terhadap berbagai aspek kehidupan manusia [11].

Gas rumah kaca merujuk pada gas-gas yang hadir di atmosfer, baik secara alami maupun sebagai hasil aktivitas manusia (antropogenik), yang mampu menyerap dan memancarkan kembali radiasi inframerah [12]. Ketika permukaan bumi menerima radiasi matahari dalam bentuk gelombang pendek, sebagian besar radiasi ini dipancarkan kembali ke atmosfer sebagai radiasi gelombang panjang (inframerah). Gas rumah kaca yang terdapat di lapisan atmosfer yang dekat dengan permukaan bumi menyerap radiasi gelombang panjang ini, menyebabkan peningkatan suhu yang tinggi yang dikenal sebagai efek rumah kaca. Peningkatan suhu ini terjadi akibat perubahan kondisi dan komposisi atmosfer yang mengelilingi planet ini [13].

Dalam era saat ini, masalah lingkungan telah menjadi pembahasan utama baik di negara-negara yang sedang berkembang maupun negara-negara maju karena adanya kerusakan lingkungan. Hal ini juga menimbulkan pertanyaan tentang pemanasan global dan perubahan iklim, yang terutama disebabkan oleh emisi gas rumah kaca, kadang-kadang terkait dengan penyebab alami seperti pergeseran benua, aktivitas gunung berapi, radiasi matahari, dan arus laut, serta aktivitas manusia langsung maupun tidak langsung yang mempengaruhi komposisi atmosfer global dan variasi lingkungan [4]. Para peneliti telah berargumen bahwa peningkatan aktivitas manusia akibat perkembangan industrialisasi, pertumbuhan populasi global, dan kebutuhan untuk mengatasi perubahan tersebut adalah penyebab utama perubahan iklim. Selain itu, aktivitas manusia seperti deforestasi pertanian dan komersial, pembakaran bahan bakar fosil, serta perubahan penggunaan lahan akibat pertumbuhan populasi juga memberikan kontribusi yang signifikan terhadap peningkatan emisi gas rumah kaca [14].

Menurut Khairunnisa Musari & Sayah [15], dalam rangka penyelesaian masalah gas rumah kaca (GRK), beberapa hal yang perlu diperhatikan adalah upaya melawan perubahan iklim, prioritas nasional, transformasi kebijakan, menciptakan lingkungan yang mendukung, dan investasi keuangan yang diperlukan, yang semuanya harus menjadi bagian dari agenda nasional. Selain itu dengan perkembangan era informasi saat ini, salah satu upaya dengan pemanfaatan teknologi atau kecerdasan buatan (*Artificial Intelligence/AI*) menjadi hal yang penting karena dapat digunakan untuk pemantauan, analisis, dan pengelolaan data terkait emisi GRK. Sementara itu, dapat membantu dalam

mengoptimalkan kebijakan dan strategi untuk mengurangi emisi GRK secara efisien. Dengan memanfaatkannya secara holistik, diharapkan dapat memberikan penyelesaian masalah GRK dan perubahan iklim dapat tercapai dengan lebih efektif dan efisien.

2.2 Text Mining

Studi literatur mengenai *text mining* telah menjadi area penelitian yang semakin populer dalam ilmu data dan pengolahan bahasa alami. *Text mining* merupakan teknik yang digunakan untuk mengekstraksi informasi berharga dari data teks yang tidak terstruktur, seperti artikel jurnal, berita, dan dokumen teks lainnya. Teknik ini mencakup beberapa tahapan, termasuk pengumpulan data teks, pemrosesan teks untuk menghapus karakter-karakter tidak penting, dan analisis data untuk mengidentifikasi pola, topik, atau informasi penting dari teks tersebut [7].

Dalam penerapannya, *text mining* telah digunakan dalam berbagai bidang. Misalnya, dalam analisis sentimen, *text mining* dapat digunakan untuk mengekstraksi sentimen atau perasaan dari teks yang ditulis oleh pengguna dalam media sosial atau ulasan produk [16]. Selain itu, dalam pengelompokan topik, *text mining* dapat membantu mengelompokkan dokumen-dokumen teks berdasarkan tema atau topik yang sama.

Dalam penelitian ini, *text mining* akan digunakan sebagai salah satu teknik dalam *preprocessing dataset* berupa teks untuk mendapatkan kata kunci yang relevan dengan permasalahan gas rumah kaca (GRK) di Indonesia. Teknik *text mining* akan membantu dalam ekstraksi fitur teks dari data teks yang tidak terstruktur, yakni **judul dan** abstrak jurnal penelitian yang terdapat pada *dataset*. Proses *text mining* akan melibatkan langkah-langkah seperti dengan teknik *stopwords*, *lemmatization*, *punctuation*, *word frequency* dari ekstraksi *unigram* yang akan menjadi kata-kata dasar dari *dataset* tersebut.

2.2.1 Natural Language Processing (NLP)

Bahasa Pemrosesan Alami atau yang biasa disebut *Natural Language Processing* (NLP) adalah bidang keilmuan yang berfokus pada pemahaman dan penerapan bahasa manusia dalam bentuk teks atau ucapan oleh komputer. NLP telah mengalami perkembangan pesat dalam beberapa dekade terakhir, berkat kemajuan teknologi komputasi dan kemampuan mesin dalam memproses data secara lebih efektif. Salah satu tantangan utama dalam NLP adalah mengenali dan memahami struktur dan makna dari teks bahasa manusia, karena bahasa manusia seringkali ambigu dan penuh dengan variasi. Beberapa aplikasi praktis dari NLP termasuk mesin penerjemah, asisten virtual, analisis sentimen, pengenalan suara, dan masih banyak lagi. NLP telah menjadi bidang penelitian yang sangat menarik dan relevan dalam era informasi digital yang sedang berkembang pesat [6].

Salah satu pendekatan populer dalam NLP adalah menggunakan teknik pembelajaran mesin (*machine learning*) untuk mengajarkan komputer bagaimana memahami dan memproses bahasa manusia. Teknik-teknik ini melibatkan penggunaan model statistik dan algoritma yang kompleks untuk mengidentifikasi pola, keterkaitan, dan makna dalam teks bahasa manusia. Contoh teknik pembelajaran mesin yang sering digunakan dalam NLP termasuk model berbasis aturan, model berbasis vektor, dan model berbasis jaringan saraf tiruan [17]. Dalam beberapa tahun terakhir, penggunaan jaringan

saraf tiruan, khususnya dalam bentuk transformasi bahasa seperti BERT dan GPT-3, telah mencatat kemajuan signifikan dalam banyak tugas NLP, seperti penerjemahan mesin dan pemahaman bahasa alami [18].

Dalam penelitian ini, teknik pemrosesan bahasa alami (Natural Language Processing/NLP) akan digunakan dalam tahapan *text mining* pada data berupa teks untuk mendapatkan kata kunci yang relevan. Proses pada tahapan tersebut meliputi *preprocessing* dengan menerapkan langkah-langkah seperti menghapus kata umum dan tanda baca, lematisasi, serta mendapatkan kata-kata yang relevan dari teks. Selain itu, teknik ekstraksi fitur seperti metode Term Frequency-Inverse Document Frequency (TF-IDF) digunakan untuk mengidentifikasi kata-kata yang paling penting dan mewakili konten teks secara numerik [19]. Hasil dari penggunaan NLP dalam pemrosesan data teks ini akan digunakan sebelum menentukan hasil interpretasi yang diharapkan pada penelitian ini.

2.2.2 Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF (*Term Frequency-Inverse Document Frequency*) merupakan salah satu teknik yang sering digunakan dalam pengolahan teks dan *Information Retrieval* (IR). Teknik ini berfungsi untuk memberikan bobot pada setiap kata dalam dokumen berdasarkan frekuensi kemunculan kata tersebut dalam dokumen itu sendiri (TF) dan kebalikannya, yakni berdasarkan keberbedaan jumlah dokumen yang mengandung kata tersebut di dalam seluruh koleksi dokumen (IDF) [20]. Dengan menggunakan TF-IDF, kata-kata yang sering muncul dalam suatu dokumen tetapi jarang muncul di dokumen lain akan memiliki bobot yang tinggi, sehingga dapat dianggap sebagai kata kunci atau kata yang lebih penting dalam konteks dokumen tersebut. Penggunaan TF-IDF telah terbukti efektif dalam berbagai aplikasi seperti klasifikasi teks, analisis sentimen, dan sistem rekomendasi, serta telah menjadi salah satu teknik yang populer dalam mengolah teks untuk mendapatkan informasi yang bermakna [21].

Adapun TF-IDF *Vectorizer* adalah salah satu fungsi yang digunakan dalam pengolahan teks untuk mengimplementasikan metode TF-IDF. TF-IDF *Vectorizer* berfungsi untuk mengubah teks menjadi vektor numerik berdasarkan nilai TF-IDF dari setiap kata dalam teks [21]. Proses ini memungkinkan data teks yang kompleks dan tidak terstruktur diubah menjadi representasi numerik yang dapat diolah lebih lanjut dengan algoritma pembelajaran mesin atau analisis statistik.

Penggunaan TF-IDF *Vectorizer* sangat berguna dalam analisis teks dan NLP, terutama dalam pemrosesan data besar yang mengandung banyak dokumen atau teks. Dengan menggunakan TF-IDF *Vectorizer*, dapat dilakukan suatu pengelompokan teks, analisis sentimen, deteksi topik, dan sistem rekomendasi secara efisien. Selain itu, dengan menerapkan TF-IDF *Vectorizer* dalam penelitian ini mengenai emisi gas rumah kaca (GRK) di Indonesia, data teks abstrak akan dikonversi menjadi representasi vektor TF-IDF dengan memperhitungkan bobot kata dalam setiap abstrak. Kata-kata kunci yang relevan dalam data teks yang berkaitan dengan GRK dan menganalisis tingkat signifikansinya dalam dokumen tersebut.

2.2.3 *Principal Component Analysis (PCA)*

PCA (*Principal Component Analysis*) adalah sebuah metode statistik yang digunakan untuk mengurangi dimensi dari data dengan mempertahankan informasi yang paling penting. PCA digunakan untuk mengidentifikasi pola dalam data, dan mengubah data tersebut menjadi bentuk yang lebih mudah dipahami dan diinterpretasikan [22].

PCA lebih sering digunakan dalam analisis data numerik, namun PCA juga dapat digunakan dalam analisis teks. Dalam analisis teks, PCA digunakan untuk mengurangi dimensi dari data teks dengan mempertahankan informasi yang paling penting [23]. PCA dapat digunakan untuk mengidentifikasi pola dalam data teks, seperti pola kata yang sering muncul bersama-sama atau pola topik yang terkait. Dengan mengurangi dimensi dari data teks, PCA dapat mempermudah analisis dan interpretasi data teks.

2.2.4 **K-Means**

K-Means adalah salah satu algoritma *clustering* yang paling populer digunakan dalam analisis data. Algoritma ini digunakan untuk membagi data menjadi beberapa kelompok atau kluster berdasarkan kesamaan fitur atau karakteristik tertentu [24]. Dalam K-Means, data dikelompokkan berdasarkan jarak antara titik data dan *centroid*, yang merupakan pusat dari setiap kluster. Tujuan dari algoritma ini adalah untuk meminimalkan varians dalam setiap kluster dan memaksimalkan jarak antara kluster.

Salah satu penerapan K-Means yang sering digunakan adalah dalam analisis teks. Dalam analisis teks, K-Means digunakan untuk mengelompokkan dokumen berdasarkan kesamaan topik atau kata kunci tertentu [25]. Misalnya, jika kita memiliki kumpulan dokumen tentang olahraga, K-Means dapat digunakan untuk mengelompokkan dokumen berdasarkan jenis olahraga atau topik tertentu seperti sepak bola, basket, atau tenis. Dalam hal ini, K-Means dapat membantu mengidentifikasi pola dan tren dalam data teks yang besar dan kompleks.

K-Means dapat digunakan untuk memproses hasil dari analisis PCA (*Principal Component Analysis*). PCA adalah teknik statistik yang digunakan untuk mengurangi dimensi data dengan mempertahankan informasi yang paling penting. Setelah PCA diterapkan pada data, K-Means dapat digunakan untuk mengelompokkan data menjadi beberapa kluster berdasarkan kesamaan fitur atau karakteristik tertentu [25]. Hasil dari analisis K-Means dapat membantu mengidentifikasi pola dan tren dalam data yang besar dan kompleks, serta memudahkan interpretasi hasil analisis.

2.3 **Teknik Kitchenham**

Teknik Kitchenham merujuk pada metode sistematis literature review (SLR) yang dikembangkan oleh Barbara Kitchenham, seorang peneliti terkemuka di bidang rekayasa perangkat lunak [26]. Pedoman Kitchenham untuk melakukan SLR telah menjadi sangat diakui dan sering disebut sebagai "teknik Kitchenham".

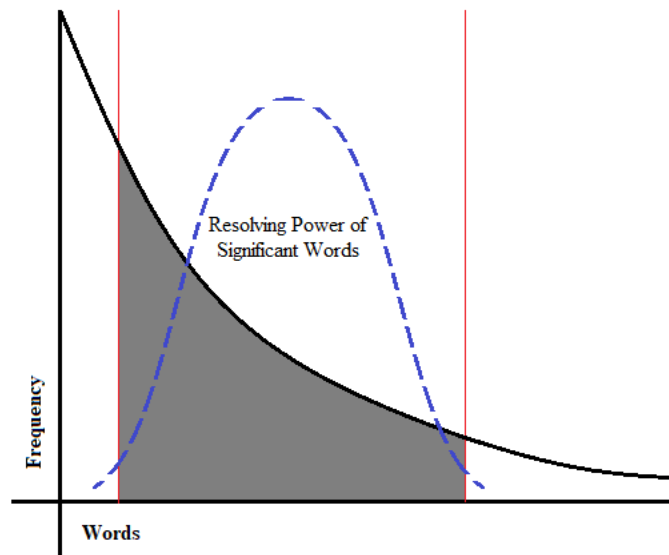
Teknik Kitchenham melibatkan proses dalam mengidentifikasi, mengevaluasi, dan mensintesis semua studi yang relevan pada pertanyaan atau topik penelitian tertentu. Proses ini meliputi pencarian yang komprehensif dan sistematis dari beberapa database dan sumber lainnya, pengembangan kriteria inklusi dan eksklusi, dan teknik analisis data untuk memastikan bahwa hasilnya dapat diandalkan dan tidak bias. Pedoman Kitchenham telah banyak diadopsi dalam penelitian rekayasa perangkat lunak dan telah berkontribusi pada pengembangan rekayasa perangkat lunak berbasis bukti (EBSE) [26, 27]. Teknik Kitchenham telah menjadi metode standar untuk melakukan SLR dalam rekayasa perangkat lunak dan telah membantu meningkatkan kualitas dan keandalan penelitian di bidang tersebut.

Dalam konteks penggunaan teknik Kitchenham, terdapat penggunaan operator logika seperti *AND*, *OR*, dan *NOT* dalam pencarian studi yang relevan. Operator logika tersebut digunakan untuk menggabungkan kata kunci pencarian dan memperluas atau mempersempit cakupan pencarian. Misalnya, jika ingin mencari studi yang relevan dengan topik "*greenhouse gas*" dapat menggunakan operator logika *AND* untuk menggabungkan kata kunci "*greenhouse*" dan "*gas*". Dengan demikian, pencarian akan memperlihatkan studi yang hanya relevan dengan kedua kata kunci tersebut.

Selain itu juga, penggunaan asterisk (*) dalam pencarian studi yang relevan juga dapat mencakup variasi kata yang terkait dengan akar kata tertentu. Misalnya, jika menggunakan kata kunci "*communicat**", pencarian akan mencakup studi yang menggunakan kata "*communicate*", "*communication*", "*communicating*", dan sebagainya. Dengan menggunakan asterisk, akan dapat memperluas cakupan pencarian dan memastikan bahwa semua studi yang relevan ditemukan, terutama jika ada variasi dalam cara menyebutkan kata kunci tertentu.

2.4 Luhn's Theory

Teori yang digagas oleh H. P. Luhn, dalam penelitiannya yaitu penciptaan otomatis abstrak literatur menggunakan mesin pemrosesan data [8]. Implementasi Luhn melibatkan analisis teks lengkap dari sebuah artikel dan pemilihan kalimat dan frasa kunci untuk membuat ringkasan singkat dari konten artikel. Mesin pemrosesan data IBM 704 digunakan untuk menghitung ukuran relatif signifikansi untuk kata dan kalimat individu, yang kemudian digunakan untuk memilih informasi paling penting untuk disertakan dalam abstrak. Implementasi ini terbukti efektif dalam menciptakan abstrak yang melayani tujuan abstrak konvensional dan dapat menghemat waktu dan usaha bagi pembaca literatur teknis.



Gambar 2.1. Diagram *Word Frequency*

Hal yang menjadi sebagai kata kunci dalam penerapan metode otomatis untuk membuat abstrak literatur menggunakan mesin pemrosesan data adalah kata-kata yang paling relevan dan penting dalam artikel yang sedang dianalisis. Kata-kata ini biasanya dipilih berdasarkan topik atau subjek dari artikel tersebut. Dalam *common words* seperti *pronouns*, *prepositions*, dan *articles* dihapus dari daftar kata-kata yang dianalisis karena kata-kata ini tidak memberikan informasi yang signifikan tentang topik atau subjek dari artikel. Sebaliknya, kata-kata yang lebih spesifik dan relevan seperti kata benda dan kata kerja digunakan sebagai kata kunci untuk memilih kalimat-kalimat kunci dalam artikel yang akan digunakan untuk membuat abstrak.

Terdapat penelitian yang dilakukan oleh Rahmah et al [10] dalam menerapkan teori Luhn, dengan Hukum Zipf yang menyatakan bahwa frekuensi kemunculan sebuah kata dalam teks berbanding terbalik dengan peringkatnya dalam daftar frekuensi kata. Dalam konteks pembuatan abstrak otomatis, hukum Zipf dapat digunakan untuk memilih kata-kata kunci yang paling relevan dan penting dalam artikel yang sedang dianalisis. Selain itu, juga membahas tentang skema Bradford [8], yang merupakan metode untuk mengelompokkan jurnal-jurnal ilmiah berdasarkan topik atau subjek yang sama, dengan adanya *lower cut* dan *upper cut* akan dipilih kata-kata yang diseleksi sebagai tema atau konsep yang dicari sebagaimana yang ditunjukkan pada Gambar 2.1 [8]. Metode tersebut digunakan untuk membantu pembaca literatur teknis menemukan jurnal-jurnal yang relevan dengan topik atau subjek yang sedang mereka teliti [10].



Universitas
Pertamina

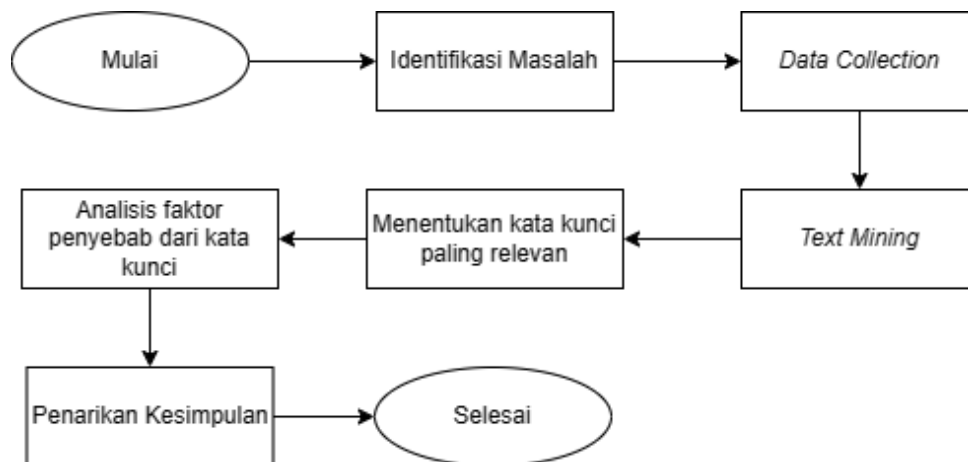
BAB III

METODE PENELITIAN

3.1 Tahapan Penelitian

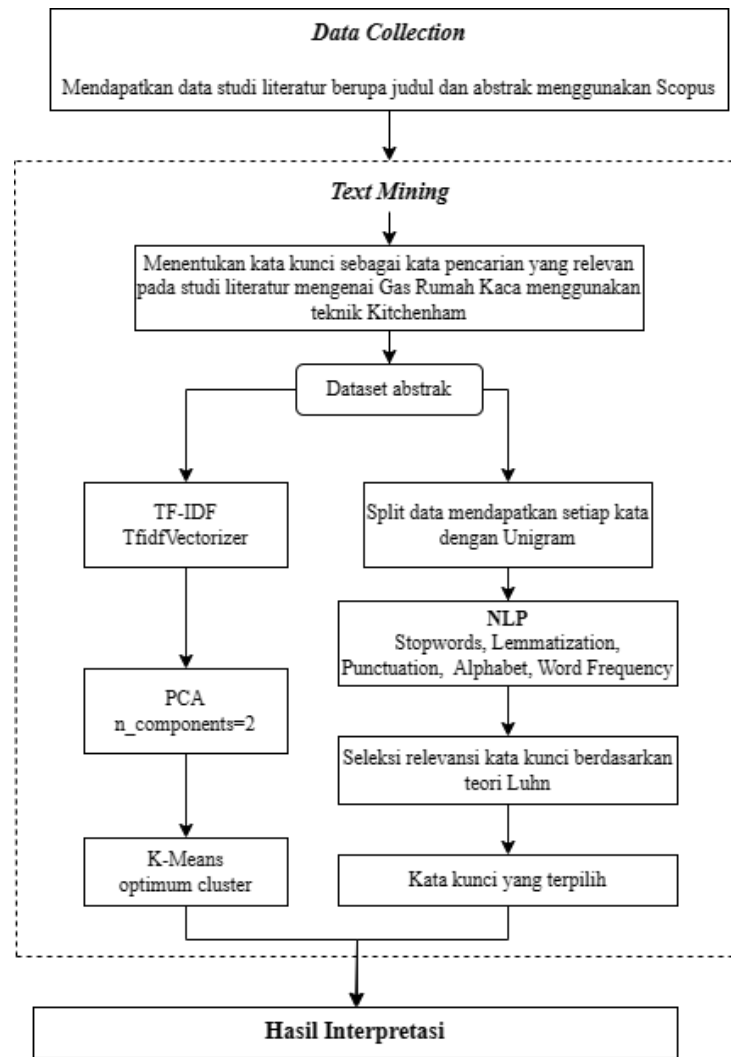
Penelitian ini mengikuti diagram alir yang terdiri dari beberapa tahapan. Tahap pertama adalah mengidentifikasi masalah terkait permasalahan Gas Rumah Kaca (GRK) di Indonesia. Selanjutnya, dilakukan tahap *data collection* dengan mengumpulkan data abstrak dari *database* jurnal penelitian dengan platform Scopus. Setelah itu, data tersebut akan diproses menggunakan teknik *text mining* untuk mendapatkan kata-kata kunci yang relevan.

Selanjutnya, berdasarkan kata-kata dasar yang dihasilkan dari *text mining*, dilakukan pemilihan kata kunci dari kata-kata dasar yang didapatkan dengan mencari relevansi dengan faktor-faktor penyebab GRK, serta hasil analisis menggunakan metode yang dipakai yaitu TF-IDF, PCA, K-Means, dan teori Luhn. Tahap akhir adalah menarik kesimpulan dari hasil penelitian. Dengan demikian, diagram alir penelitian ini menggambarkan langkah-langkah yang sistematis untuk memahami faktor penyebab permasalahan GRK di Indonesia berdasarkan analisis kata kunci.



Gambar 3.1. Diagram Alir Penelitian

Pada tahapan ketika melakukan *data collection* hingga ke hasil interpretasi, dengan menerapkan skema yang pernah digunakan oleh Lai [28] dengan skemanya yaitu "*social media to-concepts*", secara rinci akan dijelaskan di alur tahapannya sebagai berikut:



Gambar 3.2. Tahapan *Data Collection* hingga hasil interpretasi

Berdasarkan tahapan yang digambarkan pada Gambar 3.2, hal yang pertama dilakukan yaitu identifikasi masalah terkait emisi gas rumah kaca (GRK) di Indonesia yang relevan untuk pencarian data yang dilanjutkan dengan studi literatur pada *dataset* yang mendukung untuk mengidentifikasi faktor-faktor penyebab GRK. Adapun pada pemilihan data akan diambil berdasarkan abstrak yang berhubungan dengan GRK di Indonesia, kemudian dilakukan analisis lebih detail terkait relevansi, topik spesifik, dan keterbaruan untuk menentukan teknik apa yang akan digunakan dalam pencarian kata-kata kunci dari *dataset* tersebut.

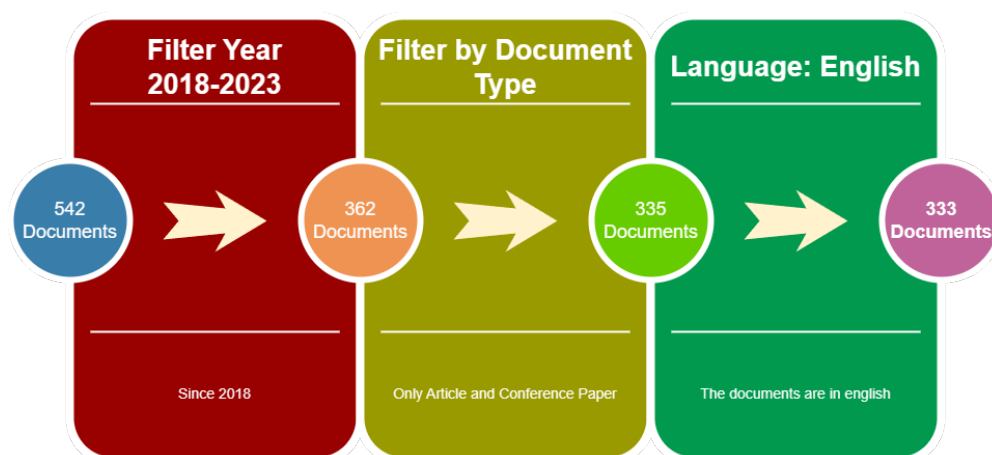
Selanjutnya, data yang telah dikumpulkan akan diproses menggunakan NLP dengan mengimplementasikan *splitting data* secara Unigram (N-grams = 1) untuk mendapatkan kata kunci per huruf dari *dataset*. Kata-kata kunci ini nantinya akan digunakan dalam analisis kualitatif. Dari banyaknya kata kunci tersebut yang memiliki frekuensi paling besar, akan dipilih kata-kata yang paling spesifik berdasarkan teori Luhn. Setelah semua tahapan di atas dilakukan, proses interpretasi akan dilakukan untuk mencari faktor-faktor penyebab GRK berdasarkan hasil yang telah diperoleh. Selanjutnya, akan ditemukan hasil upaya atau solusi dari faktor-faktor penyebab GRK yang telah dipilih berdasarkan kata kunci yang telah didapatkan.

3.2 Data Collection

Pada tahapan *data collection*, data yang diambil berdasarkan studi literatur yang telah didapatkan melalui *search engine* Scopus yang mengulas terkait GRK. Data yang diambil tersebut merupakan kumpulan abstrak di artikel atau penelitian yang dipilih berdasarkan pemilihan kata kunci dengan teknik Kitchenham [26]. Adapun dalam mendapatkan data studi literatur, diawali dengan mencari faktor-faktor penyebab gas rumah kaca terlebih dahulu dengan mendapatkan sumber referensi yang telah pernah membahas faktor penyebab gas rumah kaca. Adapun cara pencariannya terinspirasi dari penelitian yang dilakukan oleh [29] dengan melakukan teknik Kitchenham [26].

Pada pencarian studi literatur berupa jurnal atau artikel yang akan dipakai sebagai data penelitian ini, dilakukan pencarian di platform Scopus dengan implementasi Kitchenham yakni mendefinisikan objektif penelitian dengan kata kunci agar didapatkan sekumpulan dan abstrak yang relevan secara efektif. Dimulai dari aspek yang berkaitan dengan penyebab, maka kata kunci yang dipakai yaitu *"factor, feature, variable, cause, character, impact"*, kemudian dicari berdasarkan aspek hasil yang ditimbulkan yaitu dengan kata kunci *"affect, contribute, product, generate, conduct"*. Setelah itu, dari segi nama subjek yang dipakai mengenai GRK yakni *"GHG, greenhouse gas"*, dan yang berdasarkan negara atau daerah spesifik yaitu di Indonesia. Maka, rangkaian kata kunci yang akan dipakai untuk ditulis di *search engine* Scopus ialah:

(factor OR featur* OR variabl* OR caus* OR character* OR impact) AND (affect* OR contribut* OR produc* OR generat* OR conduc*) AND (ghg OR greenhouse AND gas*) AND (indonesia OR indonesian)*



Gambar 3.3. Tahapan proses pemilihan data studi literatur

Setelah mendefinisikan kata kunci pada pencarian berdasarkan objektif penelitian, terdapat tahapan yang diperlukan dalam pemilihan studi literatur yang cocok sebagaimana yang telah digambarkan pada Gambar 3.3. Pada pencarian pertama, didapatkan sebanyak 542 dokumen yang tersedia. Agar data studi literatur yang dipakai merupakan dokumen yang baru, maka diambil berdasarkan terbitan publikasi 5 tahun saat ini yaitu tahun 2018 hingga saat ini, maka dihasilkan sebanyak 362 dokumen. Setelah itu, untuk mendapatkan jenis studi literatur yang spesifik, maka dipilih jenisnya yaitu *article* dan *conference paper* dan didapatlah sebanyak 335 dokumen. Pada tahap akhir, dari dokumen studi

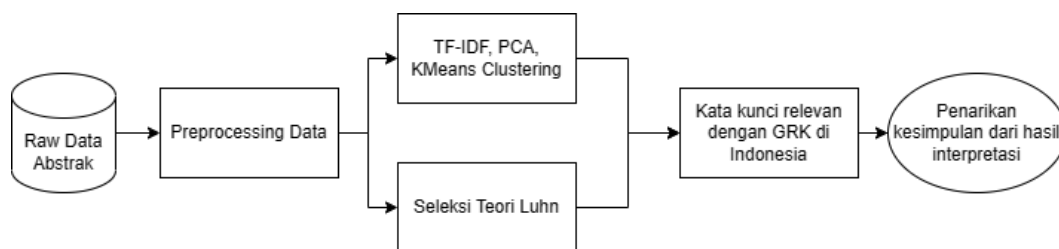
literatur yang didapat perlu dipilih hanya yang berbahasa Inggris saja, untuk pemahaman dan saat pemrosesan data nanti.

Berikut ini *query search* yang dihasilkan setelah semua pemilihan dokumen studi literatur tersebut dilakukan di platform Scopus:

TITLE-ABS-KEY (factor OR featur* OR variabl* OR caus* OR character* OR impact AND affect* OR contribut* OR produc* OR generat* OR conduc* AND ghg OR greenhouse AND gas* AND indonesia OR indonesian) AND PUBYEAR > 2017 AND PUBYEAR < 2024 AND (LIMIT-TO (LANGUAGE , "English")) AND (LIMIT-TO (DOCTYPE , "ar") OR LIMIT-TO (DOCTYPE , "cp"))*

Maka, setelah pendefinisian kata kunci dengan teknik Kitchenham, dilanjutkan dengan penyaringan dokumen di platform Scopus, sebanyak 333 dokumen, akan dipakai teks abstrak sebagai *dataset* dalam penelitian ini.

3.3 Tahapan *Text Mining* dan Interpretasinya



Gambar 3.4. Alur tahapan saat dilakukan *text mining* hingga hasil interpretasi

Pada tahapan ini, akan lebih banyak implementasi menggunakan NLP, menggunakan data yang berbentuk teks sekumpulan kalimat abstrak yang telah didapatkan, nantinya data tersebut dilakukan *preprocessing data* untuk menyortir teks agar dapat diproses dan siap untuk digunakan dalam analisis lebih lanjut secara baik. Setelah dilakukannya *preprocessing data*, akan dilakukan analisis 2 arah yaitu dengan teknik *clustering* dan teknik mencari kata kunci yang relevan.

Pada analisis *clustering*, proses tersebut melakukan representasi teks menggunakan TF-IDF, dilanjutkan dengan reduksi dimensi menggunakan PCA, kemudian melakukan K-Means clustering untuk mengelompokkan abstrak berdasarkan fitur-fitur yang dihasilkan dari PCA. Hasil akhirnya adalah pengelompokan abstrak ke dalam beberapa klaster berdasarkan kesamaan fitur-fitur TF-IDF yang telah direduksi menggunakan PCA.

Kemudian pada analisis seleksi kata-kata dasar yang menjadi kata kunci yang memiliki relevansi dengan faktor penyebab GRK, akan diseleksi kata kunci apa yang sesuai berdasarkan cara dari teori Luhn. Nantinya dari data berupa kata kunci yang terpilih atau terseleksi tersebut, serta hasil dari *clustering* sebelumnya, akan dipakai di dalam mencari hasil interpretasi berupa kata kunci apa yang sering dibahas sebagai penyebab GRK di Indonesia, dan diakhiri dengan penarikan kesimpulan.



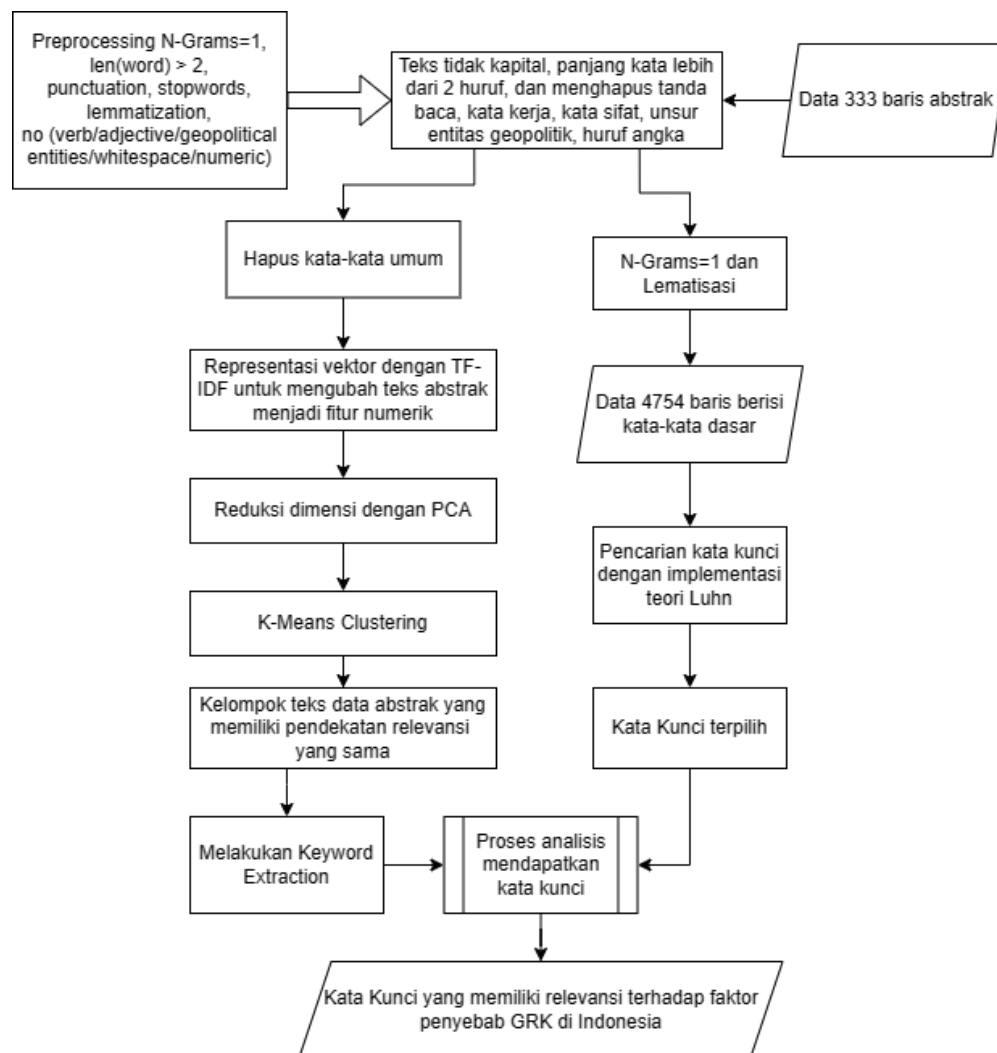
BAB IV

HASIL DAN PEMBAHASAN

4.1 *Preprocessing Data*

Data yang telah dikumpulkan berupa 333 dokumen dengan abstrak, pada penelitian ini akan dipakai sehingga berisi 333 baris. Setelah itu, dilakukan suatu *preprocessing data* pada teks tersebut. Pada setelah langkah *preprocessing data* teks tersebut, akan dibagi menjadi dua alur sebagaimana yang telah digambarkan pada Gambar 3.4, yang pertama ialah sebelum dilakukan suatu pengelompokan abstrak serupa menjadi klaster dan yang alur lainnya sebelum dilakukan pemilihan kata dasar untuk mencari kata kunci yang relevan dengan teori Luhn.

Adapun tahapan dalam *preprocessing data* hingga mendapatkan hasil berupa kata kunci yang dicari digambarkan pada gambar di bawah ini:



Gambar 4.1. Proses mendapatkan kata-kata dasar dari data studi literatur

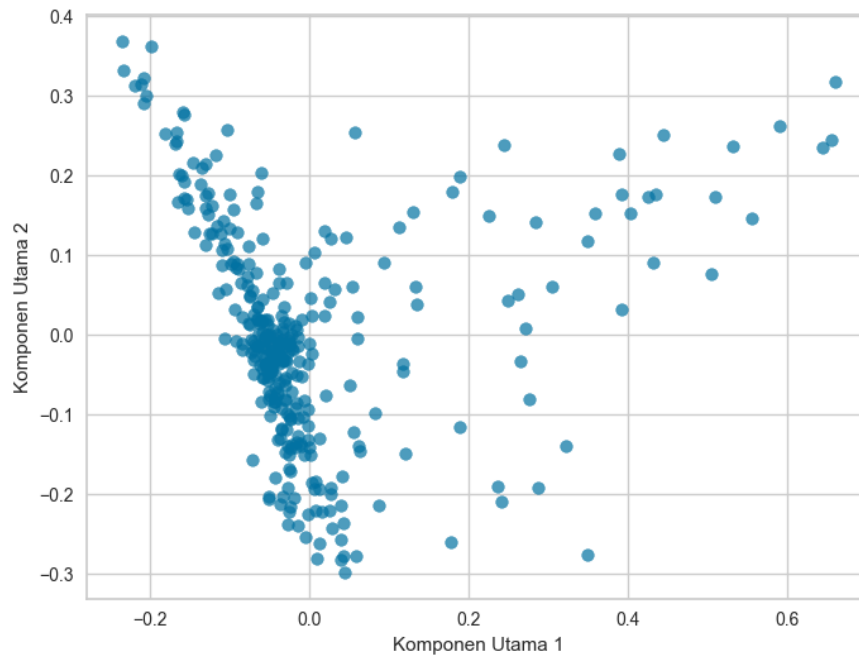
Di bawah ini, merupakan teknik berupa istilah yang digunakan pada saat melakukan *preprocessing data* sebagai berikut:

1. *Lowercase*: Membuat semua teks pada data tersebut menjadi bukan huruf kapital.
2. *Punctuation*: Menghapus semua unsur tanda baca.
3. *Stopword*: Menghapus unsur kata-kata yang tidak bermakna, seperti '*is, a, it, in, he, is*'.
4. *N-Grams=1 (Unigram)*: Mendapatkan data yang berisi dari kumpulan kata-kata dasar.
5. *len(word) > 2*: Mendapatkan kata-kata yang memiliki jumlah huruf lebih dari dua.
6. *Lemmatization*: Memilih satu dari setiap kata dasar yang memiliki bentuk lema yang sama, seperti terdapat kata '*organizing*' dan '*organize*', maka akan dipilih '*organize*'.
7. *No Verb and No Adjective*: Menghapus unsur teks yang mengandung kata kerja dan kata sifat.
8. *No Geopolitical Entities (GPE)*: Menghapus semua unsur teks yang mengandung unsur nama-nama negara, kota, atau wilayah.
9. *No Numeric*: Menghapus kata-kata yang memiliki unsur angka saja.
10. *No Whitespace*: Menghapus unsur yang hanya berisi kata spasi.

4.2 Hasil Pengelompokkan Kata Kunci berdasarkan teknik *clustering*

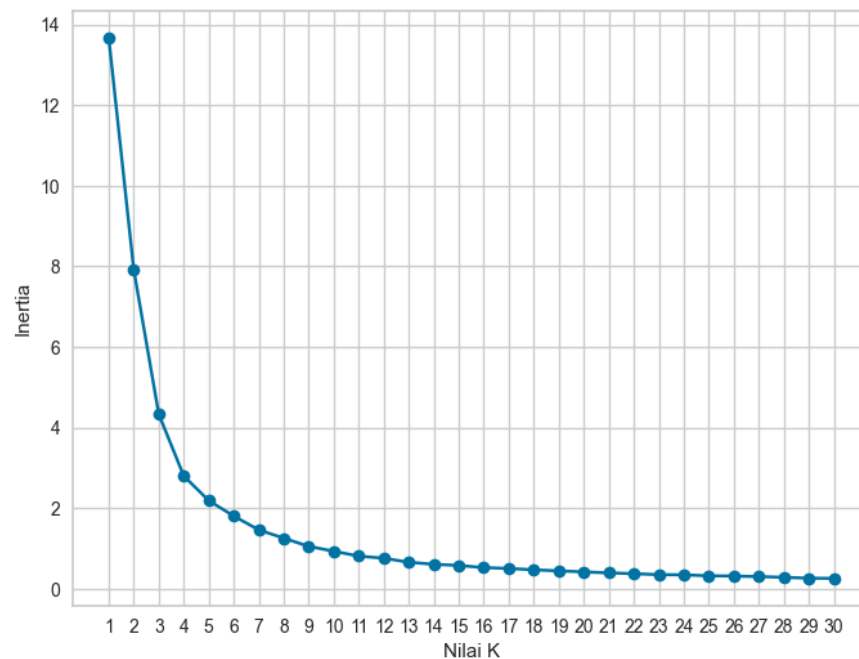
Pada tahapan ini, hasil dari *preprocessing text* sebelumnya, dilakukan suatu proses penghapusan kata-kata umum dengan *stopword* seperti yang dijelaskan pada poin sebelumnya. Kemudian, data yang telah didapatkan tersebut diproses menjadi representasi numerik dari teks abstrak menggunakan metode TF-IDF. TF-IDF digunakan untuk menghitung bobot kata-kata dalam setiap abstrak berdasarkan frekuensi kemunculan kata tersebut di abstrak tertentu dan juga frekuensi kemunculan kata tersebut di semua abstrak. Dengan menggunakan *TfidfVectorizer*, teks abstrak dikonversi menjadi representasi vektor TF-IDF dengan memperhitungkan bobot kata dalam setiap abstrak.

Setelah representasi TF-IDF dibuat, langkah selanjutnya adalah melakukan reduksi dimensi menggunakan PCA. PCA digunakan untuk mengurangi dimensi vektor TF-IDF menjadi dimensi yang lebih rendah, sehingga memungkinkan untuk memvisualisasikan data dalam ruang dua atau tiga dimensi. Dalam implementasinya, digunakan *PCA(n_components=2)* untuk mengambil dua komponen utama dan memvisualisasikan datanya dalam ruang dua dimensi. Komponen utama 1 dan komponen utama 2 pada klusterisasi dengan K-Means merujuk pada hasil reduksi dimensi menggunakan teknik Principal Component Analysis (PCA). Jika nilai komponen utama 1 (nilai x) positif, itu berarti abstrak tersebut cenderung memiliki kontribusi positif dari fitur-fitur yang berkaitan dengan komponen utama 1 tersebut. Sebaliknya, jika nilai komponen utama 1 (nilai x) negatif, itu berarti abstrak tersebut cenderung memiliki kontribusi negatif dari fitur-fitur yang berkaitan dengan komponen utama 1 tersebut.



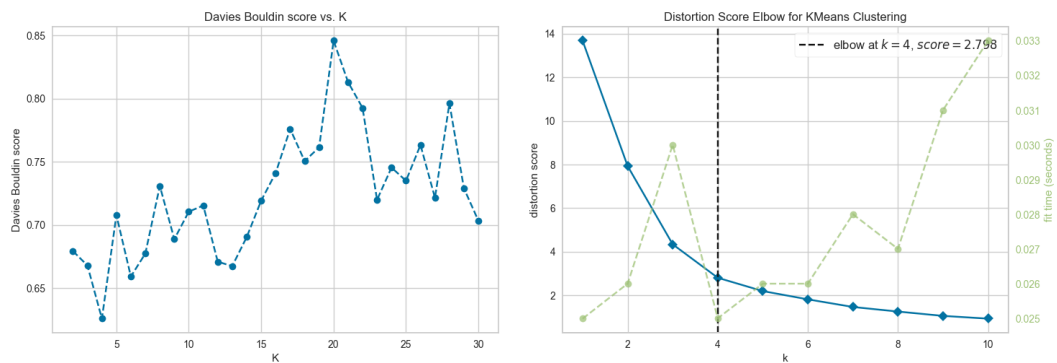
Gambar 4.2. Visualisasi hasil dari PCA

Setelah visualisasi PCA, langkah selanjutnya adalah melakukan K-Means *clustering* untuk mengelompokkan abstrak menjadi beberapa kluster berdasarkan fitur TF-IDF yang telah direduksi menggunakan PCA. Pada implementasi KMeans, dibutuhkan satu parameter yaitu nilai dari K. Untuk menentukan nilai K atau jumlah *num_clusters* yang optimal, dilakukan pencarian dengan menampilkan *Elbow Plot* menggunakan K-Means, dengan uji coba rentang kluster 1-30.



Gambar 4.3. *Elbow Plot* untuk menentukan nilai kluster dengan K-Means

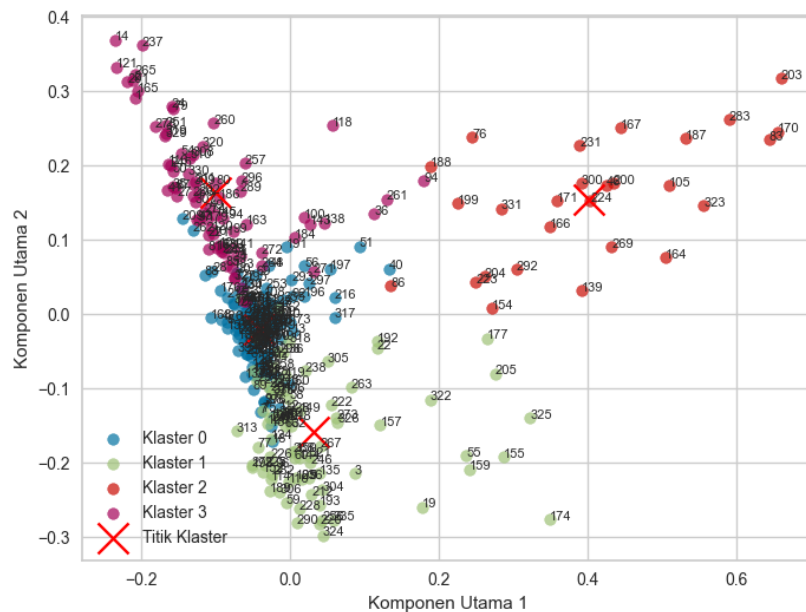
Dalam mempermudah penentuan nilai K yang optimum, dilakukan beberapa metrik evaluasi dengan *Davies-Bouldin index* dan *library ElbowVisualizer*, kedua metrik evaluasi tersebut digunakan untuk mengevaluasi nilai K yang optimal dalam algoritma K-Means. *Davies-Bouldin index* mengukur kualitas pembagian kluster berdasarkan jarak antara kluster yang berdekatan dan dispersi internal setiap kluster [30]. Semakin rendah nilai *Davies-Bouldin index*, semakin baik pembagian kluster dan semakin optimal nilai K yang digunakan. Selain itu, dengan *library ElbowVisualizer* untuk membantu dan memperjelas dari visualisasi *Elbow Plot*, pada saat titik di nilai *inertia* mulai menurun secara lambat atau berhenti menurun drastis, itulah titik yang disarankan sebagai jumlah kluster yang optimal. Hasil dari kedua evaluasi tersebut menunjukkan nilai K=4 yang merupakan nilai K yang optimal, sehingga diterapkan fungsinya yakni *KMeans(n_clusters=4)*.



Gambar 4.4. Metrik Evaluasi nilai K dengan *Davies-Bouldin index* dan *library ElbowVisualizer*

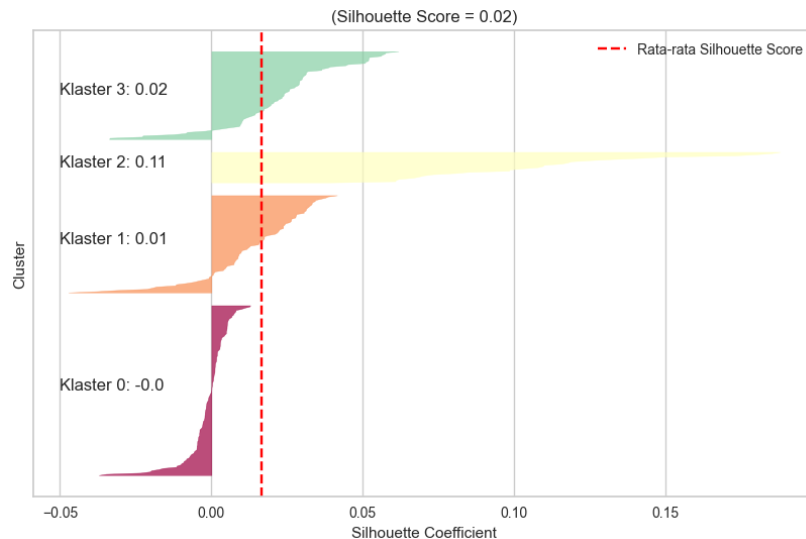
Setelah melakukan K-Means *clustering* dan menentukan jumlah kluster yang optimal berdasarkan *elbow plot* dengan didukung validasi berupa evaluasi metrik, selanjutnya dapat diketahui dalam setiap kluster dengan menggunakan atribut dari objek K-Means untuk menyimpan label kluster pada setiap sampel data abstrak dan mengelompokkannya berdasarkan warna yang berbeda.

Dalam konteks klusterisasi perwarnaan dengan K-Means, komponen utama 1 dan komponen utama 2 adalah dua dimensi baru yang dihasilkan dari PCA. Setiap baris data dalam dataset asli direpresentasikan oleh dua nilai pada komponen utama 1 dan komponen utama 2. Nilai-nilai ini menggambarkan lokasi relatif data dalam ruang dua dimensi yang telah direduksi. Pada saat PCA telah dilakukan, visualisasi data abstrak dalam ruang dua dimensi dengan *scatter plot* menunjukkan bahwa setiap titik pada *scatter plot* merepresentasikan satu abstrak, dan dapat dilihat berupa pola dan struktur data sebelum melakukan *clustering*.



Gambar 4.5. Klasterisasi Pewarnaan dengan K-Means

Selain itu, untuk mempermudah dalam melihat kluster mana yang lebih baik, dilakukanlah *Silhouette plot* yang merupakan metode untuk mengukur seberapa baik setiap data abstrak dikelompokkan oleh K-Means. Plot ini menunjukkan nilai *silhouette score* untuk setiap titik kluster pada data abstrak, nilai *silhouette score* mengukur seberapa mirip data abstrak dengan kelompoknya sendiri dibandingkan dengan kelompok lain. *Silhouette plot* juga membantu mengidentifikasi apakah ada kelompok yang tidak terdefinisi dengan baik atau apakah data abstrak tergolong dalam kelompok yang sesuai.



Gambar 4.6. *Silhouette Plot* untuk K-Means Clustering

Dari hasil tersebut, terdapat kluster yang memiliki warna yang menjulang jauh melebihi rata-rata dari *Silhouette Score* yang telah didapatkan senilai 0.02 tersebut, yang menandakan bahwa kluster tersebut memiliki *Silhouette Coefficient* yang tinggi dibandingkan dengan kluster lainnya. Kluster

dengan *Silhouette Coefficient* yang tinggi menunjukkan bahwa titik-titik data dalam klaster tersebut lebih dekat satu sama lain dan lebih terpisah dari titik-titik data di klaster lain. Hal ini, menunjukkan klaster tersebut memiliki titik-titik data yang sangat serupa dan kohesif. Dalam konteks analisis *clustering*, hal ini dapat diartikan sebagai kelompok yang lebih homogen atau memiliki karakteristik yang sama.

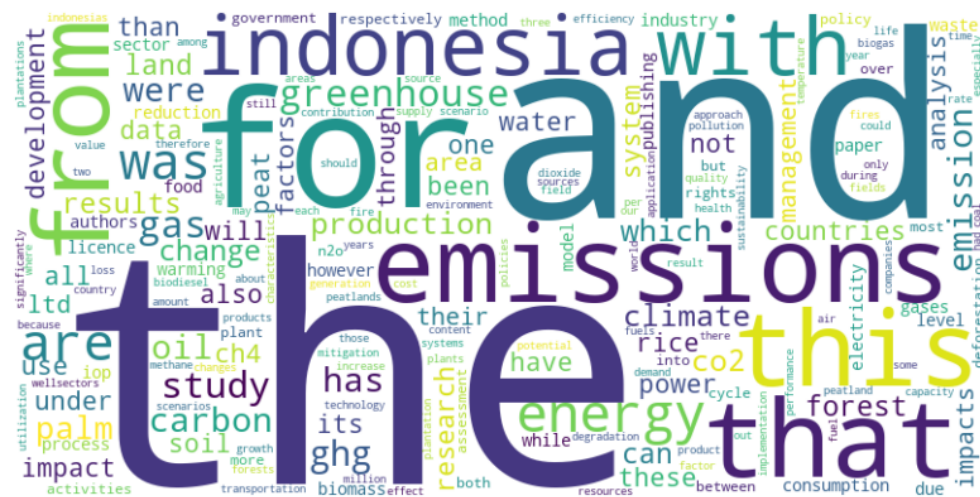
Kemudian, pada tahap akhir, dilakukan suatu *keyword extraction* untuk menentukan kata-kata kunci yang paling relevan dalam setiap klaster. Untuk setiap klaster, dilakukan *keyword extraction* masih dengan metode TF-IDF. TF-IDF digunakan untuk menentukan kata-kata kunci yang paling relevan dan memiliki bobot tertinggi dalam klaster tersebut. Setelah mendapatkan kata-kata kunci untuk setiap klaster, akan dihasilkan berupa kata-kata kunci dari proses *keyword extraction*. Berikut ini, 20 kata-kata kunci pada setiap klasternya:

- a. **Cluster 0:** *emissions, indonesia, greenhouse, carbon, ghg, study, emission, gas, climate, change, co2, results, development, research, countries, data, production, impact, food, impacts*
- b. **Cluster 1:** *emissions, forest, soil, peat, rice, ch4, indonesia, greenhouse, land, gas, study, carbon, water, ghg, n2o, emission, area, peatlands, management, palm*
- c. **Cluster 2:** *oil, palm, indonesia, production, biodiesel, emissions, industry, gas, land, greenhouse, study, products, impact, supply, ghg, change, impacts, ispo, cpo, rspo*
- d. **Cluster 3:** *energy, power, indonesia, gas, emissions, electricity, production, greenhouse, co2, study, generation, emission, plant, consumption, ghg, coal, system, use, biogas, efficiency*

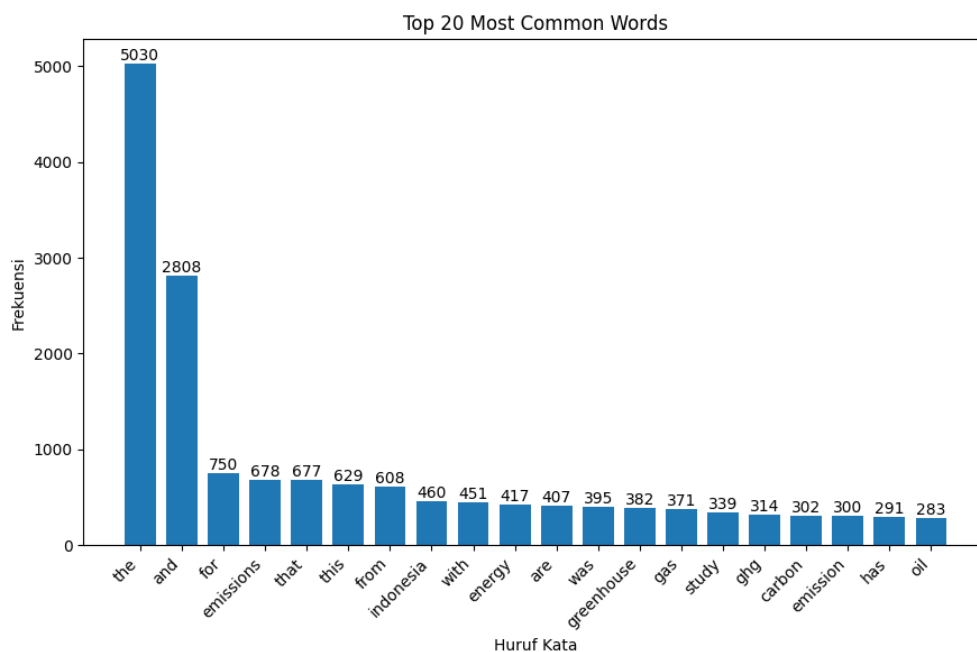
4.3 Hasil Pemilihan Kata Kunci berdasarkan Teori Luhn

Pada tahapan memilih kata kunci, berdasarkan teori Luhn, data yang telah didapatkan dari hasil *preprocessing data* sebelumnya di bagian 4.1, dilakukan suatu pemecahan kata-kata dasar dengan N-Grams = 1 (Unigram) dan melakukan teknik lematisasi yang telah dijelaskan juga di bagian 4.1.

Setelah dilakukan proses dengan teknik tersebut, didapatkan sebanyak 4754 baris data berupa kata-kata dasar berdasarkan teknik *preprocessing* yang diterapkan sebelumnya. Kemudian untuk melihat suatu frekuensi dari masing-masing kata dasar, terdapat visualisasi dengan *Word cloud* untuk menunjukkan representasi visual dari data teks yang menggambarkan frekuensi kata-kata dalam teks tersebut [31]. Dalam *word cloud*, kata-kata yang muncul lebih sering akan ditampilkan dengan ukuran yang lebih besar, sedangkan kata-kata yang muncul lebih jarang akan ditampilkan dengan ukuran yang lebih kecil.

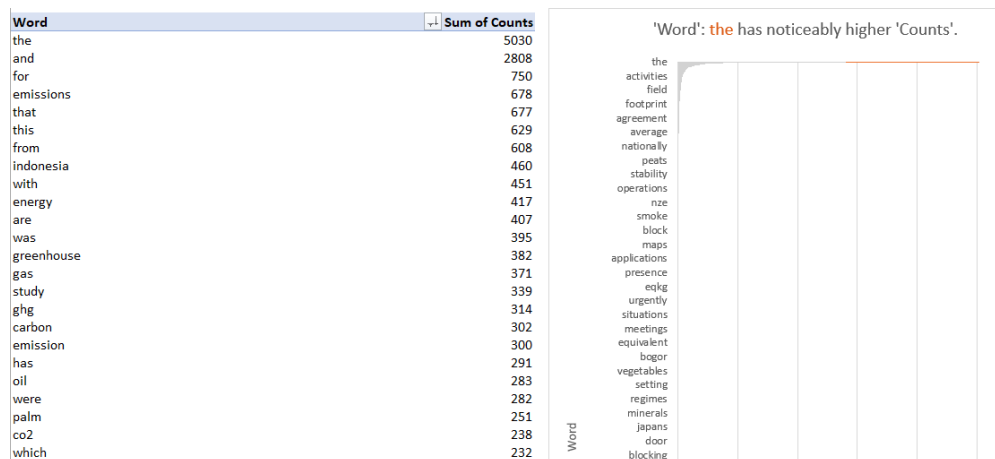


Terdapat juga visualisasi untuk melihat 20 kata yang memiliki frekuensi terbesar pada data tersebut.

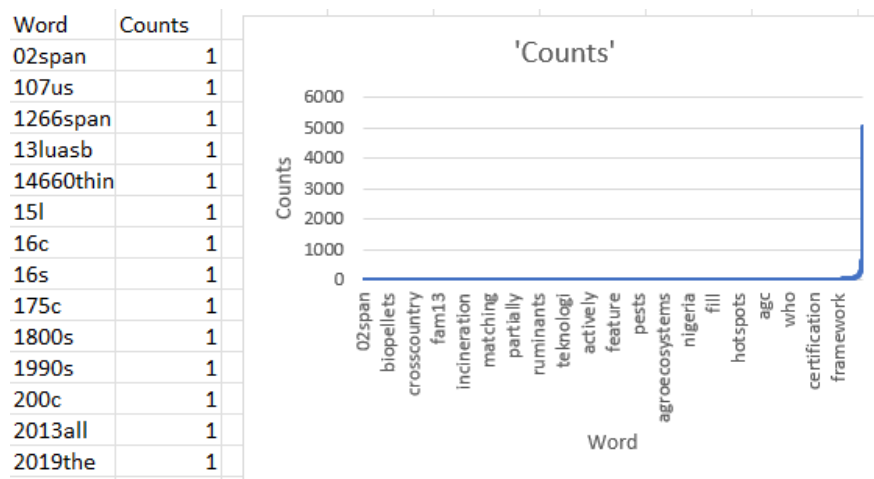


Namun, yang diproses seleksi teori Luhn bukanlah kata-kata yang memiliki frekuensi terbesar, semua data sebanyak 4754 kata sebelumnya, akan dipakai pada tahapan teori Luhn. Adapun bagian-bagian kata yang memiliki frekuensi terbesar dan frekuensi terkecil akan dibuang, karena berdasarkan teori Luhn porsi pada rentang tersebut merupakan istilah-istilah yang kurang baik dalam memilih kata kunci yang relevan secara optimal.

Pada proses seleksi tersebut, dilakukan analisis secara manual dengan melihat kata-kata dasar berdasarkan relevansi yang ingin dipilih.



Gambar 4.9. *Preview* kata-kata dengan frekuensi terbesar



Gambar 4.10. *Preview* kata-kata dengan frekuensi terendah

Agar dalam seleksi kata kunci tersebut, setelah dilakukan suatu analisis secara mendalam dengan melihat kata-kata dasar berdasarkan relevansi, bobot frekuensi, dan juga bagian interval yang cukup baik sebagai pemilihan kata kunci, agar lebih presisi, maka porsi besaran pada kata-kata di rentang frekuensi terbesar dan frekuensi terkecil akan dibuang sebesar 10% dari keseluruhan jumlah data yang berarti dihapus 475 kata di bagian rentang frekuensi terbesar dan rentang frekuensi terendah, sehingga tersisa 3804 baris data.

Setelah itu, membuang kata-kata dasar yang memiliki frekuensi lebih dari atau sama dengan 10 sehingga tersisa 169 baris data. Kemudian baru dilakukan rentang mana yang akan dipilih sebagai *upper cut* dan *lower cut* dari data tersebut. Setelah dilakukan analisis kembali, dipilih rentang sebesar 15% untuk dijadikan sebagai kata kunci yang terpilih. Dari hasil kata kunci dipilihlah sebanyak 25 kata yang memiliki relevansi berdasarkan hasil analisis secara manual dan penerapan teori Luhn.

137	central	14
138	challenges	14
139	context	14 LOWER CUT
140	contributions	14
141	control	14
142	criteria	14
143	crops	14
144	damage	14
145	difference	14
146	diversity	14
147	drainage	14
148	future	14
149	generally	14
150	governments	14
151	green	14
152	half	14
153	household	14
154	ice	14
155	petroleum	14
156	temperatures	14
157	priority	14
158	recycling	14
159	risk	14
160	root	14
161	simulation	14
162	soils	14
163	speed	14
164	states	14
165	pome	14 UPPER CUT
166	tool	14
167	usage	14

Gambar 4.11. Hasil Kata Kunci Berdasarkan Hasil Analisis Manual dan Implementasi Teori Luhn

Hasil dari 25 kata tersebut yaitu '*contributions, control, criteria, crops, damage, difference, diversity, drainage, future, generally, governments, green, half, household, ice, petroleum, temperatures, priority, recycling, risk, root, simulation, soils, speed, states*'

4.4 Hasil Kata Kunci yang Didapatkan Berdasarkan Teknik *Clustering* dan Seleksi Teori Luhn

Telah dihasilkan berupa kata-kata kunci dari hasil teknik *Clustering* dan seleksi menggunakan Teori Luhn, yaitu:

Tabel 4.1. Hasil Kata Kunci yang Terpilih Berdasarkan Teknik *Clustering* dan Implementasi Teori Luhn

Kata Kunci Terpilih	
Cluster 0 Keywords	<i>emissions, indonesia, greenhouse, carbon, ghg, study, emission, gas, climate, change, co2, results, development, research, countries, data, production, impact, food, impacts</i>
Cluster 1 Keywords	<i>emissions, forest, soil, peat, rice, ch4, indonesia, greenhouse, land, gas, study, carbon, water, ghg, n2o, emission, area, peatlands, management, palm</i>
Cluster 2 Keywords	<i>oil, palm, emissions, indonesia, production, gas, greenhouse, biodiesel, industry, ghg, study, soil, land, forest, peat, plantations, n2o, carbon, co2, impact</i>
Cluster 3 Keywords	<i>energy, power, indonesia, gas, emissions, electricity, production, greenhouse, co2, study, generation, emission, plant, consumption, ghg, coal, system, use, bio-gas, efficiency</i>
25 Kata dengan Teori Luhn	<i>contributions, control, criteria, crops, damage, difference, diversity, drainage, future, generally, governments, green, half, household, ice, petroleum, temperatures, priority, recycling, risk, root, simulation, soils, speed, states</i>

Jika dilakukan suatu pemilihan kata-kata kunci dari kedua hasil proses tersebut, berdasarkan sumber dari data studi literatur, aspek dan faktor penyebab dari Gas Rumah Kaca yang ada di Indonesia, maka dari hasil tersebut bisa dijadikan sebagai landasan secara efektif untuk melihat kondisi saat ini dengan memanfaatkan studi literatur 5 tahun terakhir yang membahas terkait subjek atau urgensi perihal GRK di Indonesia. Contohnya seperti *electricity consumption, countries emission, biodiesel production, energy use* dan kata-kata yang memiliki keterkaitan antara sesama kata-kata kunci yang bisa digabungkan menjadi frasa seperti contoh tersebut, adapun pada penelitian ini terdapat berbagai banyak frasa lainnya yang dapat dibentuk dari hasil klasterisasi dan juga kata-kata kunci dari teori Luhn (lihat halaman Lampiran B). Oleh karena itu, dari hasil penelitian ini, bisa menjadi salah satu perhatian bahwa dengan memanfaatkan teknologi saat ini seperti *Artificial Intelligence* dan juga analisis data yang didukung juga berdasarkan teori, dapat membantu dalam analisis tersebut.



BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Dari hasil analisis data teks menggunakan pengolahan dan teknik dengan *text mining*, *clustering*, dan penyeleksian menggunakan teori Luhn pada data abstrak dari studi literatur lima tahun terakhir, telah didapatkan kata-kata kunci yang memiliki relevan sebagai faktor penyebab Gas Rumah Kaca (GRK) di Indonesia. Hasil analisis ini dapat memberikan informasi yang bermanfaat terkait permasalahan GRK saat ini dan dapat menjadi landasan untuk langkah-langkah penanganan dan mitigasi lebih lanjut terhadap dampak GRK di Indonesia.

5.2 Saran

Untuk penelitian selanjutnya, disarankan agar dapat memperluas sumber data studi literatur dengan mencakup lebih banyak publikasi terkini dan diversifikasi sumber literatur, termasuk jurnal ilmiah, laporan pemerintah, dan publikasi akademis lainnya. Penggunaan sumber data yang lebih luas akan memberikan gambaran yang lebih komprehensif dan akurat terkait faktor penyebab GRK di Indonesia.

Selanjutnya, dapat dilakukan kajian lebih lanjut dengan mempertimbangkan aspek waktu, geografis, dan perubahan iklim, sehingga dapat memberikan pemahaman yang lebih holistik tentang dampak GRK di Indonesia. Dengan menggunakan analisis data yang lebih maju dan inklusif, pemanfaatan teknologi tepat, hasil yang diharapkan dapat memberikan kontribusi yang lebih signifikan dalam pengembangan strategi dan kebijakan mitigasi GRK yang lebih efektif dan berkelanjutan di Indonesia.



DAFTAR PUSTAKA

- [1] M. A. Rizaty, “Emisi Gas Rumah Kaca Indonesia Diproyeksi Terus Naik hingga 2030,” *DataIndonesia.Id*, Oct. 2022. [Online]. Available: <https://dataindonesia.id/varia/detail/emisi-gas-rumah-kaca-indonesia-diproyeksi-terus-naik-hingga-2030>
- [2] G. A. Kristanto and W. Koven, “Estimating greenhouse gas emissions from municipal solid waste management in Depok, Indonesia,” *City and Environment Interactions*, vol. 4, p. 100027, Dec. 2019. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S2590252020300088>
- [3] D. G. K. Ketaren, “PERANAN KAWASAN MANGROVE DALAM PENURUNAN EMISI GAS RUMAH KACA DI INDONESIA,” *Jurnal Kelautan dan Perikanan Terapan (JKPT)*, vol. 1, p. 73, Jan. 2023. [Online]. Available: <http://ejournal-balitbang.kkp.go.id/index.php/jkpt/article/view/12050>
- [4] J. Li, M. Irfan, S. Samad, B. Ali, Y. Zhang, D. Badulescu, and A. Badulescu, “The Relationship between Energy Consumption, CO2 Emissions, Economic Growth, and Health Indicators,” *International Journal of Environmental Research and Public Health*, vol. 20, no. 3, p. 2325, Jan. 2023. [Online]. Available: <https://www.mdpi.com/1660-4601/20/3/2325>
- [5] M. Fajri, I. Adi Nugroho, and T. Tri Hendro Atmoko Utomo, “Standar lhk, penunjang target pencapaian folu net sink 2030,” *STANDAR: Better Standard Better Living*, vol. 2, no. 3, p. 14–18, Mei 2023. [Online]. Available: <http://majalah.bsilhk.menlhk.go.id/index.php/STANDAR/article/view/124>
- [6] S. P. Tiwari, S. Prasad, and M. Thushara, “Machine learning for translating pseudocode to python: A comprehensive review,” in *2023 7th International Conference on Intelligent Computing and Control Systems (ICICCS)*, May 2023, pp. 274–280.
- [7] S. Salloum, M. Al-Emran, A. Monem, and K. Shaalan, *Using Text Mining Techniques for Extracting Information from Research Articles*, 01 2018, pp. 373–397.
- [8] H. P. Luhn, “The Automatic Creation of Literature Abstracts,” *IBM Journal of Research and Development*, vol. 2, no. 2, pp. 159–165, Apr. 1958. [Online]. Available: <http://ieeexplore.ieee.org/document/5392672/>
- [9] Computer Education Department, ACLC College of Butuan, Butuan City, Philippines, V. Z. V. Singco, J. C. Trillo, C. C. Abalorio, J. C. M. Bustillo, J. T. Bojocan, and M. C. Elape, “OCR-based Hybrid Image Text Summarizer using Luhn Algorithm with FinetuneTransformer Modelsfor Long Document,” *International Journal of Emerging Technology and Advanced Engineering*, vol. 13, no. 2, pp. 47–56, Feb. 2023. [Online]. Available: https://www.ijetae.com/files/Volume13Issue2/IJETAE_0223_07.pdf
- [10] A. Rahmah, H. B. Santoso, and Z. A. Hasibuan, “Critical Review of Technology-Enhanced Learning using Automatic Content Analysis,” *International Journal of Advanced Computer Science and Applications*, vol. 13, no. 1, 2022. [Online]. Available: <http://thesai.org/Publications/ViewPaper?Volume=13&Issue=1&Code=IJACSA&SerialNo=48>
- [11] H. Wang, J. Luo, M. Zhang, and Y. Ling, “The Impact of Transportation Restructuring on the Intensity of Greenhouse Gas Emissions: Empirical Data from China,” *International Journal of Environmental Research and Public Health*, vol. 19, no. 19, p. 12960, Oct. 2022. [Online]. Available: <https://www.mdpi.com/1660-4601/19/19/12960>

- [12] E. Purnamasari, S. Sudarno, and H. Hadiyanto, "INVENTARISASI EMISI GAS RUMAH KACA SEKTOR PERTANIAN DI KABUPATEN BOYOLALI," *Universitas Muhammadiyah Surakarta*, Apr. 2019.
- [13] R. Pratama, "EFEK RUMAH KACA TERHADAP BUMI," vol. 14, no. 2, 2019.
- [14] K. O. Yoro and M. O. Daramola, "CO₂ emission sources, greenhouse gases, and the global warming effect," in *Advances in Carbon Capture*. Elsevier, 2020, pp. 3–28. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/B9780128196571000013>
- [15] K. Musari and F. Sayah, "Green Financing through Green Sukuk in the Fight Against Climate Change: Lessons from Indonesia," 2021, publisher: Unpublished. [Online]. Available: <http://rgdoi.net/10.13140/RG.2.2.23804.26245>
- [16] A. Aqlan, D. M. Bairam, and R. L. Naik, *A Study of Sentiment Analysis: Concepts, Techniques, and Challenges*, 01 2019, pp. 147–162.
- [17] S. Chung, S. Moon, J. Kim, J. Kim, S. Lim, and S. Chi, "Comparing natural language processing (nlp) applications in construction and computer science using preferred reporting items for systematic reviews (prisma)," *Automation in Construction*, vol. 154, p. 105020, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0926580523002807>
- [18] S. Moon, S. Chi, and S.-B. Im, "Automated detection of contractual risk clauses from construction specifications using bidirectional encoder representations from transformers (bert)," *Automation in Construction*, vol. 142, p. 104465, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0926580522003387>
- [19] M. Nasimuzzaman, A. N. Merag, S. Afroj, M. M. Alam, M. H. K. Mehedi, and A. A. Rasel, "Consumer review analysis using nlp and data mining," in *2023 IEEE 13th Annual Computing and Communication Workshop and Conference (CCWC)*, March 2023, pp. 0426–0430.
- [20] Y. Setiawan, D. Gunawan, and R. Efendi, "Feature extraction tf-idf to perform cyberbullying text classification: A literature review and future research direction," in *2022 International Conference on Information Technology Systems and Innovation (ICITSI)*. IEEE, 2022, pp. 283–288.
- [21] K. L. Tan, C. P. Lee, and K. M. Lim, "A survey of sentiment analysis: Approaches, datasets, and future research," *Applied Sciences*, vol. 13, no. 7, p. 4550, Apr 2023. [Online]. Available: <http://dx.doi.org/10.3390/app13074550>
- [22] T. Kurita, "Principal component analysis (pca)," *Computer Vision: A Reference Guide*, pp. 1–4, 2019.
- [23] J. R. Beattie and F. W. L. Esmonde-White, "Exploration of Principal Component Analysis: Deriving Principal Component Analysis Visually Using Spectra," *Applied Spectroscopy*, vol. 75, no. 4, pp. 361–375, Apr. 2021. [Online]. Available: <http://journals.sagepub.com/doi/10.1177/0003702820987847>
- [24] A. F. Jahwar and A. M. Abdulazeez, "Meta-heuristic algorithms for k-means clustering: A review," *PalArch's Journal of Archaeology of Egypt/Egyptology*, vol. 17, no. 7, pp. 12 002–12 020, 2020.
- [25] A. M. Ikotun, A. E. Ezugwu, L. Abualigah, B. Abuhaija, and J. Heming, "K-means clustering algorithms: A comprehensive review, variants analysis, and advances in the era of big data," *Information Sciences*, vol. 622, pp. 178–210, Apr. 2023. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0020025522014633>

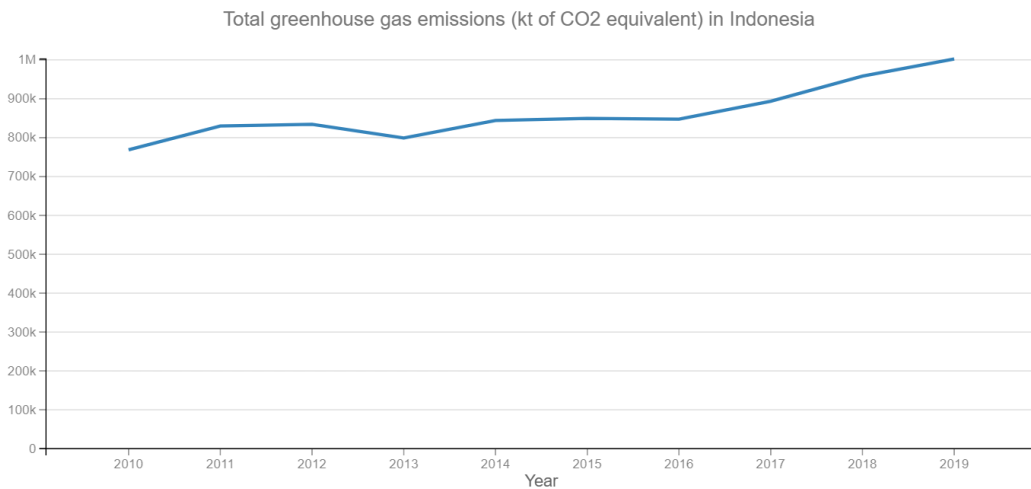
- [26] B. Kitchenham, O. Pearl Brereton, D. Budgen, M. Turner, J. Bailey, and S. Linkman, "Systematic literature reviews in software engineering – A systematic literature review," *Information and Software Technology*, vol. 51, no. 1, pp. 7–15, Jan. 2009. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0950584908001390>
- [27] S. Pizard, F. Acerenza, D. Vallespir, and B. Kitchenham, "Assessing attitudes towards evidence-based software engineering in a government agency," *Information and Software Technology*, vol. 154, p. 107101, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584922002105>
- [28] L. S. L. Lai, "CONTENT ANALYSIS OF SOCIAL MEDIA: A GROUNDED THEORY APPROACH," *Social Media Content Analysis*, vol. 16, no. 2, 2015.
- [29] A. Yunita, H. B. Santoso, and Z. A. Hasibuan, "Finding Contributing Factors of Students' Academic Achievement Using Quantitative and Qualitative Analyses-Based Information Extraction," *International Journal of Emerging Technologies in Learning (iJET)*, vol. 17, no. 16, pp. 108–125, Aug. 2022. [Online]. Available: <https://online-journals.org/index.php/i-jet/article/view/31945>
- [30] S. Lee, G. Hahn, J. Hecker, S. M. Lutz, K. Mullin, A. D. N. I. (ADNI), W. Hide, L. Bertram, D. L. DeMeo, R. E. Tanzi, C. Lange, and D. Prokopenko, "A comparison between similarity matrices for principal component analysis to assess population stratification in sequenced genetic data sets," *Briefings in Bioinformatics*, vol. 24, no. 1, p. bbac611, 12 2022. [Online]. Available: <https://doi.org/10.1093/bib/bbac611>
- [31] R. M. M. Hicke, M. Goenka, and E. Alexander, "Word Clouds in the Wild," in *2022 IEEE 7th Workshop on Visualization for the Digital Humanities (VIS4DH)*. Oklahoma City, OK, USA: IEEE, Oct. 2022, pp. 43–48. [Online]. Available: <https://ieeexplore.ieee.org/document/9974812/>
- [32] World Bank, "World development indicators databank," The World Bank, Washington, DC, Tech. Rep., Jul. 2023. [Online]. Available: <https://databank.worldbank.org/source/world-development-indicators>



LAMPIRAN A

Lampiran 1. Tren Total Emisi GRK Indonesia

Data yang diambil dari *World Development Indicators* oleh World Bank [32], menunjukkan *total greenhouse gas emissions (kt of CO₂ equivalent)* di Indonesia tahun 2010-2019. Mengindikasikan bahwa tren dari emisi yang ditimbulkan oleh GRK di Indonesia semakin meningkat secara signifikan, pada tahun 2019 dihasilkan sebesar 1.002.369,0 kt of CO₂.



Gambar 1.1. Total emisi GRK di Indonesia 2010-2019

Cluster 3:

Cluster 1:

[illegible][illegible][illegible]

```
[('emissions contribution', 0.6262), ('gas emissions', 0.6186), ('emissions petroleum', 0.6172), ('emissions oil', 0.6121), ('emissions importance', 0.61)]]
```

Gambar 2.2. Hasil frasa yang dapat dibentuk (2)

Adapun setiap elemennya, terdapat sebuah tupel yang berisi frasa dan nilai desimal, pada elemen kedua atau nilai desimal tersebut menunjukkan tingkat kepentingan atau relevansi frasa tersebut dalam teks dalam rentang nilainya 0 sampai 1. Semakin tinggi nilai desimal, semakin penting dan relevan frasa tersebut dalam konteks teks yang sedang dianalisis.