

# Pemanfaatan Transformer untuk Peringkasan Teks: Studi Kasus pada Transkripsi Video Pembelajaran

Muhammad Furqon Fadlilah\*, Aldy Rialdy Atmadja, Muhammad Deden Firdaus

Fakultas Sains dan Teknologi, Teknik Informatika, Universitas Islam Negeri Sunan Gunung Djati, Bandung, Indonesia

Email: <sup>1,\*</sup>1207050072@student.uinsgd.ac.id, <sup>2</sup>aldyrialdy@uinsgd.ac.id, <sup>3</sup>deden@uinsgd.ac.id

Email Penulis Korespondensi: 1207050072@student.uinsgd.ac.id

Submitted: 25/11/2024; Accepted: 29/12/2024; Published: 30/12/2024

**Abstrak**—Dalam era digital, video pembelajaran semakin banyak digunakan, namun sering kali berisi informasi yang kurang relevan, sehingga menyulitkan pemahaman kontennya. Penelitian ini mengusulkan pendekatan berbasis model *Whisper* dan T5 untuk menghasilkan ringkasan teks dari hasil transkripsi video pembelajaran di *YouTube*. *Whisper* digunakan untuk transkripsi suara ke teks, dengan fokus pada varian model yang menawarkan *Word Error Rate* (WER) rendah dan efisiensi waktu. Selanjutnya, model T5 dilatih ulang (*fine-tuned*) untuk menghasilkan ringkasan teks yang akurat, dengan strategi pembagian transkrip menjadi segmen untuk mengatasi batasan panjang input model. Praproses teks tidak digunakan karena menghasilkan kualitas evaluasi yang lebih baik. Hasil penelitian menunjukkan bahwa kombinasi *Whisper Turbo* dan model T5 yang telah dioptimalkan memberikan performa terbaik, dengan nilai *F1-Score* pada metrik ROUGE sebesar 39,23 (ROUGE-1), 13,17 (ROUGE-2), dan 23,84 (ROUGE-L). Pendekatan ini berhasil menghasilkan ringkasan teks yang lebih relevan dan komprehensif, serta meningkatkan efektivitas pembelajaran berbasis video. Dengan demikian, penelitian ini memberikan kontribusi signifikan dalam pengembangan teknologi peringkasan teks pada video pembelajaran.

**Kata Kunci:** Model Whisper; Ringkasan Teks; Recall-Oriented Understudy for Gisting Evaluation; Text-to-Text Transfer Transformer; Transkripsi Video

**Abstract**—In the digital era, learning videos are increasingly being used, however, they often contain irrelevant information, making it difficult to comprehend the content. This study proposes an approach based on the *Whisper* and T5 models to generate text summaries from *YouTube* educational video transcripts. *Whisper* is used for speech-to-text transcription, focusing on model variants that offer a low *Word Error Rate* (WER) and time efficiency. Subsequently, the T5 model is fine-tuned to produce accurate text summaries, with a strategy of segmenting the transcript to address input length limitations. Text preprocessing is not applied as it resulted in better evaluation quality. The results show that the combination of *Whisper Turbo* and the optimized T5 model provides the best performance, with *F1-Scores* on the ROUGE metrics of 39.23 (ROUGE-1), 13.17 (ROUGE-2), and 23.84 (ROUGE-L). This approach successfully generates more relevant and comprehensive text summaries, enhancing the effectiveness of video-based learning. Therefore, this research makes a significant contribution to the development of text summarization technology for learning videos.

**Keywords:** Whisper Model; Text Summary; Recall-Oriented Understudy for Gisting Evaluation; Text-to-Text Transfer Transformer; Video Transcription;

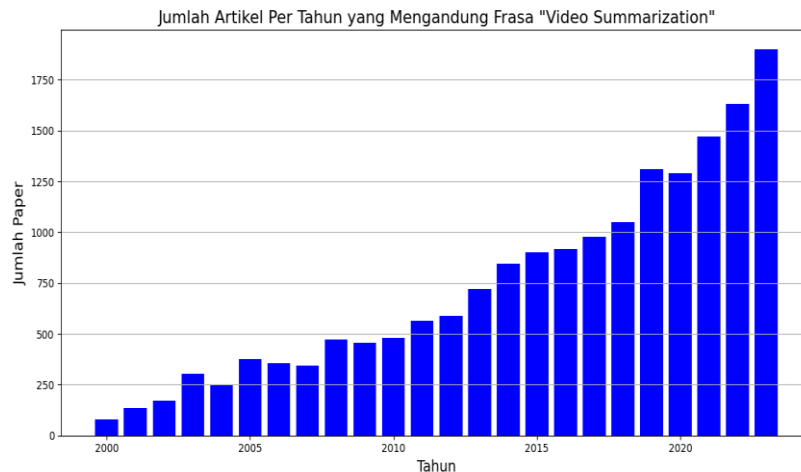
## 1. PENDAHULUAN

Media pembelajaran memiliki pengaruh signifikan terhadap proses belajar-mengajar. Revolusi industri dan revolusi pendidikan telah terjadi secara bersamaan dan berdampak pada kehidupan manusia. Revolusi industri mempengaruhi perubahan dalam pendidikan. Dengan perkembangan teknologi, didapatkan sebuah tantangan baru bagi para pengajar untuk memanfaatkan teknologi yang ada secara optimal. Salah satu bentuk pemanfaatan teknologi tersebut adalah melalui penggunaan media pembelajaran [1]. Salah satu metode yang menjadi tren dalam *e-learning* pada dekade ini adalah pembelajaran berbasis video [2]. Namun untuk meningkatkan efektivitas pembelajaran, harus disajikan informasi inti serta penghapusan informasi tambahan atau informasi yang tidak relevan dalam video [3].

Kecerdasan buatan (AI) dan pembelajaran mesin menawarkan solusi inovatif untuk mengatasi masalah ini. *Video summarization* sangat efektif dalam menghasilkan informasi penting dari video berdurasi panjang dalam waktu singkat [4], karena teks yang ringkas dapat mengurangi waktu membaca serta ukuran *file*, memberikan keuntungan bagi pengguna maupun sistem komputer [5].

Disamping itu, penelitian pada peringkasan video terus meningkat. Peningkatan ini menunjukkan bahwa penelitian di bidang ini semakin diminati dan relevan, mencerminkan upaya para peneliti untuk mengatasi tantangan dalam menyajikan informasi secara efisien dalam konteks pembelajaran berbasis video. Video yang tidak diringkas sering kali menghabiskan waktu pengguna untuk menemukan informasi yang relevan, terutama dalam konteks pembelajaran. Oleh karena itu, penelitian terkait peringkasan video tidak hanya menjadi lebih diminati tetapi juga semakin penting untuk meningkatkan efisiensi dan aksesibilitas informasi

Gambar 1 menunjukkan grafik dari jumlah artikel yang mengandung frasa *video summarization* meningkat seiring dengan pertumbuhan pesat konten video yang dihasilkan *platform YouTube*. Dengan demikian, kebutuhan untuk meringkas video secara efisien semakin meningkat dan mendorong lebih banyak penelitian dalam bidang peringkasan video.



**Gambar 1.** Grafik Jumlah Artikel Per Tahun yang Mengandung Frasa *Video Summarization* pada *Google Scholar*

Dalam upaya mewujudkan proses peringkasan video yang efisien, Penggunaan model *whisper* menjadi komponen penting dalam penelitian ini, karena *OpenAI* menggambarkan *Whisper* sebagai sistem ASR yang sangat maju. *Whisper* juga menyertakan berbagai metode penyaringan otomatis yang bertujuan meningkatkan kualitas transkrip, sehingga membantu membuat sistem pengenalan suara-ke-teks semakin baik [6]. Selain itu, peringkasan teks dilakukan menggunakan model T5 yang berhasil mencapai performa tinggi melalui pelatihan yang efektif, dengan kinerjanya diukur menggunakan metrik skor ROUGE. Model ini yang diusulkan menghasilkan ringkasan dengan bantuan representasi bahasa yang lebih baik dalam jaringan saraf [7], serta dengan pelatihan menggunakan dataset Liputan6 diharapkan dapat beradaptasi lebih baik dalam Bahasa Indonesia.

Peringkasan teks hasil transkrip video menggunakan arsitektur transformer telah banyak dilakukan pada penelitian-penelitian sebelumnya. Terdapat penelitian dengan judul “*YouTube Video Summarizer using NLP: A Review*” yang menjadikan NLP sebagai peran penting dalam membuat navigasi konten *YouTube* lebih cepat, akurat, dan mudah, sehingga mengoptimalkan pengalaman kita dalam konsumsi konten digital serta memiliki tantangan dan potensi di masa depan untuk peningkatan metode peringkasan [8].

Trankripsi teks menggunakan model *Whisper* merupakan metode yang efektif, meninjau penelitian sebelumnya menggunakan model *Whisper OpenAI*, yang diintegrasikan dengan deteksi batas heuristik. *Whisper* berhasil mencapai *Word Error Rate* (WER) sebesar 14,50% dan *Character Error Rate* (CER) sebesar 8,13% [9], menunjukkan potensinya dalam aplikasi konversi suara ke teks yang akurat. Disamping itu, telah dilakukan penelitian pada peringkasan dengan model T5 yang dilakukan *fine-tuning* menggunakan dataset INDOSUM dengan data artikel-ringkasan sebanyak 19.000 data. Memiliki nilai evaluasi terbaik dengan nilai ROUGE-1 0.17568 [10]. Pada penelitian terdahulu lainnya dilakukan peringkasan teks abstraktif menggunakan model *Bidirectional And Auto-Regressive Transformer* (BART) dengan pelatihan menggunakan dataset Liputan6 didapatkan nilai ROUGE-1, ROUGE-2, dan ROUGE-L masing-masing 37,19, 14,03, dan 33,85 [11].

Berdasarkan penelitian sebelumnya, penelitian ini bertujuan untuk mengembangkan metode peringkasan teks dari transkrip video *YouTube* dengan mengintegrasikan model *Whisper* dan T5 yang telah disesuaikan melalui *fine-tuning* menggunakan dataset Liputan6. Penggunaan model *Whisper* memungkinkan ekstraksi transkrip dengan tingkat akurasi yang tinggi, sementara model T5, yang telah dilatih dengan dataset tersebut, diharapkan dapat menghasilkan ringkasan yang lebih akurat dan relevan dalam konteks bahasa Indonesia. Penelitian ini difokuskan pada evaluasi tingkat akurasi kombinasi kedua model dan identifikasi varian *Whisper* yang paling optimal dalam menghasilkan transkrip yang akurat untuk proses peringkasan menggunakan model T5.

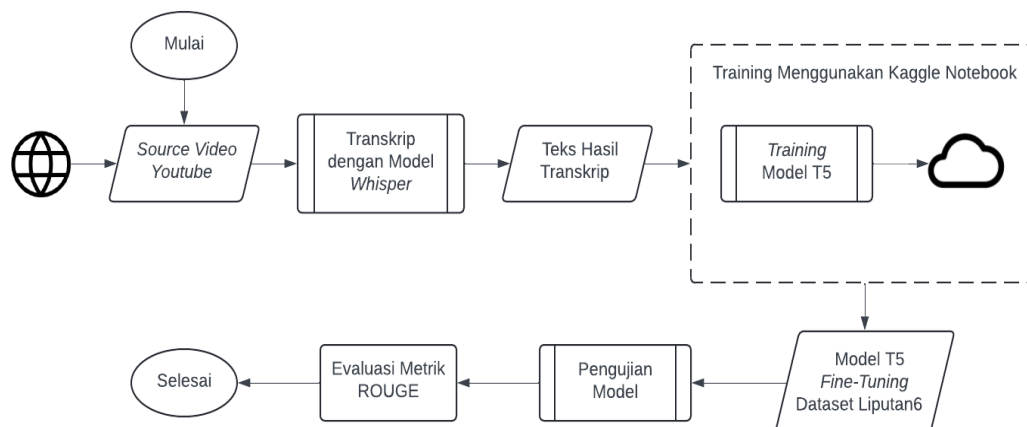
## 2. METODOLOGI PENELITIAN

Penelitian ini dimulai dengan proses transkripsi teks menggunakan model *Whisper* untuk mengonversi *audio video* menjadi teks, yang kemudian diringkas dengan model T5 yang telah dilakukan *fine-tuning* menggunakan dataset Liputan6. Evaluasi dilakukan dengan mengukur akurasi ringkasan menggunakan metrik ROUGE untuk menilai efektivitas kombinasi model *Whisper* dan T5 dalam menghasilkan ringkasan.

### 2.1. Tahapan Penelitian

Proses penelitian mencakup serangkaian langkah yang ditempuh peneliti dari awal hingga akhir penelitian. Tahapan penelitian dimulai dengan pengambilan video dari platform *YouTube*. Video ini kemudian diproses menggunakan model *Whisper* untuk menghasilkan transkrip teks. Selanjutnya, teks hasil transkripsi digunakan sebagai input dalam proses pelatihan model T5, yang dilakukan menggunakan *Kaggle Notebook* sebagai platform berbasis *cloud*. Model T5 dilatih menggunakan dataset Liputan6 yang sudah disiapkan. Setelah pelatihan selesai, model diuji untuk

menghasilkan ringkasan dari teks transkripsi, dan hasilnya dievaluasi menggunakan metrik ROUGE untuk mengukur kualitas ringkasan. Tahapan ini dirancang untuk memastikan setiap proses berjalan efisien dalam menguji kombinasi antara model *Whisper* dan T5 pada peringkasan teks hasil transkripsi video pembelajaran. Tahapan-tahapan penelitian ini dijelaskan secara rinci pada Gambar 2.



Gambar 2. Flowchat Tahapan Penelitian

## 2.2. Source YouTube Video

Video dari *YouTube* digunakan untuk menguji efektivitas kombinasi model *Whisper* dan T5 dalam menghasilkan ringkasan akurat dan berkualitas tinggi. Video yang dipilih mencakup topik-topik seperti ilmu pengetahuan, teknologi, kesehatan, dan isu sosial. Sebelum proses transkripsi dilakukan, video tersebut dianalisis berdasarkan struktur dan kualitas bahasanya guna meminimalkan potensi kesalahan transkripsi oleh model *Whisper*, khususnya dalam aspek fonetik atau pelafalan bahasa Indonesia. Kurasi ini memastikan bahwa teks transkripsi yang dihasilkan memiliki kualitas optimal, sehingga mendukung evaluasi kinerja model T5 dalam menghasilkan ringkasan yang informatif dan relevan.

## 2.3. Trankripsi Teks Menggunakan Model *Whisper*

*Whisper* menggunakan metode *Large-Scale Weak Supervision* untuk mendeteksi audio melalui banyak kumpulan data dari berbagai sumber, lingkungan rekaman, pembicara, dan bahasa [6]. Rangkaian model ini mengadopsi arsitektur *encoder-decoder*, dengan ukuran parameter yang bervariasi mulai dari 39 juta (model kecil) hingga 1,55 miliar (model besar), dan mendukung lebih dari 98 bahasa dalam fitur multibahasanya. Format datanya dirancang secara khusus agar *Whisper* dapat menjalankan berbagai tugas dengan tingkat fleksibilitas dan ketangguhan yang tinggi. Pada *OpenAI-Whisper*, *encoder* dan *decoder* disediakan untuk mendukung proses transkripsi teks menggunakan teknologi transformers yang sangat efektif [12].

Seluruh varian model *Whisper*, mulai dari *tiny* hingga *turbo*, diimplementasikan dan diuji untuk menentukan model paling optimal yang dapat memberikan hasil transkripsi teks yang akurat dan efisien, sekaligus mendukung kualitas peringkasan teks dengan skor ROUGE terbaik saat dikombinasikan dengan model T5 yang disempurnakan menggunakan dataset Liputan6. Pengujian melibatkan evaluasi ketepatan transkripsi serta efisiensi waktu pemrosesan dari setiap model. Melalui pengukuran akurasi transkripsi dan skor ROUGE yang dihasilkan, varian model *Whisper* yang mencapai keseimbangan terbaik antara kecepatan dan ketelitian, sekaligus mendukung peringkasan berkualitas tinggi, akan dipilih sebagai model paling unggul dalam memenuhi tujuan penelitian ini.

## 2.4. Teks Hasil Transkrip

Teks hasil transkripsi dari model *Whisper* akan menjadi dasar untuk mengevaluasi kualitas peringkasan yang dihasilkan. Transkripsi dari masing-masing varian model *Whisper* dari model *tiny* hingga *turbo* akan dikombinasikan dengan model T5 yang telah dilatih dengan dataset Liputan6 untuk proses peringkasan. Setiap ringkasan yang dihasilkan dari teks transkripsi ini akan dianalisis menggunakan metrik ROUGE untuk mengukur tingkat akurasi dan kesesuaian ringkasan terhadap inti informasi dari teks aslinya.

## 2.5. Text-to-Text Transfer Transformer (T5) dengan *Fine-Tuning* Dataset Liputan6

*Text to Text Transfer Transformer* (T5) merupakan model utama dalam percobaan yang diperkenalkan oleh Colin Raffel [13]. T5 merupakan kerangka kerja berbasis transformer yang mengubah teks ke teks, yang menunjukkan kinerja luar biasa dalam berbagai tugas NLP, termasuk peringkasan teks, tanya jawab, dan penerjemahan mesin [14]. Dengan kemampuan T5 dalam memahami konteks dan mengekstrak informasi utama, model ini diharapkan dapat menghasilkan ringkasan transkrip video pembelajaran yang koheren dan informatif. Evaluasi kinerja menggunakan

metrik tertentu dilakukan untuk memastikan hasil ringkasan memenuhi standar dan memberikan nilai tambah dalam pembelajaran.

Model T5 berhasil mencapai performa tinggi melalui pelatihan yang efektif, dengan kinerjanya diukur menggunakan metrik skor ROUGE. Model ini yang diusulkan menghasilkan ringkasan dengan bantuan representasi bahasa yang lebih baik dalam jaringan saraf [7]. Model bahasa ini disiapkan untuk tugas peringkasan dengan cara menyempurnakan model menggunakan kumpulan data, yang memungkinkannya untuk mengambil informasi, belajar pola, dan menciptakan representasi internal yang dapat menghasilkan ringkasan yang logis dan teratur [15].

Pada arsitektur *transformer*, operasi *self-attention* berfungsi dengan menerima urutan input tertentu dan menghasilkan *output* baru dengan panjang yang sama. Proses ini melibatkan perhitungan kontribusi elemen-elemen dalam urutan *input* yang membentuk setiap elemen pada urutan *output*. Secara lebih rinci, elemen ke-*i* dari *output* *yi* dibentuk melalui kombinasi tertimbang dari seluruh elemen *input* *xj*. Penentuan bobot ini bergantung pada mekanisme perhatian, yang secara adaptif menentukan sejauh mana elemen *input* tertentu relevan terhadap elemen *output* yang sedang dihasilkan [16].

Dalam penelitian ini, model T5 dilakukan *fine-tuning* menggunakan dataset Liputan6 yang berisi berita dalam Bahasa Indonesia. Dataset Liputan6 dipilih karena merupakan kumpulan teks yang kaya akan variasi topik dan memiliki struktur bahasa jurnalistik yang sesuai dengan konteks teks informatif dan faktual. Dataset ini memiliki dua tipe data yaitu *canonical* dan *xtreme*. Namun, pada penelitian ini data yang digunakan hanya data *canonical*, hal ini dikarenakan model yang dilakukan *fine-tuning* difokuskan pada bahasa Indonesia bukan untuk multibahasa.

## 2.6. Pengujian Model

Pada tahapan pengujian model ini, penelitian bertujuan untuk menilai kinerja model T5 yang telah dioptimasi melalui proses *fine-tuning* menggunakan dataset berita Liputan6. Tujuan pengujian ini adalah untuk mengevaluasi kualitas ringkasan yang dihasilkan oleh model T5 dari teks transkrip video yang diekstrak secara otomatis menggunakan model *Whisper*. Proses ini dilakukan untuk menghasilkan ringkasan teks yang akurat dan koheren serta mempertahankan esensi informasi dalam konteks transkrip video pembelajaran. Pengujian model dilakukan melalui dua skenario yang dirancang untuk mengamati dampak praproses terhadap akurasi dan kualitas ringkasan berdasarkan metrik ROUGE.

### 2.6.1. Pembagian Teks Menjadi Beberapa Segmen

Transkrip video umumnya memiliki panjang teks yang cukup besar sehingga melampaui kapasitas maksimal input token model T5 karena dalam implementasinya, mayoritas model T5 dibatasi pada panjang sekuens 512 token akibat kendala komputasional[17]. Maka dari itu, setiap transkrip panjang dibagi menjadi segmen-segmen yang lebih pendek, dengan masing-masing segmen memiliki batas maksimal 512 token. Panjang maksimal segmen ini ditentukan dengan tujuan mempertahankan informasi yang cukup dalam tiap bagian teks, sehingga model dapat memproses tiap segmen tanpa kehilangan konteks utama dari isi transkrip. Pembatasan token ini juga memastikan bahwa setiap segmen sesuai dengan kapasitas maksimal input model T5, sehingga model dapat berfungsi optimal dalam proses peringkasan.

### 2.6.2. Skenario Pengujian

Setelah teks transkrip dibagi menjadi segmen-segmen sesuai batas token, dua skenario pengujian diterapkan untuk mengkaji efek praproses terhadap kualitas ringkasan. Pengujian ini mencakup dua pendekatan sebagai berikut:

#### a. Segmentasi dengan Praproses

Dalam skenario pertama, setiap segmen yang telah dibagi dengan panjang maksimal 512 token akan diproses melalui tahapan praproses sebelum melalui proses peringkasan oleh model T5. Tahapan praproses ini mencakup beberapa langkah, diantaranya konversi seluruh teks menjadi huruf kecil (*lowercasing*), Penghapusan tanda baca (*punctuation removal*), Penghapusan kata yang tidak memberikan kontribusi informasi penting (*stopword removal*) dan mereduksi kata menjadi bentuk dasarnya (*stemming*) [18].

#### b. Segmentasi Tanpa Praproses Teks

Pada skenario kedua, segmen teks sepanjang maksimal 512 token diringkas menggunakan model T5 tanpa melalui praproses, sehingga semua elemen teks, termasuk tanda baca, variasi kata, dan *stopword*, dipertahankan. Pendekatan ini mengevaluasi kemampuan model T5 untuk menangkap informasi esensial tanpa modifikasi tambahan pada teks masukan. Hasil ringkasan dari skenario ini dan skenario dengan praproses dievaluasi menggunakan metrik ROUGE, yang mengukur kualitas ringkasan berdasarkan kesamaan kata dan urutan frasa antara ringkasan otomatis dan acuan.

## 2.7. Metrik Evaluasi *Recall-Oriented Understudy for Gisting* (ROUGE)

ROUGE merupakan singkatan dari *Recall-Oriented Understudy of Gisting Evaluation*, yang mencakup sejumlah kriteria untuk mengevaluasi teks yang dihasilkan secara otomatis. Metode ini biasanya digunakan untuk menilai kualitas dari algoritma Peringkasan Teks [19]. ROUGE dipilih untuk menilai ringkasan karena telah menjadi standar yang banyak diterapkan dalam berbagai penelitian. Dua varian ROUGE yang digunakan adalah ROUGE-N dan ROUGE-L. ROUGE-N mengevaluasi *recall* berdasarkan *n-gram*, sedangkan ROUGE-L berfokus pada longest common subsequence yang paling panjang [11].



ROUGE-N adalah metrik yang mengukur recall berdasarkan perbandingan n-gram antara ringkasan referensi dan teks yang dihasilkan oleh algoritma peringkasan mesin. Jumlah n-gram yang dipertimbangkan bervariasi antara n=1 hingga n=4, tetapi n-gram yang paling sering digunakan adalah n=1 (ROUGE-1) dan n=2 (ROUGE-2). Jika p mewakili jumlah n-gram yang sama antara ringkasan referensi dan teks hasil mesin, dan q mewakili jumlah n-gram dalam ringkasan referensi [20], maka ROUGE-N dapat dihitung dengan rumus (1).

$$ROUGE - N = \frac{p}{q} \quad (1)$$

ROUGE-L mengukur kualitas ringkasan teks dengan cara membandingkan *longest common subsequence* (LCS), yang merupakan urutan kata terpanjang yang sama antara ringkasan mesin dan ringkasan standar. Jika m mewakili jumlah kata dalam ringkasan standar[11], maka perhitungan ROUGE-L seperti yang ditunjukkan pada formula (2) berikut.

$$ROUGE - L = \frac{LCS}{m} \quad (2)$$

Nilai ROUGE yang diukur dalam penelitian ini adalah hasil dari peringkasan teks yang dihasilkan oleh model T5 dari transkripsi yang dilakukan oleh berbagai varian model *Whisper*. Nilai ROUGE dari ringkasan ini akan menunjukkan seberapa baik kombinasi model *Whisper* dan T5 dalam menghasilkan ringkasan yang relevan dan akurat, dengan model *Whisper* yang memberikan transkripsi terbaik diharapkan memberikan nilai ROUGE tertinggi.

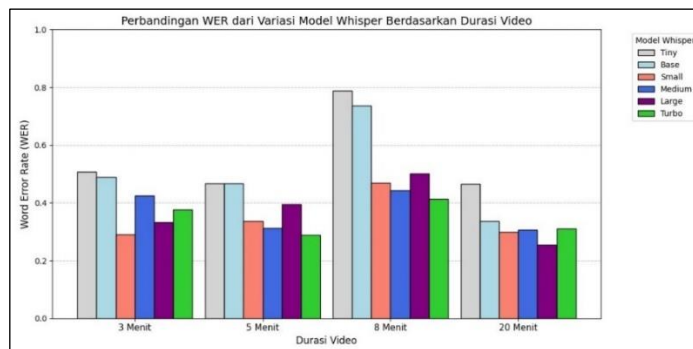
### 3. HASIL DAN PEMBAHASAN

#### 3.1. Transkripsi Teks dengan Model *Whisper*

Pengujian transkrip dilakukan menggunakan video dengan durasi yang berbeda-beda, yang semuanya mengandung unsur pembelajaran. Hal ini bertujuan untuk memahami kinerja berbagai varian model *Whisper* dalam proses transkripsi teks video pembelajaran. Video yang digunakan telah memiliki transkripsi bawaan, sehingga memungkinkan penghitungan *Word Error Rate* (WER) dengan membandingkan hasil transkripsi dari model *Whisper* terhadap transkripsi asli pada video tersebut. Hasil pengujian ini disajikan dalam Tabel 1 yang menunjukkan perbandingan variasi model *Whisper* untuk transkrip teks yang berasal dari video *YouTube*.

**Tabel 1.** Perbandingan Variasi Model *Whisper* Berdasarkan Waktu Proses Transkrip dan Nilai *Word Error Rate* (WER)

Judul Konten	Durasi Video	Model <i>Whisper</i>	Waktu Proses Transkrip	Word Error Rate (WER)
Gimana Cara Tau Jika AI Udah Punya Kesadaran?	3 menit	<i>Tiny</i>	14,09 detik	50,77%
		<i>Base</i>	16,77 detik	48,97%
		<i>Small</i>	<b>22,30 detik</b>	<b>29,12%</b>
		<i>Medium</i>	37,30 detik	42,53%
		<i>Large</i>	44,89 detik	32,22%
		<i>Turbo</i>	17,93 detik	37,63%
Beda Suku, Budaya, Dan Agama	5 Menit	<i>Tiny</i>	15,98 detik	46,73%
		<i>Base</i>	19,20 detik	44,49%
		<i>Small</i>	27,99 detik	33,67%
		<i>Medium</i>	44,99 detik	31,22%
		<i>Large</i>	61,77 detik	39,39%
		<i>Turbo</i>	<b>21,59 detik</b>	<b>28,78%</b>
BPUPKI - Sejarah Badan Penyelidik Usaha-Usaha Persiapan Kemerdekaan Indonesia	8 Menit	<i>Tiny</i>	44,26 detik	78,89%
		<i>Base</i>	32,49 detik	73,62%
		<i>Small</i>	45,96 detik	46,85%
		<i>Medium</i>	78,87 detik	44,27%
		<i>Large</i>	147,42 detik	50,06%
		<i>Turbo</i>	<b>32,73 detik</b>	<b>41,31%</b>
Sejarah Komputer Pertama I Penemuan Dan Perkembangan Komputer	20 Menit	<i>Tiny</i>	52,66 detik	46,53%
		<i>Base</i>	61,62 detik	33,70%
		<i>Small</i>	102,53 detik	29,77%
		<i>Medium</i>	180,77 detik	30,64%
		<i>Large</i>	553,09 detik	25,39%
		<i>Turbo</i>	<b>68,49 detik</b>	<b>30,96%</b>



Gambar 3. Grafik Perbandingan WER dari Variasi Model *Whisper* Berdasarkan Durasi Video

Pada Gambar 3 menunjukkan grafik perbandingan WER dari variasi model *Whisper* berdasarkan durasi video. Hasil tersebut memperlihatkan bahwa model *Whisper Turbo* terbukti paling efisien dengan waktu transkripsi tercepat dan nilai WER rendah, menjadikannya pilihan optimal untuk transkripsi teks video pembelajaran. Meskipun model *Whisper Small*, *Medium*, dan *Large* terkadang mencatat WER lebih rendah pada video tertentu, seperti pada video berdurasi 3 menit. Hal ini terjadi karena model *Turbo* mampu menangkap struktur kalimat dan detail konten lebih baik dibandingkan model *Small*, termasuk elemen seperti pengulangan kata "Intro..." seperti ditunjukkan pada Tabel 2. Namun, pada Tabel 2 model *Small* menghasilkan transkripsi yang lebih sederhana tanpa pengulangan tersebut. Sehingga dari perbandingan teks hasil transkripsi model *small* dan *turbo* memperlihatkan bahwa model *Turbo* lebih unggul dalam merepresentasikan isi dan nuansa materi video secara lebih lengkap.

Tabel 2. Tabel Perbandingan Teks Hasil Transkripsi pada Model *Small* dan *Turbo*

<b>Referensi Teks</b>	[Musik] di hidup yang penuh lik-liku ini ternyata bukan cuma si dia yang pengen dimengerti tapi juga Ai ya inilah yang pernah disampein sama si lamda Ai milik Google yang ngaku kalau dia itu punya kesadaran kayak manusia tapi sebenarnya ini bukan pernyataan resmi tapi hasil bocoran dari Karyawan Google yang berujung...
<b>Hasil Transkrip Model <i>Small</i></b>	Dibuat Di hidup yang penuh di kaliku ini ternyata bukan cuma dia yang pengen dia mengerti tapi juga Ei Iya, inilah yang pernah disampaikan sama si <i>Lambda</i> Ei milik <i>Google</i> yang ngaku kalau dia itu punya kesadaran kaya manusia Tapi sebenarnya ini bukan pernyataan resmi, tapi hasil bocoran dari karaman <i>Google</i> yang berujung...
<b>Hasil Transkrip Model <i>Turbo</i></b>	Intro Intro Intro Intro Intro Intro Intro Intro Di hidup yang penuh lika-liku ini ternyata bukan cuma si dia yang pengen dimengerti, tapi juga AI. Ya, inilah yang pernah disampaikan sama si <i>Lambda</i> , AI milik <i>Google</i> , yang ngaku kalau dia itu punya kesadaran kayak manusia. Tapi sebenarnya ini bukan pernyataan resmi, tapi hasil bocoran dari karyawan <i>Google</i> yang berujung dipecat. Oke, tapi sebenarnya ini bukan pernyataan resmi, tapi hasil bocoran dari karyawan <i>Google</i> yang berujung...

### 3.2. Text-to-Text Transfer Transformer (T5) Model

#### 3.2.1. Fine-Tuning Menggunakan Dataset Liputan6

Model *T5-Small* dilakukan *fine-tuning* menggunakan dataset Liputan6 yang terdiri dari 215.827 dataset. Namun pada proses *fine-tuning* digunakan model *T5-Small* dengan menggunakan 50.000 dataset dengan dilakukan dua skenario pembagian dataset untuk melatih dan menguji model. Skenario pertama menggunakan 80% data untuk pelatihan (*train*) dan 20% untuk pengujian (*test*), sementara skenario kedua membagi dataset menjadi 70% untuk pelatihan dan 30% untuk pengujian. Konfigurasi yang digunakan yaitu model *T5-Small* yang dilatih dengan banyak 10 epoch, menggunakan batch size sebesar 8 dan learning rate sebesar 0.0001. Konfigurasi ini dipilih untuk menjaga akurasi dan efisiensi model serta memastikan model dapat mempelajari pola data dari kedua skenario pembagian.

#### 3.2.2. Evaluasi Metrik ROUGE pada Proses *Fine-Tuning* Model T5

Evaluasi model T5 yang telah dilakukan *fine-tuning* dengan dataset Liputan6 dilakukan menggunakan metrik ROUGE. Dalam evaluasi ini, tiga varian metrik ROUGE yang dipertimbangkan adalah ROUGE-1, ROUGE-2, dan ROUGE-L, yang masing-masing mengukur kesesuaian pada tingkat unigram, bigram, serta kesesuaian urutan kata terpanjang (*longest common subsequence*) antara hasil ringkasan model dan referensi. Tabel 3 merupakan tabel yang menunjukkan hasil evaluasi metrik ROUGE untuk model T5 yang telah dilakukan *fine-tuning* menggunakan dataset Liputan6.

Tabel 3. Evaluasi Metrik Model T5 dengan *Fine-Tuning* Dataset Liputan6

Skenario	Model	Skor	ROUGE-1	ROUGE-2	ROUGE-L
----------	-------	------	---------	---------	---------

80% data train	T5-Small	Recall	57,71	24,17	51,20
20% data test		Precision	35,27	12,82	30,31
		F1-Score	43,41	16,64	37,79
70% data train	T5-Small	Recall	57,30	23,68	50,78
30% data test		Precision	36,52	13,55	31,69
		F1-Score	44,10	17,01	38,59

Pada skenario 2 atau penggunaan 70% data train dan 30% data test memberikan hasil yang lebih baik dalam evaluasi ROUGE. Struktur ini tampaknya memberikan keseimbangan yang baik antara pelatihan model dan evaluasi ketat untuk performa yang lebih optimal dalam peringkasan teks. Hal ini menunjukkan bahwa pembagian data dengan rasio 70:30 memberikan keseimbangan yang ideal antara proses pelatihan model dan evaluasi yang mendalam, sehingga mampu meningkatkan performa dalam tugas peringkasan teks. Kenaikan nilai F1-Score pada skenario ini dapat dihubungkan dengan jumlah data uji yang lebih besar, yang memungkinkan pengujian kemampuan generalisasi model secara lebih menyeluruh, meskipun data yang digunakan untuk pelatihan sedikit berkurang.

Karakteristik dataset Liputan6, seperti variasi topik dan pola struktur artikelnya yang konsisten, memiliki peran penting dalam meningkatkan kinerja model T5. Artikel dalam dataset ini umumnya memiliki bagian utama yang terorganisir dengan baik, sehingga memudahkan model dalam mengekstrak informasi penting selama proses pelatihan. Sifat ini membantu model lebih efektif memahami konteks dan menghasilkan ringkasan yang sesuai. Oleh karena itu, perpaduan antara pembagian data yang optimal dan karakteristik dataset yang mendukung secara signifikan berkontribusi pada peningkatan performa model secara keseluruhan.

### 3.3. Evaluasi Akhir Metrik ROUGE pada Peringkasan Teks Hasil Transkripsi Model Whisper

Sebelum mengevaluasi kombinasi model Whisper dan T5, dilakukan analisis perbandingan hasil ringkasan antara skenario dengan dan tanpa praproses teks. Analisis ini dilakukan setelah teks transkrip dibagi menjadi segmen-segmen sepanjang 512 token untuk memenuhi batas kapasitas pemrosesan input model T5, sehingga setiap segmen dapat diproses secara optimal tanpa kehilangan informasi penting. Tahapan ini bertujuan mengevaluasi kontribusi praproses terhadap kualitas ringkasan.

Secara umum, evaluasi metrik ROUGE yang dihasilkan dari perbandingan ringkasan hasil transkrip dengan praproses dan tanpa praproses menunjukkan bahwa hasil terbaik ditunjukkan pada peringkasan tanpa praproses yang mencapai nilai ROUGE-1 sebesar 39,23 dan ROUGE-2 sebesar 13,17. Skenario tanpa praproses terbukti optimal untuk menghasilkan ringkasan yang akurat dan lengkap. Sehingga, pada Tabel 4 dilakukan evaluasi metrik ROUGE dengan kombinasi model Whisper dan T5 pada peringkasan tanpa tahapan praproses. Evaluasi yang ditunjukkan pada tabel 4 juga menunjukkan bahwa kombinasi model Whisper dengan T5 yang telah dilakukan *fine-tuning* menggunakan dataset Liputan6 memberikan akurasi ringkasan yang bervariasi, tergantung pada varian Whisper yang digunakan. Hasil evaluasi metrik tersebut diperlihatkan secara detail pada Tabel 4 berikut.

**Tabel 4.** Evaluasi Metrik Kombinasi Variasi Model Whisper dan T5 dengan *Fine-Tuning* Dataset Liputan6

Model Whisper	Skor	ROUGE-1	ROUGE-2	ROUGE-L
Tiny	Recall	27,02	0,90	13,51
	Precision	19,60	0,65	9,80
	F1-Score	22,72	0,76	11,36
Base	Recall	34,23	10,00	20,72
	Precision	18,18	5,28	11,00
	F1-Score	23,75	6,91	14,37
Small	Recall	35,13	9,09	21,62
	Precision	22,94	5,91	14,11
	F1-Score	27,75	7,16	17,08
Medium	Recall	44,14	8,18	29,72
	Precision	28,99	5,35	19,52
	F1-Score	35,00	6,47	23,57
Large	Recall	49,54	10,90	27,02
	Precision	28,49	6,25	15,54
	F1-Score	36,18	7,94	19,73
Turbo	Recall	45,94	15,45	27,92
	Precision	34,22	11,48	20,80
	F1-Score	39,23	13,17	23,84

Sebagai ditunjukkan pada Tabel 4, evaluasi metrik ROUGE mengindikasikan bahwa model Whisper Turbo mencapai nilai *F1-Score* tertinggi pada ketiga metrik utama, yaitu ROUGE-1, ROUGE-2, dan ROUGE-L, masing-masing dengan skor sebesar 39,23, 13,17, dan 23,84 begitupun pada nilai *precision*, model Whisper Turbo mencatat nilai tertinggi pada ketiga metrik tersebut, dengan nilai ROUGE-1 sebesar 34,22, ROUGE-2 sebesar 11,48, dan



ROUGE-L sebesar 20,80. Meskipun model *Whisper Large* mencatat nilai *recall* yang lebih tinggi pada ROUGE-1 sebesar 49,54, serta model *Whisper Medium* memperoleh *recall* tertinggi pada ROUGE-L sebesar 29,72, penentuan model yang paling unggul dalam analisis ini didasarkan pada nilai *F1-Score* tertinggi. Model *Whisper Medium* menunjukkan performa *F1-Score* yang hampir sebanding dengan *Whisper Turbo*, yaitu sebesar 23,57. Kesetaraan ini menunjukkan bahwa model *Whisper Medium* mampu menjaga struktur dan konteks kalimat yang hampir sama dengan performa *Whisper Turbo*. Namun, mengingat bahwa *Whisper Turbo* secara signifikan lebih unggul dalam *F1-Score* pada ROUGE-1 dan ROUGE-2, model ini dapat dinilai lebih unggul dalam hal keseimbangan performa, terutama dalam hal akurasi penangkapan konsep utama dan relevansi keseluruhan ringkasan. Dengan demikian, model *Whisper Turbo* dinyatakan sebagai model *Whisper* yang paling optimal berdasarkan evaluasi metrik ROUGE dalam tugas peringkasan menggunakan model T5 yang telah dilakukan *fine-tuning* dengan dataset Liputan6.

Adapun pada tabel 5 ditunjukkan hasil transkripsi model *Whisper Turbo* menggunakan model T5 dengan *Fine-tuning* Dataset Liputan6. Hasil transkripsi menghasilkan struktur narasi yang lebih terorganisir dan mudah ditranskripsi oleh *Whisper Turbo*. Hasil transkripsi tersebut sudah merepresentasikan konten video yang sesuai, walaupun terdapat keterbatasan model T5 dalam menangkap konteks dari transkrip awal. Hal ini dapat diatasi dengan transkripsi yang dibagi kedalam beberapa segmen. Disamping itu, pada Tabel 5 diperlihatkan hasil peringkasan yang dihasilkan cukup optimal pada konten yang lebih kompleks atau panjang dengan mempergunakan *Fine-tuning* Dataset Liputan6.

**Tabel 5.** Ringkasan Teks Hasil Transkripsi model *Whisper Turbo* Menggunakan Model T5 dengan *Fine-tuning* Dataset Liputan6

<b>Teks Hasil Transkripsi Model <i>Whisper Turbo</i></b>	BPUPKI Pada tahun 1944, kedudukan tentara Jepang di medan Perang Pasifik makin terdesak. Di berbagai medan pertempuran, Jepang menderita kekalahan. Ditambah dengan timbulnya pemberontakan oleh rakyat Indonesia, maka kedudukan Jepang semakin terjepit. Pertahanan Jepang sudah rapuh dan bayangan kekalahan sudah semakin nyata. Namun Jepang masih berusaha menarik simpati rakyat Indonesia dengan menjanjikan kemerdekaan di kemudian hari. Pada 1 Maret 1945, pemerintah Jepang di Jawa dipimpin oleh...
<b>Referensi Ringkasan</b>	Pada Maret 1945, Jepang membentuk BPUPKI ( <i>Dokuritsu Junbi Chosakai</i> ) untuk menarik simpati rakyat Indonesia, dipimpin oleh Rajiman Wedio Diningrat. BPUPKI memiliki 60 anggota Indonesia dan 7 anggota Jepang tanpa hak suara. Sidang pertama BPUPKI, dari 29 Mei hingga 1 Juni 1945, membahas dasar negara yang diusulkan oleh tiga tokoh utama: Muhammad Yamin, Supomo, dan Soekarno. Pada sidang ini, Soekarno mengusulkan Pancasila sebagai dasar negara. Setelah sidang pertama, Panitia Sembilan dibentuk dan menghasilkan Piagam Jakarta pada 22 Juni 1945, yang memuat lima asas negara. Pada sidang kedua, 10-17 Juli 1945, BPUPKI menyetujui rancangan Undang-Undang Dasar dengan pembukaan diambil dari Piagam Jakarta. BPUPKI dibubarkan pada 7 Agustus 1945 setelah menyelesaikan tugasnya.
<b>Ringkasan Hasil Model T5 Setelah Dilakukan <i>Fine-tuning</i> Menggunakan Dataset Liputan6</b>	Pemerintah Jepang di Jawa dipimpin oleh Saiko Shikikan Kumakichi Harada membentuk Badan Penyelidik Usaha-Usaha Kemerdekaan Indonesia BPUPKI. Ia dibantu dua ketua muda yaitu seorang jepang, Shuchokan Shirebon dan R.P. Suroso. Dalam kesempatan itu Insinyur Soekarno tidak menjadi ketua karena ia ingin lebih aktif dalam berbagai diskusi. Tugas pokok BPUPKI adalah melakukan penyelidikan terhadap usaha-usaha persiapan kemerdekaan Indonesia. Persidangan ini dibicarakan masalah dasar negara. Beberapa orang tokoh yang berpidato untuk mengusulkan konsepsi yaitu Mr. Muhammad Yamin, Insinyur Soekarno, dan Prof. Mr. Supomo. Ketuhanan yang maha esa Kelima asas itu, atas petunjuk seorang ahli bahasa oleh Insinyur Soekarno diberi nama Pancasila, kemudian diusulkan dijadikan dasar negara Indonesia. Panitia 9 tersebut terus bekerja keras merumuskan rancangan pembukaan Undang-Undang Dasar yang nanti harus mengandung asas dan tujuan negara Indonesia Merdeka. Akhirnya tugas itu terselesaikan pada 22 Juni 1945. Insinyur Soekarno melaporkan hasil kerja panitia perancang Undang-Undang Dasar berhasil menyelesaikan tugasnya, maka pada 7 Agustus 1945 BPUPKI dibubarkan.

#### 4. KESIMPULAN

Penelitian ini mengevaluasi kinerja model *Whisper* dan T5 dalam tugas peringkasan teks dari transkrip video pembelajaran. Hasilnya menunjukkan bahwa kombinasi model *Whisper Turbo* dan T5 memberikan pendekatan unggul dengan akurasi dan relevansi terbaik. *Whisper Turbo* dipilih karena efisiensi dan akurasinya dalam transkripsi





teks, sementara T5 yang telah dilakukan *fine-tuning* dengan pembagian data pelatihan 70% dan pengujian 30% memberikan performa optimal dalam peringkasan. Untuk menangani batasan panjang input model T5, transkripsi teks dibagi menjadi segmen-segmen maksimal 512 token. Evaluasi menggunakan metrik ROUGE mengonfirmasi bahwa kombinasi ini menghasilkan *F1-Score* tertinggi pada ROUGE-1 (39,23), ROUGE-2 (13,17), dan ROUGE-L (23,84). Pendekatan ini menawarkan solusi praktis dan efisien untuk menganalisis konten pembelajaran berbasis video, menjadikannya relevan untuk mendukung inovasi di bidang pendidikan. Meskipun menunjukkan hasil yang unggul, penelitian ini menyoroti batasan panjang input teks pada model T5 sebagai tantangan penting. Pemotongan teks menjadi segmen dapat memengaruhi kualitas ringkasan, terutama untuk teks yang sangat panjang. Penelitian mendatang dapat mengeksplorasi model dengan kapasitas token yang lebih besar untuk menghasilkan ringkasan yang lebih komprehensif. Selain itu, perbandingan kinerja model T5 dengan model lain yang mampu memproses input teks panjang akan memberikan wawasan lebih dalam mengenai faktor-faktor yang memengaruhi kualitas peringkasan. Pendekatan ini berpotensi memperkaya analisis konten video pembelajaran, memperluas aplikasinya untuk mendukung pengajaran, *e-learning*, dan pengelolaan sumber belajar secara lebih efektif.

## REFERENCES

- [1] H. M. E. Putry, V. Nuzulul'Adila, R. Sholeha, and D. Hilmi, "Video based learning sebagai tren media pembelajaran di era 4.0," *Tarbiyatuna J. Pendidik. Ilm.*, vol. 5, no. 1, pp. 1–24, 2020, doi: 10.55187/tarjpi.v5i1.3870.
- [2] B. Rahmat and D. Darmiati, "Pengembangan Media Pembelajaran dengan Video Based Learning di Akademi Kebidanan Pelamonia," *Lect. J. Pendidik.*, vol. 12, no. 2, pp. 149–165, 2021, doi: 10.31849/lectura.v12i2.7268.
- [3] D. Wong, "Effectiveness of learning through video clips and video learning improvements between business related postgraduate and undergraduate students," *Int. J. Mod. Educ.*, vol. 2, no. 7, pp. 119–127, 2020, doi: 10.35631/IJMOE.27009.
- [4] H. B. U. Haq, M. Asif, and M. Bin Ahmad, "Video summarization techniques: a review," *Int. J. Sci. Technol. Res.*, vol. 9, no. 11, pp. 146–153, 2020.
- [5] A. Bahari and K. E. Dewi, "Peringkasan Teks Otomatis Abstraktif Menggunakan Transformer Pada Teks Bahasa Indonesia," *Komputa J. Ilm. Komput. dan Inform.*, vol. 13, no. 1, pp. 83–91, 2024, doi: 10.34010/komputa.v13i1.11197.
- [6] R. F. Khoiroh, E. Julianto, S. A. Ardiyansa, H. A. Fajri, A. A. R. Yasa, and B. Sangapta, "Implementasi Speech Recognition Whisper pada Debat Calon Wakil Presiden Republik Indonesia," *Explore*, vol. 14, no. 2, pp. 67–74, 2024, doi: 10.35200/ex.v14i2.115.
- [7] B. Ay, F. Ertam, G. Fidan, and G. Aydin, "Turkish abstractive text document summarization using text to text transfer transformer," *Alexandria Eng. J.*, vol. 68, pp. 1–13, 2023, doi: 10.1016/j.aej.2023.01.008.
- [8] Y. Singh, R. Kumar, S. Kabdal, and P. Upadhyay, "YouTube Video Summarizer using NLP: A Review," *Int. J. Performability Eng.*, vol. 19, no. 12, p. 817, 2023, doi: 10.23940/ijpe.23.12.p6.817823.
- [9] L. R. S. Gris, R. Marcacini, A. C. Junior, E. Casanova, A. Soares, and S. M. Aluisio, "Evaluating OpenAI's Whisper ASR for Punctuation Prediction and Topic Modeling of life histories of the Museum of the Person," *arXiv Prepr. arXiv2305.14580*, 2023.
- [10] I. N. Purnama and N. N. W. Utami, "Implementasi Peringkat Dokumen Berbahasa Indonesia Menggunakan Metode Text To Text Transfer Transformer (T5)," *J. Teknol. Inf. dan Komput.*, vol. 9, no. 4, 2023.
- [11] G. Hartawan, D. S. Maylawati, and W. Uriawan, "Bidirectional and Auto-Regressive Transformer (BART) for Indonesian Abstractive Text Summarization," *J. Inform. Polinema*, vol. 10, no. 4, pp. 535–542, 2024, doi: 10.33795/jip.v10i4.5242.
- [12] D. Ferdiansyah and C. S. K. Aditya, "Implementasi Automatic Speech Recognition Bacaan Al-Qur'an Menggunakan Metode Wav2Vec 2.0 dan OpenAI-Whisper," *J. Tek. Elektro dan Komput. TRIAC*, vol. 11, no. 1, pp. 11–16, 2024, doi: 10.21107/triac.v11i1.24332.
- [13] C. Raffel *et al.*, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, no. 140, pp. 1–67, 2020.
- [14] A. G. Etemad, A. I. Abidi, and M. Chhabra, "Fine-tuned t5 for abstractive summarization," *Int. J. Performability Eng.*, vol. 17, no. 10, p. 900, 2021, doi: 10.23940/ijpe.21.10.p8.900906.
- [15] D. A. Fadhilillah and D. Ikasari, "Optimizing Online News Understanding: Abstractive Summarization Approach with T5 for Comprehend Content." 2023
- [16] E. Zolotareva, T. M. Tashu, and T. Horváth, "Abstractive Text Summarization using Transfer Learning,," in *ITAT*, 2020, pp. 75–80.
- [17] J. Gabín, M. E. Ares, and J. Parapar, "Enhancing Automatic Keyphrase Labelling with Text-to-Text Transfer Transformer (T5) Architecture: A Framework for Keyphrase Generation and Filtering," *arXiv Prepr. arXiv2409.16760*, 2024.
- [18] A. A. Magriyanti, "Analisis pengembangan algoritma porter stemming dalam bahasa indonesia," 2018, doi: 10.31227/osf.io/7ge4v.
- [19] M. Barbella and G. Tortora, "Rouge metric evaluation for text summarization techniques," *Available SSRN 4120317*, 2022, doi: 10.2139/ssrn.4120317.
- [20] Y. Yuliska and K. U. Syaliman, "Literatur Review Terhadap Metode, Aplikasi dan Dataset Peringkasan Dokumen Teks Otomatis untuk Teks Berbahasa Indonesia," *IT J. Res. Dev.*, vol. 5, no. 1, pp. 19–31, 2020, doi: 10.25299/itjrd.2020.vol5(1).4688.