

BAB 1 PENDAHULUAN

1.1 Latar Belakang

Dokumen seringkali digunakan para ahli dalam dua pengertian. Pertama, berarti sumber tertulis bagi informasi sejarah sebagai kebalikan dari kesaksian lisan, artefak, peninggalan-peninggalan terlukis, dan petilasan-petilasan arkeologis. Pengertian kedua diperuntukkan bagi surat-surat resmi dan surat-surat negara seperti surat perjanjian, undang-undang, hibah, konsesi, dan lainnya (Gottschalk, 1986). Tetapi dalam pembuatannya memerlukan ketelitian yang cukup baik agar dokumen yang dibuat tidak perlu dilakukan revisi. Merevisi dokumen merupakan pekerjaan yang memakan waktu karena pengerjaannya masih dilakukan secara manual oleh manusia.

Salah satu jenis dokumen yang cukup familiar di lingkungan akademis adalah dokumen jurnal. Menurut Lasa Hs (2009), “Jurnal ilmiah adalah publikasi ilmiah yang berisi tentang kegiatan terkait ilmu pengetahuan seperti pengamatan empiris, pengembangan gagasan, sampai kumpulan dari pengetahuan baru.” Dan dalam perilisannya, jurnal ilmiah memiliki beberapa tahap yang cukup panjang dan memakan waktu, salah satu tahapannya adalah “*journal review*” *journal review* membutuhkan waktu yang cukup lama karena pengerjaannya masih dilakukan secara manual oleh pihak tim perilisan jurnal, sehingga dapat disimpulkan bahwa lamanya proses *journal review* disebabkan oleh banyaknya revisi pada sebuah jurnal.

Belum lama ini, fenomena ChatGPT sedang meningkat popularitasnya. Hal ini disebabkan oleh kemudahan akses dan kecanggihan yang dimiliki oleh ChatGPT, yang mampu membantu menyelesaikan banyak masalah. Teknologi dibalik ChatGPT adalah sebuah *language model* yang berukuran sangat besar untuk saat ini, yaitu *Generative Pre-Trained Transformers* (GPT) yang pertama kali dikenalkan pada 2018 oleh OpenAI pada sebuah artikel berjudul “*Improving Language Understanding by Generative Pre-Training*” GPT memiliki beberapa iterasi dan saat ini versi GPT sudah mencapai versi ke 4, GPT memiliki kemampuan dalam memahami konteks pada sebuah teks secara mendalam sehingga model tersebut dapat membantu menyelesaikan banyak permasalahan terkait dengan pemodelan bahasa.

Penelitian sebelumnya telah menunjukkan bahwa penggunaan teknologi NLP (*Natural Language Processing*) untuk memeriksa dan memperbaiki kesalahan penulisan sudah dilakukan. Sebagai contoh, penelitian yang dilakukan oleh (Musyafa, Gao, Solyman, Wu, & Khan, 2022). Dalam penelitiannya mengusulkan *Grammatical error correction* (GEC) otomatis berbasis arsitektur *Transformer* untuk memperbaiki kesalahan tata bahasa Indonesia dan juga dapat digunakan untuk memperbaiki teks *low resources languages* dan mencapai hasil yang signifikan mengungguli model GEC sebelumnya dengan skor F1 rata-rata 0,7194 dan skor BLEU 78,13, yang relatif efektif dalam tugas koreksi kesalahan tata bahasa.

1.2 Rumusan Masalah

Berdasarkan latar belakang yang telah dituliskan, pada penelitian ini mengangkat beberapa rumusan masalah, diantaranya:

- Apakah bisa menerapkan model GPT-2 untuk tugas koreksi kesalahan penulisan dokumen jurnal.
- Apa saja jenis kesalahan penulisan yang dapat dikoreksi menggunakan GPT-2.
- Seberapa efektif dan efisien penggunaan GPT-2 dalam mengoreksi kesalahan penulisan pada dokumen jurnal 6-8 halaman.

1.3 Tujuan Penelitian

Beberapa tujuan yang ingin dicapai setelah pengerjaan tugas akhir diantaranya:

- Menerapkan model GPT-2 untuk tugas koreksi penulisan dokumen jurnal.
- Mengidentifikasi terkait kesalahan penulisan apa saja yang dapat diperbaiki oleh GPT-2.
- Mengukur kinerja GPT-2 pada tugas koreksi penulisan dokumen jurnal.

1.4 Ruang Lingkup

Penelitian ini memiliki ruang lingkup sebagai berikut:

- Penerapan model akan dilakukan pada dokumen jurnal dengan 6-8 halaman dalam bahasa Indonesia.
- Bagian jurnal yang akan dikoreksi oleh model adalah hanya pada bagian teks.
- Teks yang akan dikoreksi dalam satu waktu oleh model tidak kurang dari 55 kata dan tidak lebih dari 512 token.
- Model akan menerima *input* dokumen jurnal dengan format PDF.
- *Dataset* yang dikumpulkan berasal dari jurnal yang ada di internet dengan tema atau topik pembahasan yang bebas (Tidak ditentukan).
- Varian jenis kesalahan penulisan yang ada pada *dataset* meliputi penggantian sinonim, pemberlakuan typo, pemberlakuan seluruh data menjadi *lowercase*, penerapan spasi setelah tanda baca yang tidak tepat, dan penggantian kata secara acak pada setiap paragraf atau kalimat pada *dataset input*.

1.5 State of The Art

Tabel 1. 1 State of The Art

Judul Jurnal	Pembahasan
<p data-bbox="300 1265 643 1301"><i>Attention Is All You Need</i></p> <p data-bbox="300 1420 419 1451">Peneliti:</p> <p data-bbox="300 1496 794 1697">Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Jakob Uszkoreit, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin.</p> <p data-bbox="300 1816 403 1848">Tahun:</p> <p data-bbox="300 1892 371 1924">2017</p>	<p data-bbox="813 1265 1050 1301"><u>Hasil Penelitian:</u></p> <p data-bbox="813 1346 1355 1765">Artikel ini memperkenalkan arsitektur <i>Transformer</i> model transduksi urutan pertama yang sepenuhnya didasarkan pada perhatian (<i>attention</i>), menggantikan <i>layers reccurent</i> yang paling umum digunakan dalam arsitektur pengode-penyandi (<i>encoder-decoder</i>) dengan <i>multi-headed self-attention</i>.</p> <p data-bbox="813 1883 1331 1919"><u>Alasan Menjadi Tinjauan Penelitian:</u></p>

<p>Publisher: NIPS</p>	<p>Sebagai panduan dalam memahami arsitektur <i>Transformer</i>, sebab bagian arsitektur <i>decoder</i> digunakan pada GPT</p>
<p><i>Improving Language Understanding by Generative Pre-training</i></p> <p>Peneliti: Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever</p> <p>Tahun: 2018</p> <p>Publisher: OpenAI</p>	<p><u>Hasil Penelitian:</u></p> <p>Artikel ini membahas tentang bagaimana pemahaman bahasa alami yang kuat melalui tahap <i>generative pre-training</i> dan <i>discriminative fine tuning</i>. Hasil menunjukkan peningkatan kinerja yang signifikan memang memungkinkan, dan memberikan petunjuk tentang model (<i>Transformer</i>) dan <i>dataset</i> (teks dengan ketergantungan jarak jauh) mana yang paling cocok dengan pendekatan ini.</p> <p><u>Alasan Menjadi Tinjauan Penelitian:</u></p> <p>Sebagai panduan dasar untuk pemahaman terkait <i>GENERATIVE PRE-TRAINED TRANSFORMERS</i>.</p>
<p><i>Language Models are Unsupervised Multitask Learners</i></p> <p>Peneliti: Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever</p>	<p><u>Hasil Penelitian:</u></p> <p>Artikel ini menjelaskan ketika model bahasa besar dilatih pada <i>dataset</i> yang cukup besar dan beragam, model tersebut mampu berperforma dengan baik di banyak domain dan <i>dataset</i>. Keragaman tugas yang dapat dilakukan oleh model dalam pengaturan <i>zero-shot</i></p>

<p>Tahun: 2019</p> <p>Publisher: OpenAI</p>	<p>menunjukkan bahwa model berkapasitas tinggi yang dilatih untuk memaksimalkan peluang dari korpus teks yang cukup beragam mulai belajar bagaimana melakukan sejumlah tugas secara mengejutkan tanpa perlu supervisi eksplisit.</p> <p><u>Alasan Menjadi Tinjauan Penelitian:</u></p> <p>Penelitian ini membantu menjelaskan mengenai model yang mampu melakukan berbagai tugas secara <i>zero-shot</i>, yaitu tanpa pelatihan khusus pada tugas tertentu. Hal ini menunjukkan fleksibilitas dan adaptabilitas model dalam memanfaatkan pengetahuannya untuk berbagai konteks dan tugas.</p>
<p><i>The survey: Text generation models in deep learning</i></p> <p>Peneliti: Touseef Iqbal, Shaima Qureshi</p> <p>Tahun: 2020</p>	<p><u>Hasil Penelitian:</u></p> <p>Jurnal ini membahas secara singkat perkembangan yang terjadi dalam pemodelan <i>Deep Generative</i>. Jurnal ini menyimpulkan bahwa di dalam bidang data kontinu (gambar), dominasi dipegang oleh GAN sementara di dalam bidang data diskrit (teks), dominasi dipegang oleh Variational Auto-Encoders.</p>

<p><i>Publisher:</i></p> <p>King Saud University</p>	<p><u>Alasan Menjadi Tinjauan Penelitian:</u></p> <p>Membantu dalam memahami pengertian mengenai model <i>generative</i> teks.</p>
<p><i>Automatic Correction of Indonesian Grammatical Errors Based on Transformers</i></p> <p>Penulis:</p> <p>Ahmad Musyafa, Ying Gao, Aiman Solyman, Chaojie Wu, Siraj Khan</p> <p>Tahun:</p> <p>2022</p> <p><i>Publisher:</i></p> <p>MDPI</p>	<p><u>Hasil Penelitian:</u></p> <p>Artikel ini membahas model GEC berbasis <i>neural-based</i> pertama untuk bahasa Indonesia yang juga bisa memperbaiki teks bahasa-bahasa dengan sumber daya terbatas lainnya. Selain itu, Untuk mengatasi masalah kekurangan data dalam GEC bahasa Indonesia, artikel ini mengusulkan metode <i>confusion semi supervised</i> untuk menghasilkan data pelatihan paralel dari korpus monolingual di luar domain. Hasil eksperimen menunjukkan keunggulan dengan skor F1 0.7194 dan BLEU 78.13, efektif dalam koreksi kesalahan tata bahasa, tetapi masih belum menangani kesalahan semantik dan sintetik.</p> <p><u>Alasan Menjadi Tinjauan Penelitian:</u></p> <p>Artikel ini memberi pengetahuan mengenai tugas koreksi kesalahan tata bahasa Indonesia.</p>

1.6 Sistematika Penulisan

Secara garis besar penulisan laporan tugas akhir ini terbagi dalam beberapa bab yang tersusun antara lain sebagai berikut:

1. BAB I PENDAHULUAN

Bab ini berisi latar belakang, rumusan masalah, tujuan penelitian, ruang lingkup, *state of the art*, dan sistematika penulisan.

2. BAB II LANDASAN TEORI

Pada bab ini diuraikan tentang teori yang berhubungan dengan menunjang penulisan tugas akhir.

3. BAB III METODOLOGI PENELITIAN

Pada bab berisi penjelasan tentang metodologi penelitian yang digunakan, permasalahan yang ada serta pengerjaan yang dilakukan. Implementasi yang dilakukan menggunakan model GPT-2.

4. BAB IV HASIL DAN PEMBAHASAN

Pada bab ini dijelaskan tentang *software/hardware environment* dan hasil implementasi model GPT-2 *small* untuk tugas koreksi dokumen jurnal.

5. BAB V PENUTUP

Pada bab ini berisi tentang kesimpulan yang diperoleh dari hasil pengujian model, serta saran perbaikan dan pengembangan model ke depannya.