

KLASIFIKASI HALAMAN WEB BERBASIS *MACHINE LEARNING* UNTUK OPTIMASI SEO MENGGUNAKAN FITUR NUMERIK DAN SEMANTIK BERBASIS INDOBERT

SITI NURADILLA



**PROGRAM STUDI MAGISTER STATISTIKA DAN SAINS DATA
SEKOLAH SAINS DATA, MATEMATIKA, DAN INFORMATIKA
INSTITUT PERTANIAN BOGOR
BOGOR
2026**

@Hak cipta milik IPB University

IPB University



IPB University
Bogor Indonesia

- Hak Cipta Dilindungi Undang-undang
1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
 2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

Perpustakaan IPB University



@Hak cipta milik IPB University

Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

KLASIFIKASI HALAMAN WEB BERBASIS *MACHINE LEARNING* UNTUK OPTIMASI SEO MENGGUNAKAN FITUR NUMERIK DAN SEMANTIK BERBASIS INDOBERT

SITI NURADILLA

Tesis
sebagai salah satu syarat untuk memperoleh gelar
Magister pada
Program Studi Statistika dan Sains Data

**PROGRAM STUDI MAGISTER STATISTIKA DAN SAINS DATA
SEKOLAH SAINS DATA, MATEMATIKA, DAN INFORMATIKA
INSTITUT PERTANIAN BOGOR
BOGOR
2026**

@Hak cipta milik IPB University

IPB University





PERNYATAAN MENGENAI TESIS DAN SUMBER INFORMASI SERTA PELIMPAHAN HAK CIPTA

Dengan ini saya menyatakan bahwa tesis dengan judul “Klasifikasi Halaman Web Berbasis *Machine Learning* untuk Optimasi SEO Menggunakan Fitur Numerik dan Semantik Berbasis IndoBERT” adalah karya saya dengan arahan dari dosen pembimbing dan belum diajukan dalam bentuk apa pun kepada perguruan tinggi mana pun. Sumber informasi yang berasal atau dikutip dari karya yang diterbitkan maupun tidak diterbitkan dari penulis lain telah disebutkan dalam teks dan dicantumkan dalam Daftar Pustaka di bagian akhir tesis ini.

Dengan ini saya melimpahkan hak cipta dari karya tulis saya kepada Institut Pertanian Bogor.

Bogor, Mei 2026

Siti Nuradilla
M0501241010

Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumunkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

RINGKASAN

SITI NURADILLA. Klasifikasi Halaman Web Berbasis *Machine Learning* untuk Optimasi SEO Menggunakan Fitur Numerik dan Semantik Berbasis IndoBERT. Dibimbing oleh BUDI SUSETYO dan CICI SUHAENI.

Meningkatnya volume data teks telah memperkuat kebutuhan akan metode yang efektif untuk mengekstraksi informasi yang bermakna, khususnya pada *search engine optimization* (SEO). Pada proses optimasi SEO, ekstraksi makna semantik menjadi krusial karena relevansi halaman tidak hanya dipengaruhi oleh panjang karakter komponen *on-page*, namun juga koherensi antar komponennya, meliputi *title*, *meta description*, H1, dan *address*. Namun, evaluasi SEO masih dilakukan secara manual, sehingga kurang efisien dan rentan terhadap bias subjektif. Penelitian sebelumnya hanya berfokus pada indikator permukaan seperti kepadatan kata kunci dan bergantung pada dataset empiris, sehingga belum menjamin keandalan model pada kondisi data yang berbeda. Oleh karena itu, diperlukan pembangkitan data sintesis yang dapat merepresentasikan beragam skenario distribusi data, sehingga ketahanan dan konsistensi performa model dapat diuji secara lebih komprehensif.

Untuk mengatasi keterbatasan tersebut, penelitian ini bertujuan untuk mengevaluasi kemampuan model GPT yang di-*fine-tune* dalam membangkitkan data sintesis yang merepresentasikan karakteristik SEO *on-page*. Selanjutnya, penelitian ini juga menganalisis karakteristik model klasifikasi halaman web dengan memanfaatkan fitur numerik berupa panjang karakter komponen SEO *on-page*, serta fitur semantik berupa tingkat koherensi antar komponen yang diperoleh melalui proses *embedding* menggunakan IndoBERT. Koherensi semantik diukur menggunakan *cosine similarity* antar komponen SEO *on-page* untuk merepresentasikan keselarasan konteks antar komponen. Berdasarkan fitur tersebut, halaman web diklasifikasikan menjadi *SEO-friendly* dan *non-SEO-friendly* menggunakan Random Forest, XGBoost, LightGBM, dan TabNet. Penelitian menggunakan delapan dataset, terdiri atas satu dataset empiris berisi 10.791 halaman web dan tujuh dataset sintesis dengan variasi tingkat koherensi (20%–80%). Data empiris dikumpulkan menggunakan Screaming Frog SEO Spider untuk memperoleh komponen SEO *on-page* seperti *title*, *meta description*, *address*, dan *heading*, sedangkan Semrush digunakan untuk memperoleh data peringkat halaman pada SERP yang dimanfaatkan dalam proses pelabelan kelas *SEO-friendly* dan *non-SEO-friendly*. Proses penelitian meliputi prapemrosesan, representasi teks, penghitungan koherensi, serta pelatihan dan evaluasi model, sementara kualitas data sintesis dievaluasi menggunakan metrik *diversity*, *novelty*, dan *duplication*.

Hasil eksplorasi menunjukkan bahwa halaman *SEO-friendly* memiliki pola yang lebih konsisten dalam panjang teks dan koherensi semantik dibandingkan halaman *non-SEO-friendly*. Selain itu, model GPT mampu menghasilkan data sintesis dengan kualitas yang baik, ditunjukkan oleh nilai *diversity* yang tinggi (0,9–1), *novelty* pada rentang 0,7–0,85, serta tingkat duplikasi yang sangat rendah (<0,06%). Temuan ini menunjukkan bahwa data sintesis yang dihasilkan mampu merepresentasikan variasi kondisi data dan mendukung proses klasifikasi dengan lebih komprehensif.



Pada tahap pemodelan, dilakukan perbandingan antara Random Forest, XGBoost, LightGBM, dan TabNet pada data empiris dan data sintetis. Pada data sintetis, performa model menunjukkan pola yang berbeda pada setiap tingkat koherensi. XGBoost dan LightGBM cenderung lebih kompetitif pada koherensi rendah hingga menengah, sedangkan TabNet mulai menunjukkan performa yang lebih tinggi pada koherensi tinggi. Pada data empiris, Random Forest memperoleh performa yang tinggi dengan *balanced accuracy* sebesar 0,8677, diikuti oleh XGBoost (0,8654) dan LightGBM (0,8605). Temuan ini sejalan dengan kajian awal (*baseline*) yang menunjukkan bahwa Random Forest sangat efektif dalam menangkap pola dominan pada distribusi data nyata yang relatif stabil. Namun, Random Forest cenderung mengalami penurunan performa pada struktur semantik antar kelas yang kurang tegas. Sebaliknya, model boosting yaitu XGBoost dan LightGBM menunjukkan performa yang lebih konsisten di berbagai tingkat koherensi, dengan nilai *balanced accuracy* yang kompetitif serta variasi performa yang relatif kecil. Performa terbaik secara umum dicapai pada skenario koherensi menengah (40%), di mana perbedaan karakteristik antar kelas menjadi lebih jelas. Sementara itu, pada koherensi rendah (20%–30%) dan tinggi (70%–80%), performa model cenderung menurun akibat meningkatnya ambiguitas atau homogenitas distribusi data.

Hasil uji statistik menunjukkan bahwa performa klasifikasi berbeda pada setiap tingkat koherensi dan bergantung pada jenis model yang digunakan. Uji lanjut perbandingan nilai tengah berganda dengan penyesuaian Holm memperlihatkan bahwa keunggulan model bersifat kontekstual, di mana Random Forest unggul pada data empiris, XGBoost dan LightGBM menunjukkan performa yang kompetitif pada koherensi rendah hingga menengah, sedangkan TabNet unggul pada koherensi tinggi. Dalam konteks implementasi sistem evaluasi SEO otomatis pada lingkungan data yang dinamis, model berbasis *boosting*, khususnya XGBoost, menunjukkan performa yang relatif stabil baik pada data empiris maupun berbagai tingkat koherensi data sintetis. Temuan ini menunjukkan bahwa model berbasis *boosting* memiliki potensi yang baik untuk digunakan pada kondisi data SEO *on-page* yang bervariasi.

Kata kunci: *IndoBERT, koherensi semantik, machine learning, search engine optimization.*

SUMMARY

SITI NURADILLA. Machine Learning-Based Web Page Classification for SEO Optimization Using Numerical and Semantic Features with IndoBERT. Supervised by BUDI SUSETYO and CICI SUHAENI.

The increasing volume of textual data has strengthened the need for effective methods to extract meaningful information, particularly in the context of search engine optimization (SEO). In SEO optimization, semantic information extraction is crucial because page relevance is influenced not only by the character length of on-page components, but also by the coherence among components, including title, meta description, H1, and address. However, SEO evaluation is still commonly performed manually, making it less efficient and prone to subjective bias. Previous studies have primarily focused on surface-level indicators such as keyword density and relied heavily on empirical datasets, which limits the reliability of models under varying data conditions. Therefore, synthetic data generation is required to represent diverse data distribution scenarios, enabling a more comprehensive evaluation of model robustness and performance consistency.

To address these limitations, this study aims to evaluate the capability of a fine-tuned GPT model in generating synthetic data that represent the characteristics of SEO on-page components. Furthermore, this study analyzes the characteristics of web page classification models by utilizing numerical features in the form of character lengths of SEO on-page components, as well as semantic features represented by the coherence level among components obtained through an embedding process using IndoBERT. Semantic coherence is measured using cosine similarity among SEO on-page components to represent contextual alignment between components. Based on these features, web pages are classified into SEO-friendly and non-SEO-friendly categories using Random Forest, XGBoost, LightGBM, and TabNet.

This study utilizes eight datasets, consisting of one empirical dataset containing 10,791 web pages and seven synthetic datasets with varying coherence levels (20%–80%). The empirical data were collected using Screaming Frog SEO Spider to obtain SEO on-page components such as title, meta description, address, and heading, while Semrush was used to extract SERP ranking data utilized in the labeling process for SEO-friendly and non-SEO-friendly classes. The research process includes preprocessing, text representation, coherence calculation, as well as model training and evaluation, while the quality of synthetic data is evaluated using diversity, novelty, and duplication metrics.

Exploratory results indicate that SEO-friendly pages exhibit more consistent patterns in text length and semantic coherence compared to non-SEO-friendly pages. Additionally, the GPT model successfully generates high-quality synthetic data, as indicated by high diversity scores (0.9–1.0), novelty values ranging from 0.7 to 0.85, and extremely low duplication rates (<0.06%). These findings suggest that the generated synthetic data effectively represent diverse data conditions and support more comprehensive classification experiments.

In the modeling stage, Random Forest, XGBoost, LightGBM, and TabNet were compared using both empirical and synthetic datasets. On synthetic datasets, model performance patterns varied across different coherence levels. XGBoost and



LightGBM tended to perform more competitively at low to medium coherence levels, whereas TabNet achieved higher performance at high coherence levels. On the empirical dataset, Random Forest achieved a high balanced accuracy of 0.8677, followed by XGBoost (0.8654) and LightGBM (0.8605). These findings are consistent with the baseline analysis, which showed that Random Forest is highly effective in capturing dominant patterns within relatively stable real-world data distributions. However, Random Forest tended to experience performance degradation when semantic boundaries between classes became less distinct. In contrast, boosting-based models, namely XGBoost and LightGBM, demonstrated more consistent performance across various coherence levels, achieving competitive balanced accuracy values with relatively low performance variation. Overall, the best performance was generally observed at medium coherence levels (40%), where semantic relationships among classes became more distinguishable. Meanwhile, at low coherence levels (20%–30%) and high coherence levels (70%–80%), model performance tended to decline due to increasing ambiguity or homogeneity in data distributions.

Statistical analysis results showed that classification performance differed across coherence levels and depended on the type of model used. Holm-adjusted multiple comparison tests further revealed that model superiority was contextual, where Random Forest performed best on empirical data, XGBoost and LightGBM showed competitive performance at low to medium coherence levels, and TabNet achieved the best performance at high coherence levels. In the context of implementing automated SEO evaluation systems in dynamic data environments, boosting-based models, particularly XGBoost, demonstrated relatively stable performance across both empirical data and multiple synthetic coherence scenarios. These findings suggest that boosting-based models have strong potential for implementation in on-page SEO evaluation tasks involving diverse data characteristics.

Keywords: *IndoBERT, machine learning, search engine optimization, semantic coherence.*



Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

© Hak Cipta milik IPB, tahun 2026
Hak Cipta dilindungi Undang-Undang

Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan atau menyebutkan sumbernya. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik, atau tinjauan suatu masalah, dan pengutipan tersebut tidak merugikan kepentingan IPB.

Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apa pun tanpa izin IPB.



@Hak cipta milik IPB University

IPB University



IPB University
— Bogor Indonesia —

Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

KLASIFIKASI HALAMAN WEB BERBASIS *MACHINE LEARNING* UNTUK OPTIMASI SEO MENGGUNAKAN FITUR NUMERIK DAN SEMANTIK BERBASIS INDOBERT

SITI NURADILLA

Tesis
sebagai salah satu syarat untuk memperoleh gelar
Magister pada
Program Studi Statistika dan Sains Data

**PROGRAM STUDI MAGISTER STATISTIKA DAN SAINS DATA
SEKOLAH SAINS DATA, MATEMATIKA, DAN INFORMATIKA
INSTITUT PERTANIAN BOGOR
BOGOR
2026**



@Hak cipta milik IPB University

Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumunkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.



Judul Tesis : Klasifikasi Halaman Web Berbasis *Machine Learning* untuk Optimasi SEO Menggunakan Fitur Numerik dan Semantik Berbasis IndoBERT
Nama : Siti Nuradilla
NIM : M0501241010

Disetujui oleh

Pembimbing 1:
Dr. Ir. Budi Susetyo, M.S.

Pembimbing 2:
Cici Suhaeni, S.Si., M.Si., Ph.D.

Diketahui oleh

Ketua Program Studi:
Dr. Agus Mohamad Soleh, S.Si., M.T.
NIP. 19750315 199903 1 004

Dekan Sekolah Sains Data, Matematika, dan Informatika:
Prof. Dr. Ir. Agus Bueno, M.Si., M.Kom.
NIP. 19660702 199302 1 001

Tanggal Ujian:
7 Mei 2026

Tanggal Pengesahan:



PRAKATA

Puji syukur penulis panjatkan ke hadirat Allah Subhanahu wa Ta'ala atas limpahan rahmat, taufik, dan karunia-Nya sehingga penulis dapat menyelesaikan tesis ini dengan baik. Karya ini disusun sebagai bagian dari pemenuhan tugas akhir dengan judul: “Klasifikasi Halaman Web Berbasis *Machine Learning* untuk Optimasi SEO Menggunakan Fitur Numerik dan Semantik Berbasis IndoBERT”.

Pada proses penyusunan tesis ini, sangat banyak bantuan dari berbagai pihak. Penulis menyampaikan terima kasih dan penghargaan setinggi-tingginya kepada:

Bapak Dr. Ir. Budi Susetyo, M.S. dan Ibu Cici Suhaeni, S.Si., M.Si., Ph.D selaku dosen pembimbing yang dengan penuh dedikasi telah memberikan arahan, bimbingan, motivasi, serta berbagai masukan yang membangun selama proses penelitian hingga penyusunan karya ilmiah ini.

Bapak Prof. Dr. Ir. Hari Wijayanto, M.Si. selaku penguji luar komisi pembimbing yang telah memberikan saran, masukan, dan wawasan yang sangat berharga dalam menyempurnakan karya ilmiah ini.

3. Bapak Dr. Agus Mohamad Soleh, S.Si., M.T. selaku pimpinan sidang tesis sekaligus Ketua Program Studi Statistika dan Sains Data atas arahan, saran, dan masukan yang bermanfaat bagi penyempurnaan karya ilmiah in

4. Seluruh dosen Program Studi Statistika dan Sains Data IPB University yang telah memberikan ilmu pengetahuan, pengalaman, serta pembelajaran yang bermanfaat selama penulis menempuh pendidikan.

5. Seluruh tenaga kependidikan Program Studi Statistika dan Sains Data IPB University yang telah membantu penulis dalam proses administrasi dan pelayanan akademik.

6. Kedua orang tua serta keluarga penulis yang senantiasa memberikan doa, kasih sayang, dukungan moral maupun material, serta semangat selama proses akademik dan penyusunan karya ilmiah ini.

7. Teman dekat penulis yang selalu hadir memberikan dukungan, bantuan, motivasi, serta menemani perjalanan penulis sejak awal perkuliahan hingga terselesaikannya karya ilmiah ini.

8. Penulis juga mengucapkan terima kasih kepada pihak Beasiswa Pendidikan Indonesia yang telah memfasilitasi dalam pelaksanaan Pendidikan dan penelitian.

Penulis menyadari bahwa karya ilmiah ini masih memiliki keterbatasan. Oleh karena itu, penulis membuka diri terhadap saran dan kritik yang membangun demi kesempurnaan karya ini di masa mendatang.

Akhir kata, semoga karya ini dapat memberikan manfaat bagi semua pihak yang membutuhkan dan berkontribusi nyata bagi pengembangan ilmu pengetahuan dan perumusan kebijakan berbasis data di Indonesia.

Bogor, Mei 2026

Siti Nuradilla

DAFTAR ISI

DAFTAR TABEL	x
DAFTAR GAMBAR	x
DAFTAR LAMPIRAN	xi
I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	4
1.3 Tujuan	4
1.4 Manfaat	4
1.5 Ruang Lingkup	5
II TINJAUAN PUSTAKA	6
2.1 <i>Search Engine Optimization</i> (SEO)	6
2.2 Representasi Semantik dan Pembangkitan Data Teks	7
2.3 Model Klasifikasi Berbasis <i>Machine Learning</i>	12
III METODE	21
3.1 Data	21
3.2 Prosedur Analisis Data	22
IV HASIL DAN PEMBAHASAN	29
4.1 Eksplorasi Karakteristik Data	29
4.2 Analisis Pembangkitan Data Sintetis	32
4.3 Analisis Pemodelan Klasifikasi pada Data Empiris dan Data Sintetis	38
V SIMPULAN DAN SARAN	50
5.1 Simpulan	50
5.2 Saran	50
DAFTAR PUSTAKA	51
LAMPIRAN	57
RIWAYAT HIDUP	67



DAFTAR TABEL

1	<i>Hyperparameter</i> utama XGBoost	16
2	<i>Confusion matrix</i>	20
3	Fitur-fitur data SEO	22
4	Fitur dengan penambahan koherensi antar komponen SEO <i>on-page</i>	22
5	Distribusi sampel data koherensi rendah dan koherensi tinggi per domain	25
6	Skema proporsi pelabelan fitur target pada data sintesis	27
7	Hasil <i>intra-duplication</i> data sintesis pada koherensi rendah	35
8	Hasil <i>inter-duplication</i> data sintesis pada koherensi rendah	36
9	Hasil <i>intra-duplication</i> data sintesis pada koherensi tinggi	36
10	Hasil <i>inter-duplication</i> data sintesis pada koherensi tinggi	36
11	Metrik evaluasi untuk setiap teknik penanganan ketidakseimbangan data	38
12	Hasil uji ART ANOVA	46
13	Hasil uji perbandingan nilai tengah berganda dengan penyesuaian Holm	47

DAFTAR GAMBAR

1	Arsitektur IndoBERT	8
2	Arsitektur GPT	10
3	Arsitektur TabNet	18
4	Prosedur analisis data secara menyeluruh	23
5	Distribusi panjang teks berdasarkan status SEO	29
6	Distribusi skor similaritas berdasarkan status SEO pada data empiris	31
7	Perbandingan skor similaritas pada data sintesis hasil pembangkitan	32
8	Distribusi <i>diversity</i> pada (a) koherensi tinggi dan (b) koherensi rendah	33
9	Boxplot <i>diversity</i> berdasarkan komponen SEO <i>on-page</i> pada (a) data koherensi tinggi dan (b) data koherensi rendah	33
10	Distribusi <i>novelty</i> pada (a) koherensi tinggi dan (b) koherensi rendah	34
11	Boxplot <i>novelty</i> berdasarkan komponen SEO <i>on-page</i> pada (a) data koherensi tinggi dan (b) data koherensi rendah	34
12	Perbandingan persentase konten SEO- <i>friendly</i> dan non-SEO- <i>friendly</i>	39
13	Perbandingan <i>balanced accuracy</i> dan standar deviasi model pada berbagai dataset	40
14	Perbandingan <i>precision</i> dan standar deviasi model pada berbagai dataset	43
15	Perbandingan <i>recall</i> dan standar deviasi model pada berbagai dataset	44
16	Perbandingan <i>F1-Score</i> dan standar deviasi model pada berbagai dataset	45
17	Plot interaksi antara tingkat koherensi dan model terhadap <i>balanced accuracy</i>	47

Hak Cipta Dilindungi Undang-undang
1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumunkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.

DAFTAR LAMPIRAN

1	Contoh dataset koherensi tinggi	58
2	Contoh dataset koherensi rendah	59
3	Contoh data <i>fine-tuning</i> berbasis instruksi dalam format JSONL untuk koherensi tinggi	60
4	Contoh data <i>fine-tuning</i> berbasis instruksi dalam format JSONL untuk koherensi rendah	61
5	<i>Template prompt</i> untuk pembangkitan data skema koherensi tinggi	62
6	<i>Template prompt</i> untuk pembangkitan data skema koherensi rendah	63
7	<i>Hyperparameter</i> terbaik dari setiap model dan skenario	64

Hak Cipta Dilindungi Undang-undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber :
 - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah
 - b. Pengutipan tidak merugikan kepentingan yang wajar IPB University.
2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin IPB University.