

**IMPLEMENTASI ALGORITMA BERT UNTUK KLASIFIKASI TOPIK
PENELITIAN PADA UNIVERSITAS MULTIMEDIA NUSANTARA**



SKRIPSI

Charlie Frederico
00000043442

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG
2024**

**IMPLEMENTASI ALGORITMA BERT UNTUK KLASIFIKASI TOPIK
PENELITIAN PADA UNIVERSITAS MULTIMEDIA NUSANTARA**



Diajukan sebagai salah satu syarat untuk memperoleh
Gelar Sarjana Komputer (S.Kom.)

Charlie Frederico

00000043442

UMMN

UNIVERSITAS

MULTIMEDIA

NUSANTARA

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA**

TANGERANG

2024

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : Charlie Frederico
Nomor Induk Mahasiswa : 00000043442
Program Studi : Informatika

Skripsi dengan judul:

Implementasi Algoritma BERT untuk Klasifikasi Topik Penelitian di Universitas Multimedia Nusantara

merupakan hasil karya saya sendiri bukan plagiat dari karya ilmiah yang ditulis oleh orang lain, dan semua sumber baik yang dikutip maupun dirujuk telah saya nyatakan dengan benar serta dicantumkan di Daftar Pustaka.

Jika di kemudian hari terbukti ditemukan kecurangan/ penyimpangan, baik dalam pelaksanaan Skripsi maupun dalam penulisan laporan Skripsi, saya bersedia menerima konsekuensi dinyatakan TIDAK LULUS untuk Tugas akhir yang telah saya tempuh.

UMM
UNIVERSITAS
MULTIMEDIA
NUSANTARA

Tangerang, 22 Mei 2024



(Charlie Frederico)




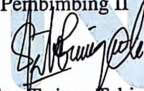
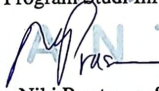
HALAMAN PENGESAHAN

Skripsi dengan judul
**IMPLEMENTASI ALGORITMA BERT UNTUK KLASIFIKASI TOPIK
PENELITIAN PADA UNIVERSITAS MULTIMEDIA NUSANTARA**

oleh
Nama : Charlie Frederico
NIM : 00000043442
Program Studi : Informatika
Fakultas : Fakultas Teknik dan Informatika

Telah diujikan pada hari Selasa, 04 Juni 2024
Pukul 10.00.s/s 12.00 dan dinyatakan
LULUS

Dengan susunan penguji sebagai berikut

Ketua Sidang	Penguji
	
(Sy Yuliani Yakub, S.Kom., M.T. PhD)	(Aditiyawan, S.Komp., M.Si)
NIDN: 0411037904	NIDN: 8994550022
Pembimbing I	Pembimbing II
	
(Eunike Endariahna Surbakti, S.Kom., M.T.I)	(Fenina Adline Twince Tobing, S.Kom., M.Kom)
NIDN: 0322099401	NIDN: 0406058802
Ketua Program Studi Informatika,	
	
(Dr. Eng. Niki Prastomo, S.T., M.Sc.)	
NIDN: 0419128203	

iii

Implementasi Algoritma BERT..., Charlie Frederico, Universitas Multimedia Nusantara

N U S A N T A R A

**HALAMAN PERSETUJUAN PUBLIKASI KARYA ILMIAH UNTUK
KEPENTINGAN AKADEMIS**

Yang bertanda tangan di bawah ini:

Nama : Charlie Frederico
NIM : 00000043442
Program Studi : Informatika
Jenjang : S1
Jenis Karya : Skripsi

Menyatakan dengan sesungguhnya bahwa:

- Saya bersedia memberikan izin sepenuhnya kepada Universitas Multimedia Nusantara untuk mempublikasikan hasil karya ilmiah saya di repositori Knowledge Center, sehingga dapat diakses oleh Civitas Akademika/Publik. Saya menyatakan bahwa karya ilmiah yang saya buat tidak mengandung data yang bersifat konfidensial dan saya juga tidak akan mencabut kembali izin yang telah saya berikan dengan alasan apapun.
- Saya tidak bersedia karena dalam proses pengajuan untuk diterbitkan ke jurnal/konferensi nasional/internasional (dibuktikan dengan *letter of acceptance*)**.

Tangerang, 22 Mei 2024

Yang menyatakan

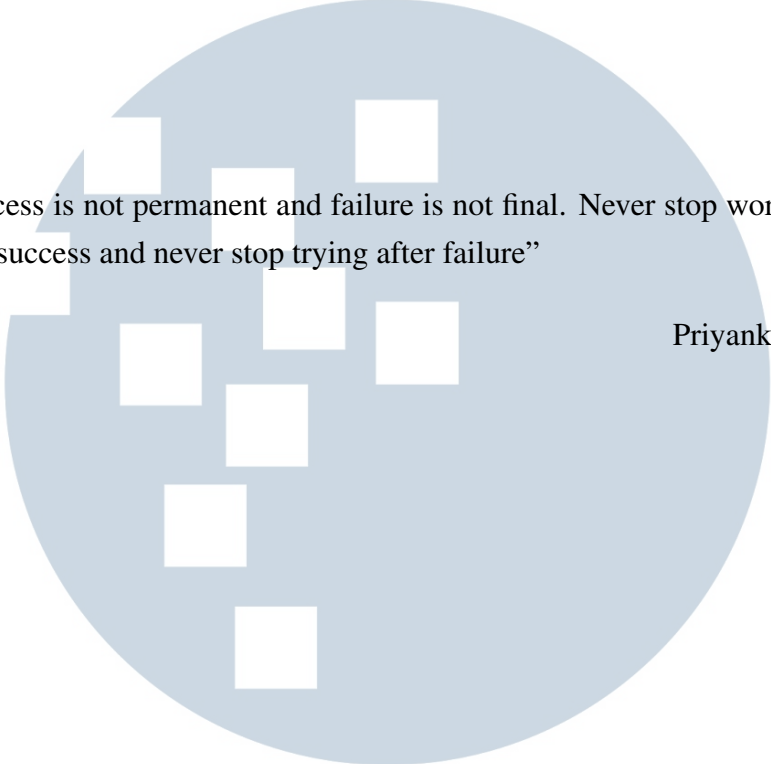


Charlie Frederico

UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

** Jika tidak bisa membuktikan LoA jurnal/HKI selama enam bulan ke depan, saya bersedia mengizinkan penuh karya ilmiah saya untuk diunggah ke KC UMN dan menjadi hak institusi UMN.

Halaman Persembahan / Motto



”Success is not permanent and failure is not final. Never stop working after success and never stop trying after failure”

Priyanka Ghediya

UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

KATA PENGANTAR

Puji Syukur atas berkat dan rahmat kepada Tuhan Yang Maha Esa, atas selesainya penulisan laporan Skripsi ini dengan judul: Implementasi Algoritma BERT untuk Klasifikasi Topik Penelitian di Universitas Multimedia Nusantara dilakukan untuk memenuhi salah satu syarat untuk mencapai gelar Sarjana Komputer Jurusan Informatika Pada Fakultas Teknik dan Informatika Universitas Multimedia Nusantara. Saya menyadari bahwa, tanpa bantuan dan bimbingan dari berbagai pihak, dari masa perkuliahan sampai pada penyusunan skripsi ini, sangatlah sulit bagi saya untuk menyelesaikan skripsi ini. Oleh karena itu, saya mengucapkan terima kasih kepada:

1. Bapak Dr. Ninok Leksono, selaku Rektor Universitas Multimedia Nusantara.
2. Bapak Dr. Eng. Niki Prastomo, S.T., M.Sc., selaku Dekan Fakultas Teknik dan Informatika Universitas Multimedia Nusantara.
3. Bapak Dr. Eng. Niki Prastomo, S.T., M.Sc., selaku Ketua Program Studi Informatika Universitas Multimedia Nusantara.
4. Ibu Eunike Endariahna Surbakti, S.Kom., M.T.I, sebagai Pembimbing pertama yang telah banyak meluangkan waktu untuk memberikan bimbingan, arahan dan motivasi atas terselesainya skripsi ini.
5. Ibu Fenina Adline Twince Tobing, S.Kom., M.Kom, sebagai Pembimbing kedua yang telah banyak membantu dan memberikan bimbingan atas terselesainya Skripsi ini.
6. Orang Tua yang telah memberikan bantuan dukungan material dan moral, sehingga penulis dapat menyelesaikan skripsi ini.

Semoga skripsi ini bermanfaat, baik sebagai sumber informasi maupun sumber inspirasi, bagi para pembaca.

Tangerang, 22 Mei 2024



Charlie Frederico

IMPLEMENTASI ALGORITMA BERT UNTUK KLASIFIKASI TOPIK PENELITIAN PADA UNIVERSITAS MULTIMEDIA NUSANTARA

Charlie Frederico

ABSTRAK

Terdapat lebih dari 3000 jurnal akademik dan artikel penelitian dari Universitas Multimedia Nusantara yang diterbitkan pada tahun 2018 hingga 2023. Untuk memudahkan pengkategorian, dibangun sistem klasifikasi teks yang menggunakan model *Bidirectional Encoder Representations from Transformers* (BERT) untuk mengkategorikan judul atau abstrak jurnal ke dalam 17 kategori UN SDG. UN SDG merupakan parameter utama yang digunakan untuk meningkatkan akreditasi fakultas dan program studi oleh Universitas Multimedia Nusantara. Terdapat 76.958 dataset berbahasa Inggris yang digunakan selama *training*. Model terbaik diperoleh dengan menggunakan *preprocessing* dengan *library NLTK* tanpa metode *sampling* yang menggunakan 70% *data training* dan 30% *data testing*. Parameter model mencakup 4 *epochs*, *learning rate* $2e-5$, dan *batch size* 32. Model ini mendapatkan nilai *precision* 0.99, *recall* 0.82, *f1-score* 0.87, dan *akurasi* 90.68%. Model ini dipilih untuk pengklasifikasian teks karena mendapatkan hasil terbaik dibandingkan model lainnya dan telah di demonstrasikan kepada pihak LPPM UMN.

Kata kunci: *Bidirectional Encoder Representations from Transformers*, *Deep Learning*, Jurnal Akademik, Klasifikasi Teks, *Natural Language Processing*



Implementation of BERT Algorithm for Classification of Research Topics at Multimedia Nusantara University

Charlie Frederico

ABSTRACT

There are more than 3,000 academic journals and research articles from Universitas Multimedia Nusantara, published between 2018 and 2023. To facilitate categorization, a text classification system was built using the Bidirectional Encoder Representations from Transformers (BERT) model to categorize the titles or abstracts of the journals into 17 UN SDG categories. The UN SDGs are crucial for enhancing the accreditation of faculties and study programs at the university. A dataset of 76,958 English-language records was used for training. The best model, obtained using NLTK preprocessing without sampling methods, utilized 70% training data and 30% testing data. Model parameters included 4 epochs, a learning rate of $2e-5$, and a batch size of 32. It achieved a precision of 0.99, recall of 0.82, f1-score of 0.87, and accuracy of 90.68%. This model is chosen for text classification because it outperforms other models and has been demonstrated to the LPPM UMN.

Keywords: Academic Journal, Bidirectional Encoder Representations from Transformers, Deep Learning, Natural Language Processing, Text Classification



DAFTAR ISI

HALAMAN JUDUL	i
PERNYATAAN TIDAK MELAKUKAN PLAGIAT	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERSETUJUAN PUBLIKASI ILMIAH	iv
HALAMAN PERSEMBAHAN/MOTO	v
KATA PENGANTAR	vi
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI	ix
DAFTAR GAMBAR	x
DAFTAR TABEL	xi
DAFTAR KODE	xii
DAFTAR LAMPIRAN	xiii
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	3
1.3 Batasan Permasalahan	3
1.4 Tujuan Penelitian	3
1.5 Manfaat Penelitian	4
1.6 Sistematika Penulisan	4
BAB 2 LANDASAN TEORI	6
2.1 Text Classification	6
2.2 Algoritma BERT	7
2.3 Natural Language Processing	8
2.4 Pemodelan Algoritma BERT	8
2.4.1 Pre-Training	9
2.4.2 Fine-Tuning	9
2.4.3 Input Representation	10
2.4.4 Model Architecture	10
2.4.5 Training	11
2.4.6 Evaluation	11
2.5 United Nations (UN) Sustainable Development Goals (SDG)	12
BAB 3 METODOLOGI PENELITIAN	13
3.1 Metodologi Penelitian	13
3.2 Perancangan Model	14
BAB 4 HASIL DAN DISKUSI	20
4.1 Spesifikasi Sistem	20
4.2 Implementasi Model	21
4.2.1 Load Dataset	21
4.2.2 Data Labelling	21
4.2.3 Preprocessing	22
4.2.4 Training Model	32
4.2.5 Evaluate Model	36
4.3 Uji Coba dan Evaluasi	41
4.3.1 Hasil Uji Coba	42
4.3.2 Evaluasi	53
4.3.3 Hasil Tanggapan LPPM UMN	54
BAB 5 SIMPULAN DAN SARAN	56
5.1 Simpulan	56
5.2 Saran	56
DAFTAR PUSTAKA	58

DAFTAR GAMBAR

Gambar 3.1	<i>Flowchart</i> secara keseluruhan	15
Gambar 3.2	<i>Flowchart</i> Modul Preprocessing	16
Gambar 3.3	<i>Flowchart</i> Modul Training Model	17
Gambar 3.4	<i>Flowchart</i> Modul Evaluate Model	19
Gambar 4.1	Tampilan Dataset	21
Gambar 4.2	Tampilan 17 kategori SDG beserta jumlah sampel	22
Gambar 4.3	Hasil <i>training model</i>	35
Gambar 4.4	Hasil <i>evaluate model</i>	37
Gambar 4.5	Contoh hasil dari <i>confusion matrix</i>	39
Gambar 4.6	Grafik <i>Classification Report</i> Skenario 1	43
Gambar 4.7	Grafik Akurasi Skenario 1	44
Gambar 4.8	Grafik <i>Classification Report</i> Skenario 2	46
Gambar 4.9	Grafik Akurasi Skenario 2	47
Gambar 4.10	Grafik <i>Classification Report</i> Skenario 3	49
Gambar 4.11	Grafik Akurasi Skenario 3	50
Gambar 4.12	Grafik <i>Classification Report</i> Skenario 4	52
Gambar 4.13	Grafik Akurasi Skenario 4	53



DAFTAR TABEL

Tabel 1.1	Perbandingan Algoritma	5
Tabel 4.1	Contoh teks menggunakan proses <i>case folding</i> pada <i>bert-base-uncased</i>	24
Tabel 4.2	Contoh teks menggunakan tokenisasi	25
Tabel 4.3	Contoh teks menggunakan <i>remove stopwords</i>	29
Tabel 4.4	Contoh teks menggunakan <i>lemmatization</i>	30
Tabel 4.5	Contoh teks menggunakan seluruh langkah <i>preprocessing</i>	31
Tabel 4.6	Hasil dari Confusion Matrix	40
Tabel 4.7	Hasil <i>Classification Report</i> Skenario 1	43
Tabel 4.8	Hasil <i>Classification Report</i> Skenario 2	45
Tabel 4.9	Hasil <i>Classification Report</i> Skenario 3	48
Tabel 4.10	Hasil <i>Classification Report</i> Skenario 4	51
Tabel 4.11	Hasil uji coba model keempat pada setiap skenario	54



DAFTAR KODE

4.1	Potongan kode pembuatan <i>onehotencoder</i> dan <i>import dataset</i>	22
4.2	Potongan code proses <i>case folding</i> menggunakan <i>bert-base-uncased</i>	23
4.3	Potongan kode untuk <i>library nltk</i> dan penggunaan <i>remove stopwords</i>	28
4.4	Potongan kode untuk <i>library nltk</i> dan penggunaan <i>lemmatization</i>	29
4.5	Potongan kode pembuatan tokenisasi dan konversi teks	31
4.6	Potongan kode inialisasi model bert dan penyesuaian <i>device</i>	32
4.7	Potongan kode <i>function training model</i>	33
4.8	Potongan kode <i>function evaluate</i> dan <i>compile training model</i>	34
4.9	Potongan kode <i>classification report evaluate model</i>	36
4.10	Potongan kode menghitung akurasi dan <i>confusion matrix</i>	38
4.11	Potongan kode penggunaan metode <i>undersampling</i>	47
4.12	Potongan kode penggunaan metode <i>oversampling</i>	50



DAFTAR LAMPIRAN

Lampiran 1	Form Bimbingan	61
Lampiran 2	Transkrip Wawancara 1	63
Lampiran 3	Transkrip Wawancara 2	65
Lampiran 4	Validasi Artikel Oleh LPPM UMN	66
Lampiran 5	Hasil Turnitin	69

