

**IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE
SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS**



SKRIPSI

**George Marcellino Jo
00000045841**

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

**IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE
SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS**



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

George Marcellino Jo

00000045841

UMMN

UNIVERSITAS

MULTIMEDIA

NUSANTARA

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

TANGERANG

2024

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : George Marcellino Jo

NIM : 00000045841

Program Studi : Informatika

Menyatakan dengan sesungguhnya bahwa Skripsi saya yang berjudul:
**IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE
SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS**

merupakan hasil karya saya sendiri, bukan merupakan hasil plagiat, dan tidak pula dituliskan oleh orang lain; Semua sumber, baik yang dikutip maupun dirujuk, telah saya cantumkan dan nyatakan dengan benar pada bagian Daftar Pustaka.

Jika di kemudian hari terbukti ditemukan kecurangan/penyimpangan, baik dalam pelaksanaan skripsi maupun dalam penulisan laporan karya ilmiah, saya bersedia menerima konsekuensi untuk dinyatakan **TIDAK LULUS**. Saya juga bersedia menanggung segala konsekuensi hukum yang berkaitan dengan tindak plagiarisme ini sebagai kesalahan saya pribadi dan bukan tanggung jawab Universitas Multimedia Nusantara.

Tangerang, 18 Juni 2024



(George Marcellino Jo)

HALAMAN PENGESAHAN

Skripsi dengan judul

IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS

oleh

Nama : George Marcellino Jo
NIM : 00000045841
Program Studi : Informatika
Fakultas : Fakultas Teknik dan Informatika

Telah diujikan pada hari Selasa, 4 Juni 2024

Pukul 13.00 s/s 15.00 dan dinyatakan

LULUS

Dengan susunan penguji sebagai berikut

Ketua Sidang



(Eunike Endariahna Surbakti, S.Kom.,
M.T.I.)

NIDN: 0322099401

Penguji



(Alexander Waworuntu, S.Kom., M.T.I.)

NIDN: 0309068503

Pembimbing



14 Juni 2024
(Arya Wicaksana, S.Kom., M.Eng.Sc. (OCA, CEH, CEI)

NIDN: 0315109103

Pjs. Ketua Program Studi Informatika,



(Dr. Eng. Niki Prastomo, S.T., M.Sc.)

NIDN: 0419128203

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

Yang bertanda tangan di bawah ini:

Nama : George Marcellino Jo

NIM : 00000045841

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, 18 Juni 2024

Yang menyatakan



George Marcellino Jo

U M M N
U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A

** 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

"The more you fear something to happen, the more likely it is to occur"

No.1 Murphy Law



KATA PENGANTAR

Puji Syukur atas berkat dan rahmat kepada Tuhan Yang Maha Esa, atas selesainya penulisan laporan Skripsi ini dengan judul: IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS 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: Mengucapkan terima kasih

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 Pjs. Ketua Program Studi Informatika Universitas Multimedia Nusantara.
4. Bapak Arya Wicaksana, S.Kom., M.Eng.Sc. (OCA, CEH, CEI, sebagai Pembimbing pertama yang telah banyak meluangkan waktu untuk memberikan bimbingan, arahan dan motivasi atas terselesainya tesis ini.
5. Orang Tua, dan keluarga saya yang telah memberikan bantuan dukungan material dan moral, sehingga penulis dapat menyelesaikan tesis ini.

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

Tangerang, 18 Juni 2024



George Marcellino Jo

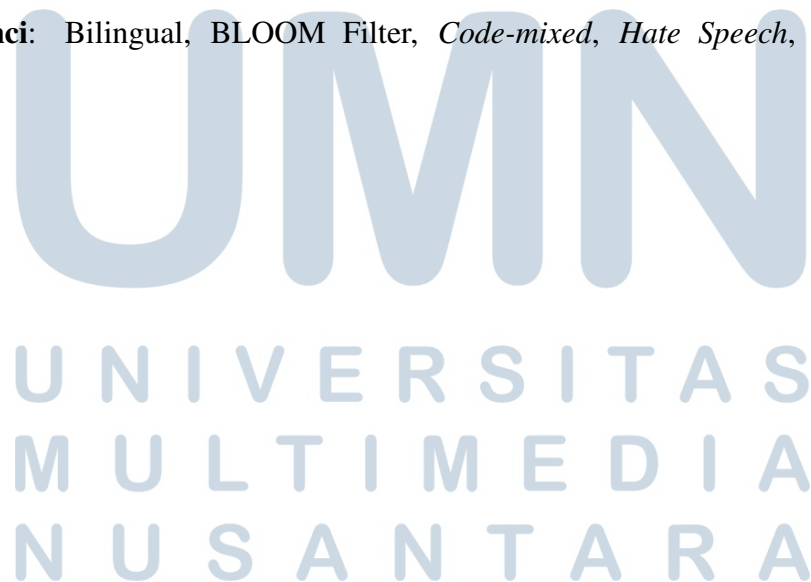
IMPLEMENTASI BLOOM FILTER UNTUK PENDETEKSI HATE SPEECH BILINGUAL BAHASA INDONESIA DAN INGGRIS

George Marcellino Jo

ABSTRAK

Hate speech adalah ungkapan atau pidato yang menghina individu atau kelompok berdasarkan (dugaan) keanggotaan dalam suatu kelompok sosial yang dikenali melalui ciri-ciri seperti ras, etnis, jenis kelamin, orientasi seksual, agama, usia, disabilitas fisik atau mental, dan faktor-faktor lainnya. Dikarenakan pertumbuhan pesat internet yang memungkinkan untuk terjadinya komunikasi dari berbagai belahan dunia dengan mudah tanpa ada batasan jarak, interferensi bahasa lain ke dalam bahasa Indonesia pun tidak dapat dihindari, dan membuat terciptanya penggunaan bahasa bilingual di kehidupan sehari-hari. BLOOM-560M akan digunakan untuk mendeteksi apakah sebuah teks bilingual Indonesia-Inggris merupakan teks *hate speech* atau *non-hate speech*. Hasil dari penelitian mendeteksi *hate speech* atau *non-hate speech* pada teks bilingual dengan menggunakan BLOOM Filter pada jumlah *epoch* 6 dengan menggunakan nilai *learning rate* sebesar $2e-5$ dan nilai *epsilon* $1e-8$, didapatkan *accuracy* 94,28%, *precision* 88,63%, *recall* 92,85%, dan *F1-score* 90,69%. Model juga mendapatkan hasil yang cukup memuaskan dalam mendeteksi teks manual yang di-*forward pass* ke dalam model dengan baik, 5 dari 6 teks yang dimasukkan dapat diprediksi oleh model dengan benar.

Kata kunci: Bilingual, BLOOM Filter, *Code-mixed*, *Hate Speech*, *Machine Learning*



**BLOOM FILTER IMPLEMENTATION FOR DETECTING BILINGUAL
HATE SPEECH OF INDONESIA AND ENGLISH**

George Marcellino Jo

ABSTRACT

Hate speech is an expression or speech that insults individuals or groups based on (alleged) membership in a social group identified by characteristics such as race, ethnicity, gender, sexual orientation, religion, age, physical or mental disability, and other factors. Due to the rapid growth of the internet, which enables communication from various parts of the world easily without distance limitations, language interference into Indonesian cannot be avoided, leading to the use of bilingual language in daily life. BLOOM-560M will be used to detect whether a bilingual Indonesian-English text is hate speech or non-hate speech. The results of the research on detecting hate speech or non-hate speech in bilingual texts using the BLOOM Filter, with 6 epochs, a learning rate of $2e-5$, and an epsilon value of $1e-8$, showed an accuracy of 94.28%, precision of 88.63%, recall of 92.85%, and an F1-score of 90.69%. The model achieved satisfactory results and was able to effectively detect manually input texts during the forward pass. The model correctly predicted 5 out of the 6 texts that were inputted.

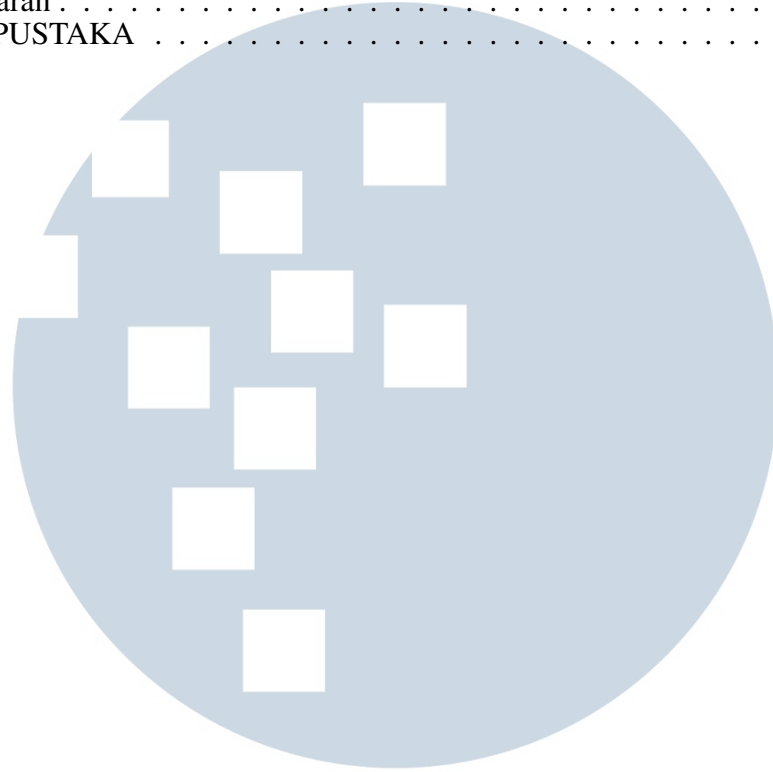
Keywords: *Bilingual, BLOOM Filter, Code-mixed, Hate Speech, Machine Learning*



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 TABEL	xi
DAFTAR GAMBAR	xii
DAFTAR KODE	xiii
DAFTAR LAMPIRAN	xiv
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	2
1.3 Batasan Permasalahan	3
1.4 Tujuan Penelitian	3
1.5 Manfaat Penelitian	3
1.6 Sistematika Penulisan	4
BAB 2 LANDASAN TEORI	5
2.1 Bilingual Hate Speech	5
2.2 BLOOM Filter	6
2.3 Metrik Evaluasi	8
BAB 3 METODOLOGI PENELITIAN	10
3.1 Pengumpulan Data	11
3.2 Data Pre-Process	11
3.3 Model Training	13
3.4 Model Testing	17
BAB 4 HASIL DAN DISKUSI	19
4.1 Spesifikasi Sistem	19
4.2 BLOOM Filter Model	19
4.2.1 Import library	19
4.2.2 Use Google Colab resource	20
4.2.3 Load Training Dataset	21
4.2.4 Initialize BLOOM Tokenizer	21
4.2.5 Padding Dataset	22
4.2.6 Masking Dataset	23
4.2.7 Split Training Data and Validation Data	23
4.2.8 Validation Data and Training Data Into Tensor Object	24
4.2.9 Initiate Optimizer and Scheduler	25
4.2.10 Train Model and Validate Model	25
4.2.11 Save Model	30
4.2.12 Evaluate Model	31
4.3 Pengujian dan Evaluasi	35
4.3.1 Skenario pengujian dan evaluasi	35
4.3.2 Hasil dari pengujian dan evaluasi	36
4.3.3 Evaluasi hasil pengujian	37

BAB 5	SIMPULAN DAN SARAN	40
5.1	Simpulan	40
5.2	Saran	40
DAFTAR PUSTAKA		41



UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

DAFTAR TABEL

Tabel 4.1	Tabel hasil pelatihan model BLOOM yang sudah di- <i>fine tuning</i> dengan data pelatihan	28
Tabel 4.2	Tabel hasil pengujian model BLOOM yang sudah di- <i>fine tuned</i> dengan data bilingual	36
Tabel 4.3	Tabel hasil pengujian model BLOOM yang sudah di- <i>fine tuned</i> berdasarkan skenario.	36



DAFTAR GAMBAR

Gambar 2.1	Inisialisasi <i>array</i>	6
Gambar 2.2	BLOOM menyimpan lebih dari satu indeks	7
Gambar 3.1	<i>Workflow</i> Penelitian	10
Gambar 3.2	<i>Tokenized text</i>	11
Gambar 3.3	<i>Padding and Masking</i>	12
Gambar 3.4	<i>Data Pre-process</i>	13
Gambar 3.5	<i>DataLoader output</i>	14
Gambar 3.6	<i>Train model</i>	15
Gambar 3.7	<i>Fine-tuning process</i>	16
Gambar 3.8	<i>Validate trained model</i>	17
Gambar 3.9	<i>Evaluate model</i>	18
Gambar 4.1	Tampilan <i>training loss graph</i>	29
Gambar 4.2	Hasil dari <i>confusion matrix</i> dalam bentuk plot	38



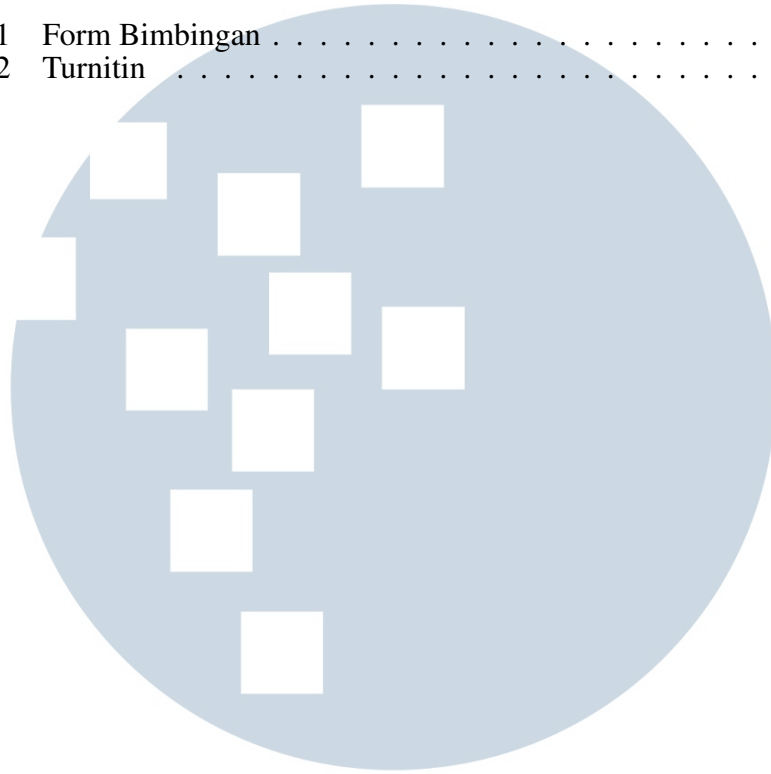
DAFTAR KODE

Kode 4.1	Potongan kode <i>import library</i>	19
Kode 4.2	Potongan kode menggunakan GPU	20
Kode 4.3	Potongan kode memasukkan <i>training data</i> ke model	21
Kode 4.4	Potongan kode pemanggilan BLOOM <i>tokenizer</i>	21
Kode 4.5	Potongan kode penggunaan BLOOM <i>tokenizer</i>	21
Kode 4.6	Potongan kode penggunaan BLOOM <i>tokenizer</i> untuk semua <i>dataset</i>	22
Kode 4.7	Potongan kode proses <i>padding</i> untuk <i>dataset</i>	23
Kode 4.8	Potongan kode proses <i>masking</i> untuk <i>dataset</i>	23
Kode 4.9	Potongan kode proses pembagian data <i>training</i> dan <i>validation</i> . .	24
Kode 4.10	Potongan kode proses pengubahan data <i>training</i> dan <i>validation</i> menjadi objek tensor	24
Kode 4.11	Potongan kode inisialisasi <i>optimizer</i> dan <i>scheduler</i>	25
Kode 4.12	Potongan kode tahapan awal pelatihan model	26
Kode 4.13	Potongan kode proses inti dalam pelatihan model	27
Kode 4.14	Potongan kode <i>training loss graph</i>	29
Kode 4.15	Potongan kode penyimpanan model	30
Kode 4.16	Potongan kode memasukkan data <i>testing</i>	31
Kode 4.17	Potongan kode memasukkan model dan <i>tokenizer</i>	31
Kode 4.18	Potongan kode pembangunan proses <i>pre-process</i> untuk fungsi <i>testing</i>	31
Kode 4.19	Potongan kode pembangunan proses prediksi dalam fungsi <i>testing</i>	32
Kode 4.20	Potongan kode pembangunan proses evaluasi dalam fungsi <i>testing</i>	33
Kode 4.21	Potongan kode pembangunan proses menampilkan prediksi dan <i>confusion matrix</i> dalam fungsi <i>testing</i>	34



DAFTAR LAMPIRAN

Lampiran 1	Form Bimbingan	44
Lampiran 2	Turnitin	45



UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA