

**IDENTIFIKASI DUPLIKAT PERTANYAAN MENGGUNAKAN METODE  
*GLOBAL VECTOR (GLOVE)***

**SKRIPSI**



disusun oleh

**Yohanes Eudes Anjas Susetya**

**17.11.0922**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2020**

**IDENTIFIKASI DUPLIKAT PERTANYAAN MENGGUNAKAN METODE  
*GLOBAL VECTOR (GLOVE)***

**SKRIPSI**

untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Yohanes Eudes Anjas Susetya**

**17.11.0922**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2020**

**PERSETUJUAN**

**SKRIPSI**

**IDENTIFIKASI DUPLIKAT PERTANYAAN MENGGUNAKAN  
METODE *GLOBAL VECTOR* (GLOVE)**

yang dipersiapkan dan disusun oleh

**Yohanes Eudes Anjas Susetya**

**17.11.0922**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 23 November 2020

**Dosen Pembimbing,**

**Mardhiya Hayaty, S.T., M.Kom.**

**NIK. 190302108**

**PENGESAHAN**  
**SKRIPSI**  
**IDENTIFIKASI DUPLIKAT PERTANYAAN MENGGUNAKAN**  
**METODE *GLOBAL VECTOR* (GLOVE)**

yang dipersiapkan dan disusun oleh

**Yohanes Eudes Anjas Susetya**

**17.11.0922**

telah dipertahankan di depan Dewan Penguji  
pada tanggal 18 November 2020

**Susunan Dewan Penguji**

**Nama Penguji**

**Tanda Tangan**

**Barka Satya, M.Kom.**  
**NIK. 190302126**

**Nuraini, M.Kom.**  
**NIK. 190302066**

**Mardhiya Hayaty, S.T., M.Kom.**  
**NIK. 190302108**

Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 18 November 2020

**DEKAN FAKULTAS ILMU KOMPUTER**

**Krisnawati, S.Si, M.T.**  
**NIK. 190302038**

## PERNYATAAN

Saya yang bertandatangan di bawah ini menyatakan bahwa skripsi ini merupakan karya saya sendiri (ASLI) dan isi pada skripsi ini tidak terdapat karya yang pernah diajukan oleh orang lain untuk memperoleh gelar akademis di suatu institusi pendidikan tinggi manapun, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis dan/atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Segala sesuatu yang terkait dengan masalah dan karya yang telah dibuat adalah menjadi tanggungjawab saya pribadi.

Yogyakarta, 21 November 2020



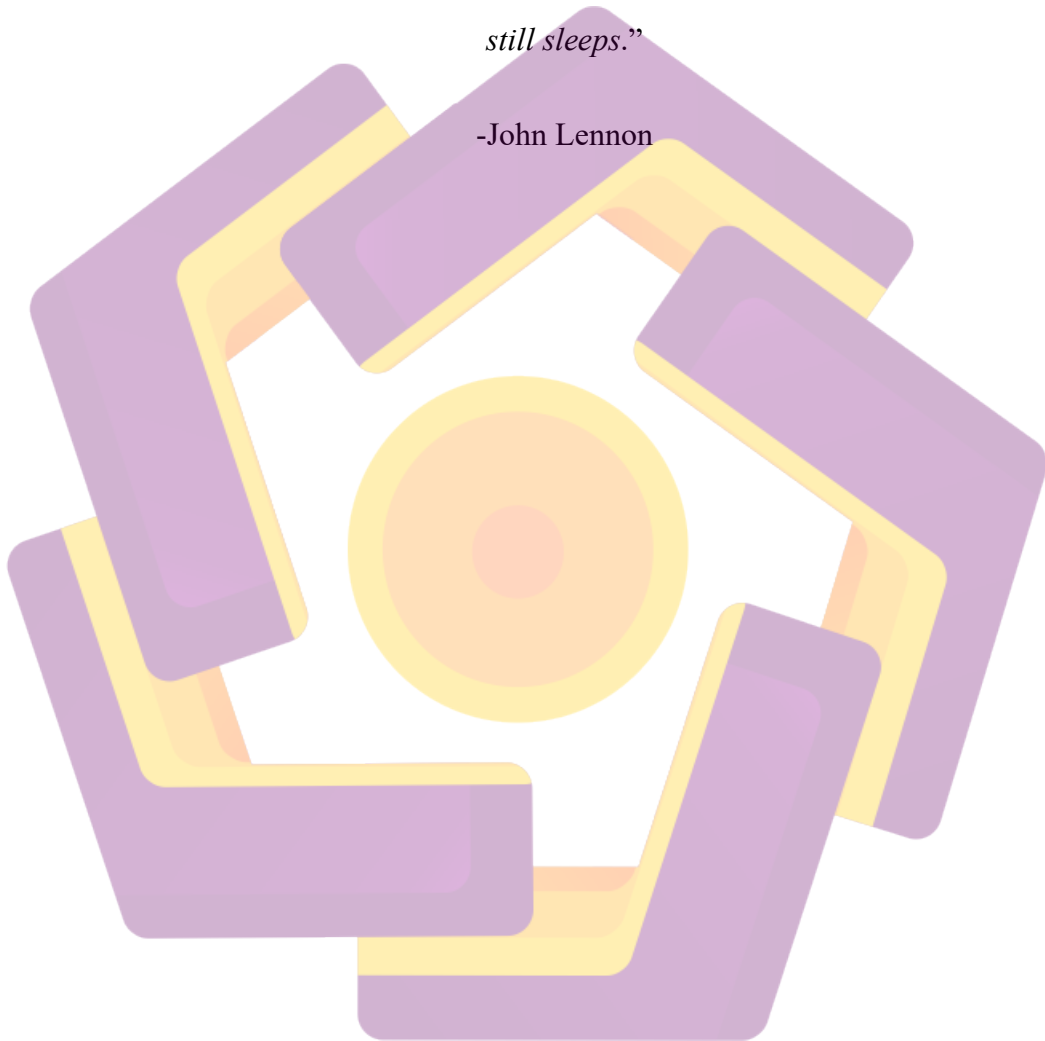
Yohanes Eudes Anjas Susetya

NIM 17.11.0922

**MOTTO**

*“When you do something noble and beautiful and nobody noticed, do not be sad.  
For the sun every morning is a beautiful spectacle and yet most of the audience  
still sleeps.”*

-John Lennon



## PERSEMBAHAN

Puji Tuhan berkat kerja keras serta doa, skripsi ini dapat diselesaikan dengan baik. Segala puji dan syukur kepada Tuhan Yang Maha Esa yang tiada henti memberikan berkat serta kekuatan untuk dapat menyelesaikan skripsi. Dengan ini saya mempersembahkan skripsi ini kepada semua pihak yang terlibat secara langsung atau tidak langsung, yaitu untuk :

1. Kedua orang tua saya (Theopilus Suryantoro Budi Susilo dan Theresia Wagiyem).
2. Adik saya (Wilfridus Bimo Ardiyantoro) dan bude (Yustina Widaryanti).
3. Dosen pembimbing saya Ibu Mardhiya Hayaty, S.T., M.Kom., yang telah sabar dalam membimbing saya dari awal sampai skripsi ini selesai.
4. Dosen Universitas AMIKOM Yogyakarta yang telah memberikan banyak ilmu selama perkuliahan berlangsung.
5. Teman sepenelitian (Dhimas, Lutfi, Andri, Bangdim) yang selalu membantu, menemani dan memberikan dukungan kepada saya untuk dapat menyelesaikan skripsi ini.
6. Teman sekelas 17-IF-01 yang telah memberikan dukungan selama perkuliahan.
7. Teman sorak-sorai yang senantiasa memberikan semangat (Devita, Gilang, Dendy, Bagas, Lordjek, Pandu, Aldo, Bapok, Roy, Arfian dan Dito).
8. Tim valorant penghuni *iron* yang senantiasa memberikan semangat dan petunjuk dalam setiap *round game* (Bram, Emo, Parmin, Fariz).

## KATA PENGANTAR

Penulis panjatkan puji dan syukur kehadiran Tuhan Yang Maha Esa atas berkat dan rahmat-Nya sehingga penulis dapat menyelesaikan skripsi dengan baik yang berjudul **“IDENTIFIKASI DUPLIKAT PERTANYAAN MENGGUNAKAN METODE *GLOBAL VECTOR (GLOVE)*”** disusun sebagai salah satu syarat utama untuk menyelesaikan program sarjana pada Universitas AMIKOM Yogyakarta. Penyelesaian skripsi ini juga tidak lepas dari bantuan berbagai pihak, karena itu pada kesempatan ini penulis ingin menyampaikan rasa hormat dan terima kasih kepada :

1. Prof. Dr. M. Suyanto, MM. selaku Rektor Universitas AMIKOM Yogyakarta.
2. Ibu Krisnawati, S.Si, M.T. selaku Dekan Fakultas Ilmu Komputer Universitas AMIKOM Yogyakarta.
3. Bapak Sudarmawan, M.T. selaku Ketua Program Studi Informatika Universitas AMIKOM Yogyakarta.
4. Ibu Mardhiya Hayaty, S.T., M.Kom. selaku dosen pembimbing yang selalu bijaksana dalam memberikan bimbingan serta nasehat dan waktunya selama pengerjaan skripsi ini.
5. Bapak Barka Satya, M.Kom. dan Ibu Nuraini, M.Kom. selaku dosen penguji. Terimakasih atas segala saran yang diberikan selama pengujian untuk memperbaiki penelitian menjadi lebih baik lagi.

Penulis menyadari bahwa skripsi ini masih ada kekurangan. Maka, penulis menerima segala kritik dan saran yang membangun serta teguran dari semua pihak. Semoga skripsi ini bisa bermanfaat baik bagi penulis serta pembaca. Atas saran dan kritik penulis ucapkan terima kasih.

Yogyakarta, 23 November 2020

Yohanes Eudes Anjas Susetya

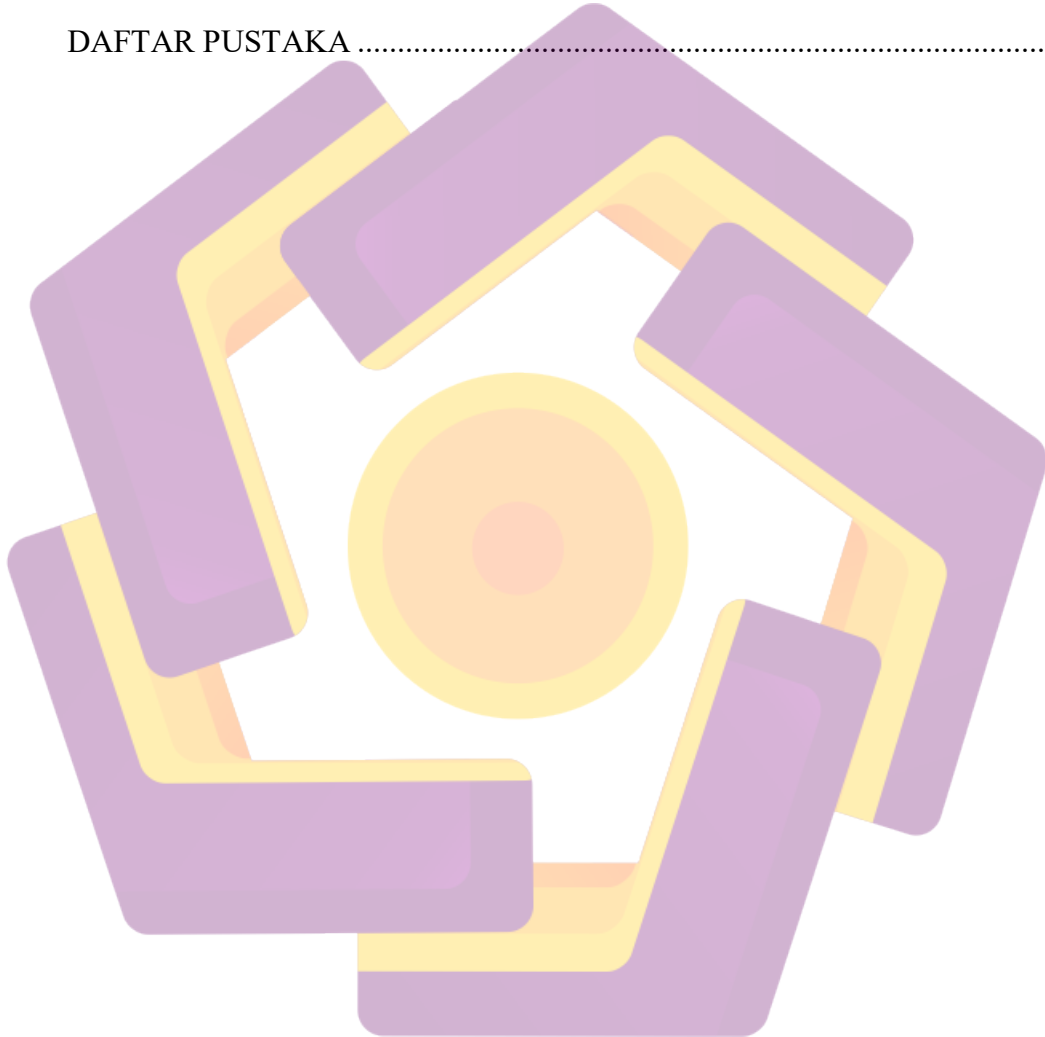


## DAFTAR ISI

JUDUL .....	i
LEMBAR PERSETUJUAN.....	ii
LEMBAR PENGESAHAN .....	iii
PERNYATAAN.....	iv
MOTTO.....	v
PERSEMBAHAN .....	vi
KATA PENGANTAR .....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	xi
DAFTAR GAMBAR .....	xii
INTISARI.....	xiii
<i>ABSTRACT</i> .....	xiv
<b>BAB I PENDAHULUAN</b> .....	<b>1</b>
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah .....	2
1.3 Batasan Masalah.....	2
1.4 Maksud dan Tujuan Penelitian.....	3
1.5 Manfaat Penelitian.....	3
1.6 Metode Penelitian.....	3
1.6.1 Pengumpulan Data.....	3
1.6.2 <i>Annotation</i> .....	4
1.6.3 <i>Preprocessing Data</i> .....	4
1.6.4 Ekstraksi Fitur.....	4
1.6.5 Implementasi Algoritma Klasifikasi.....	5
1.6.6 Evaluasi.....	5
1.7 Sistematika Penulisan.....	5
<b>BAB II LANDASAN TEORI</b> .....	<b>7</b>
2.1 Tinjauan Pustaka .....	7
2.2 Dasar Teori.....	10

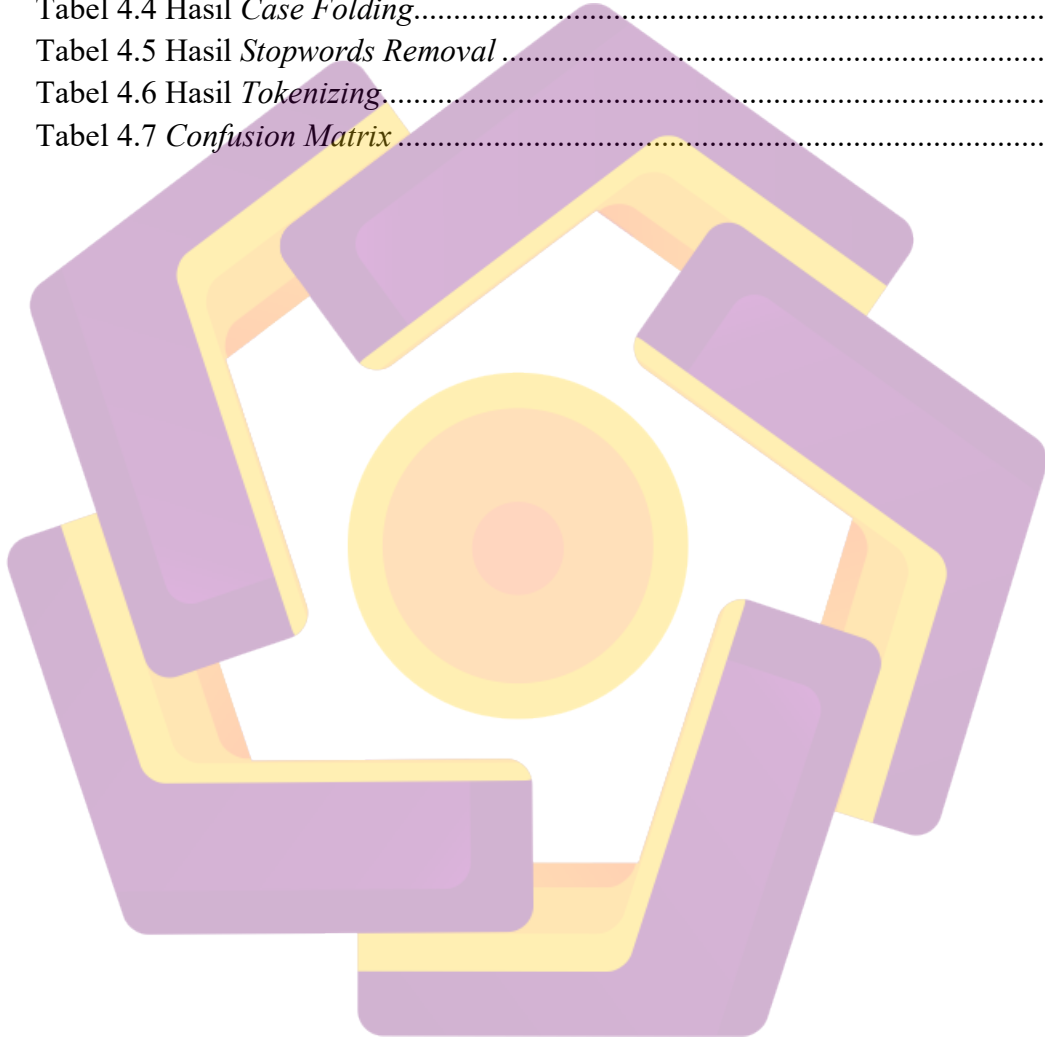
2.2.1	<i>Natural Language Processing</i> .....	10
2.2.2	<i>Deep Learning</i> .....	10
2.2.3	<i>Data Mining</i> .....	11
2.2.4	<i>Text Mining</i> .....	11
2.2.5	<i>Word Embedding</i> .....	13
2.2.6	<i>Indexing</i> .....	15
2.2.7	<i>Global Vector (GloVe)</i> .....	16
2.2.8	<i>Preprocessing Data</i> .....	17
2.2.9	<i>Optimasi Adam</i> .....	18
2.2.10	<i>Batch Size dan Epoch</i> .....	20
2.2.11	<i>Long Short Term Memory</i> .....	20
2.2.12	<i>Bidirectional Long Short Term Memory</i> .....	21
2.2.13	<i>Confusion Matrix</i> .....	23
BAB III METODE PENELITIAN.....		25
3.1	Tahapan Penelitian .....	25
3.2	Alat Penelitian .....	25
3.2.1	Perangkat Keras ( <i>Hardware</i> ).....	26
3.2.2	Perangkat Lunak ( <i>Software</i> ) .....	26
3.3	Pengumpulan Data.....	26
3.4	<i>Annotation</i> .....	26
3.5	<i>Preprocessing Data</i> .....	27
3.6	<i>Data Training dan Data Test</i> .....	28
3.7	Ekstraksi Fitur .....	28
3.8	Implementasi Algoritma Klasifikasi.....	29
3.9	Evaluasi .....	29
BAB IV HASIL DAN PEMBAHASAN .....		30
4.1	Dataset.....	30
4.1.1	Pengumpulan Data.....	30
4.1.2	<i>Output Hasil Scraping</i> .....	31
4.1.3	Contoh <i>Dataset</i> .....	32
4.2	<i>Preprocessing Data</i> .....	35
4.3	Ekstraksi Fitur.....	39

4.4	Implementasi Algoritma Bi-LSTM.....	42
4.5	Hasil <i>Training Data</i> .....	46
4.6	Pengukuran Algoritma.....	47
BAB V PENUTUP.....		50
5.1	Kesimpulan.....	50
5.2	Saran.....	50
DAFTAR PUSTAKA .....		51



**DAFTAR TABEL**

Tabel 2.1 Perbandingan Penelitian.....	9
Tabel 2.2 <i>Confusion Matrix</i> .....	23
Tabel 4.1 Deskripsi atribut dataset.....	33
Tabel 4.2 Contoh dataset pasangan pertanyaan .....	34
Tabel 4.3 Hasil Data <i>Cleaning</i> .....	35
Tabel 4.4 Hasil <i>Case Folding</i> .....	36
Tabel 4.5 Hasil <i>Stopwords Removal</i> .....	37
Tabel 4.6 Hasil <i>Tokenizing</i> .....	38
Tabel 4.7 <i>Confusion Matrix</i> .....	48



## DAFTAR GAMBAR

Gambar 2.1 Ilustrasi <i>Deep Learning</i> .....	11
Gambar 2.2 Alur proses <i>text mining</i> .....	13
Gambar 2.3 <i>Word Embedding</i> .....	15
Gambar 2.4 Arsitektur <i>Bidirectional Long Short Term Memory</i> .....	22
Gambar 3.1 Diagram Alur Tahapan Penelitian.....	25
Gambar 4.1 <i>Scraping</i> data menggunakan <i>quora-scrape</i> .....	30
Gambar 4.2 Hasil <i>Scraping</i> menggunakan <i>quora-scrape</i> .....	31
Gambar 4.3 Contoh isi data <i>file output</i> .....	32
Gambar 4.4 <i>Dataset</i> .....	33
Gambar 4.5 Visual Perbandingan Label .....	34
Gambar 4.6 <i>Script Case Folding</i> .....	36
Gambar 4.7 <i>Script Stopwords Removal</i> .....	37
Gambar 4.8 <i>Script Tokenizing</i> .....	38
Gambar 4.9 Contoh token yang dihasilkan.....	39
Gambar 4.10 Contoh <i>Word Index</i> .....	39
Gambar 4.11 <i>Script pre-trained GloVe</i> .....	40
Gambar 4.12 Contoh Hasil <i>Word Embedding GloVe</i> .....	41
Gambar 4.13 <i>Script Pad Sequences</i> .....	41
Gambar 4.14 Contoh Hasil <i>Pad Sequences</i> .....	42
Gambar 4.15 <i>Script</i> Simpan Data.....	42
Gambar 4.16 <i>Script Load</i> dan <i>Split Data</i> .....	43
Gambar 4.17 <i>Script Build Model</i> .....	43
Gambar 4.18 Total Parameter .....	44
Gambar 4.19 <i>Script Training Model Bi-LSTM</i> .....	45
Gambar 4.20 <i>Architecture Model Bi-LSTM</i> .....	45
Gambar 4.21 Hasil <i>Training</i> dan <i>Validasi</i> .....	46
Gambar 4.22 Hasil <i>Confusion Matrix</i> .....	47
Gambar 4.23 <i>Script Confusion Matrix</i> .....	48

## INTISARI

Forum diskusi *online* merupakan tempat untuk berkumpul pengguna internet yang di dalamnya kita bisa berbagi informasi dan berdiskusi antar pengguna forum mengenai suatu topik tertentu. Mengingat semakin bertambahnya pengguna internet saat ini juga membuat forum diskusi semakin ramai dan aktif dalam melakukan tanya jawab. Dengan banyaknya pengguna menyebabkan setiap pengguna di dalam forum menanyakan pertanyaan dengan maksud yang sama namun dengan kata yang berbeda. Oleh sebab itu penelitian ini menggunakan data pertanyaan forum tanya jawab untuk melakukan identifikasi duplikat pertanyaan yang diperoleh dari situs web Quora Indonesia.

Pada penelitian ini data diklasifikasikan menggunakan *Bidirectional Long Short Term Memory* atau Bi-LSTM. Pertanyaan yang diidentifikasi akan memiliki 2 label yaitu *duplicate* dan *not\_duplicate*. Untuk merepresentasikan kata ke dalam vektor, penelitian ini menggunakan *word embedding*. Penelitian ini bertujuan untuk mengetahui berapa hasil akurasi dari *word embedding* yang digunakan dalam identifikasi duplikat pasangan pertanyaan menggunakan Bahasa Indonesia.

Metode *word embedding* yang digunakan adalah *Global Vector* (GloVe). Dari percobaan *word embedding* metode *Global Vector* (GloVe) yang dilakukan dengan pembagian data *training* 80% dan data *test* 20% dari total jumlah dataset 10000 mendapatkan hasil akurasi 95% dengan presisi 96%, *recall* 93% dan *f1-score* 94%.

**Kata-kunci:** Quora, Identifikasi Pertanyaan, *Word Embedding*, *GloVe*, *Bi-LSTM*

## **ABSTRACT**

*An online discussion forum is a place for internet users to gather in which we can share information and discuss among forum users on a particular topic. Considering the increasing number of internet users nowadays it also makes discussion forums more crowded and active in conducting questions and answers. With so many users it causes every user in the forum to ask questions with the same intent but with different words. Therefore this study uses question and answer forum question data to identify duplicate questions obtained from the Quora Indonesia website.*

*In this study, the data were classified using Bidirectional Long Short Term Memory or Bi-LSTM. Identified questions will have 2 labels namely duplicate and not\_duplicate. To represent words into vectors, this study uses word embedding. This study aims to determine the accuracy results of the word embedding used in the identification of duplicate question pairs using Indonesian.*

*The word embedding method used is Global Vector (GloVe). From the word embedding experiment with the Global Vector (GloVe) method, which was carried out by sharing 80% training data and 20% test data of the total number of 10000 datasets, the results obtained 95% accuracy with 96% precision, 93% recall and 94% f1-score.*

**Keywords:** *Quora, Identification questions pairs, Word Embedding, GloVe, Bi-LSTM*