

**METODE "BIDIRECTIONAL ENCODER REPRESENTATIONS FROM  
TRANSFORMERS" (BERT) DALAM PERHITUNGAN  
*SEMANTIC SIMILARITY***

**SKRIPSI**

Untuk memenuhi sebagian persyaratan  
Mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Dhimas Yoga Pratama**

**17.11.0972**

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

**METODE "BIDIRECTIONAL ENCODER REPRESENTATIONS FROM  
TRANSFORMERS" (BERT) DALAM PERHITUNGAN  
*SEMANTIC SIMILARITY***

**SKRIPSI**

Untuk memenuhi sebagian persyaratan  
Mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Dhimas Yoga Pratama**

**17.11.0972**

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

**PERSETUJUAN**

**SKRIPSI**

**METODE "BIDIRECTIONAL ENCODER REPRESENTATIONS FROM  
TRANSFORMERS" (BERT) DALAM PERHITUNGAN  
*SEMANTIC SIMILARITY***

yang dipersiapkan dan disusun oleh

**Dhimas Yoga Pratama**

**17.11.0972**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 25 November 2020

**Dosen Pembimbing,**

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

**NIK. 190302108**

**PENGESAHAN****SKRIPSI****METODE "BIDIRECTIONAL ENCODER REPRESENTATIONS FROM  
TRANSFORMERS" (BERT) DALAM PERHITUNGAN  
*SEMANTIC SIMILARITY***

yang dipersiapkan dan disusun oleh

**Dhimas Yoga Pratama**

**17.11.0972**

telah dipertahankan di depan Dewan Penguji  
pada tanggal 18 November 2020

**Susunan Dewan Penguji**

**Nama Penguji**

**Tanda Tangan**

**Ainul Yaqin, M.Kom**

**NIK. 190302255**

**Bayu Setiaji, M.Kom**

**NIK. 190302216**

**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, 25 November 2020



Dhimas Yoga Pratama

NIM 17.11.0972

## MOTTO



*“Sesungguhnya bersama kesulitan pasti ada kemudahan. Maka apabila engkau telah selesai (dari suatu urusan), tetaplah bekerja keras (untuk urusan yang lain)”*

*(QS 94:6-7)*

*“The best thing you can ever do is follow your dreams”*

*(James Owen Sullivan)*

## PERSEMBAHAN

Alhamdulillahirabbil'alamin, segala puji bagi Allah SWT atas izin dan rahmat-Nya serta doa, usaha, juga kerja keras skripsi ini dapat diselesaikan dengan baik. Pada penyusunan skripsi ini penulis tentu dibantu oleh berbagai pihak, untuk itu dengan ini penulis mempersembahkan ucapan terimakasih sedalam-dalamnya kepada :

1. Ibu saya Dini Sintya Anggraeni dan bapak saya Widayanto yang senantiasa memberikan doa serta dukungan.
2. Segenap keluarga yang telah hadir sebagai lingkungan dimana penulis dididik, tumbuh dan berkembang hingga dewasa.
3. Dosen pembimbing saya Ibu Mardhiya Hayaty, S.T., M.Kom. yang telah membimbing saya dengan sabar dan baik dari awal hingga skripsi ini selesai.
4. Dosen Universitas AMIKOM Yogyakarta yang telah memberikan banyak ilmu yang bermanfaat selama perkuliahan.
5. Alde Satriayu Putri Nadeak yang selalu mendampingi, memberikan dukungan, motivasi serta doa.
6. Rekan sepenelitian saya Anjas, Andri, Lutfi, Dimas Midyan yang selalu membantu, menemani, dan berbagi ilmu selama pengerjaan skripsi ini.
7. Segenap teman-teman IF-01 khususnya Aldo, Parmin, Zul, Roy, Irfan, Dedi, Dendy, Dandy, Jeki, Fariz, dan Yoga yang telah menemani, berbagi canda tawa, duka, dan juga sejumlah uang selama 3 tahun masa perkuliahan.
8. Pihak lain yang tidak dapat saya sebutkan satu persatu.

## KATA PENGANTAR

Penulis panjatkan puji dan syukur kehadirat Allah SWT atas izin, berkah, dan rahmat-Nya sehingga penulis dapat menyelesaikan skripsi ini dengan baik yang berjudul **“METODE *BIDIRECTIONAL ENCODER REPRESENTATION FROM TRANSFORMERS* (BERT) DALAM PERHITUNGAN *SEMANTIC SIMILARITY*”** disusun sebagai salah satu syarat utama untuk menyelesaikan program sarjana di Universitas AMIKOM Yogyakarta. Penyelesaian skripsi ini juga tidak lepas dari bantuan berbagai pihak, karena itu pada kesempatan ini penulis ingin menyampaikan rasa hormat dan terimakasih 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 telah memberikan bimbingan dengan baik selama penyusunan skripsi ini.
5. Bapak Ainul Yaqin, M.Kom. dan Bapak Bayu Setiaji, M.Kom. selaku dosen penguji.

Penulis menyadari bahwa skripsi ini masih terdapat kekurangan. Maka penulis menerima segala kritik dan saran yang membangun serta teguran dari semua pihak. Semoga skripsi ini dapat bermanfaat bagi penulis, pembaca, dan perkembangan *Deep Learning* di Indonesia. Atas saran dan kritik penulis ucapkan terimakasih.

Yogyakarta, 25 November 2020



Dhimas Yoga Pratama

## DAFTAR ISI

PERSETUJUAN.....	ii
PENGESAHAN.....	iii
PERNYATAAN.....	iv
MOTTO.....	v
PERSEMBAHAN.....	vi
KATA PENGANTAR.....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	xi
DAFTAR GAMBAR.....	xii
INTISARI.....	xiv
ABSTRACT.....	xv
BAB I.....	1
PENDAHULUAN.....	1
<b>1.1. Latar Belakang</b> .....	1
<b>1.2. Rumusan Masalah</b> .....	3
<b>1.3. Batasan Masalah</b> .....	3
<b>1.4. Tujuan Penelitian</b> .....	3
<b>1.5. Manfaat Penelitian</b> .....	4
<b>1.6. Metodologi Penelitian</b> .....	4
<b>1.6.1 Pengumpulan <i>Dataset</i></b> .....	4
<b>1.6.2 Pengisian Data Pertanyaan Uji dan <i>Labelling</i></b> .....	5
<b>1.6.3 <i>Preprocessing Data</i></b> .....	5
<b>1.6.4 Tahap <i>Fine-Tuning</i></b> .....	5
<b>1.6.5 Pengujian</b> .....	6
<b>1.6.6 Evaluasi</b> .....	6
<b>1.7. Sistematika Penulisan</b> .....	6

BAB II.....	9
LANDASAN TEORI.....	9
2.1 <b>Tinjauan Pustaka</b> .....	9
2.2 <b>Natural Language Processing</b> .....	12
2.3 <b>Data Mining</b> .....	13
2.4 <b>Text Mining</b> .....	14
2.5 <b>Deep Learning</b> .....	15
2.6 <b>Word Embedding</b> .....	17
2.7 <b>Semantic Similarity</b> .....	18
2.8 <b>Regresi Logistik (Logistic Regression)</b> .....	18
2.9 <b>Scraping</b> .....	20
2.10 <b>Text Preprocessing</b> .....	21
2.10.1 <b>Case Folding</b> .....	21
2.10.2 <b>Tokenization</b> .....	22
2.10.3 <b>Stopword Removal</b> .....	22
2.10.4 <b>Stemming</b> .....	22
2.11 <b>Bidirectional Encoder Representations from Transformers (BERT)</b> .....	23
2.11.1 <b>Arsitektur Model</b> .....	24
2.11.2 <b>Representasi Input/Output</b> .....	24
2.11.3 <b>Pre-training BERT</b> .....	26
2.11.4 <b>Fine-tuning BERT</b> .....	27
2.12 <b>Optimasi Adam</b> .....	28
2.13 <b>Confusion Matrix</b> .....	30
BAB III.....	33
METODE PENELITIAN.....	33
3.1 <b>Tahapan Penelitian</b> .....	33
3.2 <b>Alat Penelitian</b> .....	33
3.3 <b>Pengumpulan Data</b> .....	34
3.4 <b>Dataset</b> .....	35
3.5 <b>Preprocessing Data</b> .....	35
3.6 <b>Mengambil Model Pre-Train</b> .....	36
3.7 <b>Melakukan Training</b> .....	36

<b>3.8 Pengujian dan Evaluasi</b> .....	37
BAB IV .....	39
HASIL DAN PEMBAHASAN .....	39
<b>4.1 Pengumpulan Data</b> .....	39
<b>4.1.1 Proses Scraping</b> .....	39
<b>4.1.2 File output hasil Scraping</b> .....	40
<b>4.1.3 Contoh isi file output</b> .....	41
<b>4.1.4 Contoh dataset</b> .....	41
<b>4.2 Persiapan</b> .....	43
<b>4.3 Mempersiapkan Data Pre-Train</b> .....	45
<b>4.4 Mempersiapkan Dataset untuk Training</b> .....	47
<b>4.5 Melakukan Permodelan</b> .....	49
<b>4.6 Melakukan Testing</b> .....	53
<b>4.7 Mengukur Akurasi Testing</b> .....	56
BAB V .....	63
PENUTUP .....	63
<b>5.1 Kesimpulan</b> .....	63
<b>5.2 Saran</b> .....	63
DAFTAR PUSTAKA .....	65

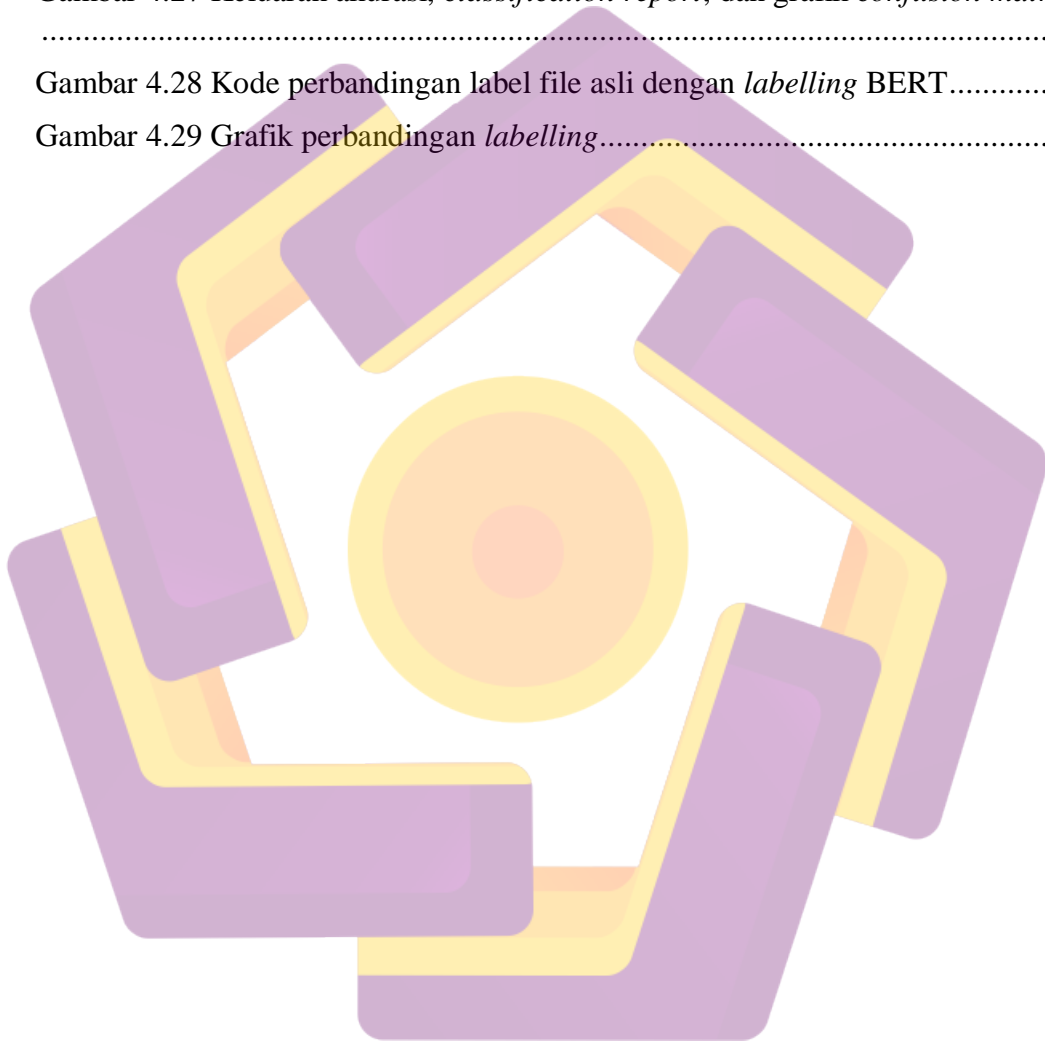
## DAFTAR TABEL

Tabel 2.1 Tabel Perbandingan Penelitian.....	11
Tabel 2.1 Proses <i>case folding</i> .....	21
Tabel 2.2 Proses <i>tokenization</i> .....	22
Tabel 2.3 Proses <i>stopword removal</i> .....	22
Tabel 2.4 Proses <i>stemming</i> .....	23
Tabel 2.5 Confusion Matrix .....	30
Tabel 3.1 Tabel alat penelitian .....	34
Tabel 3.1 Proses <i>case folding</i> .....	35
Tabel 3.2 Proses <i>tokenization</i> .....	36
Tabel 4.1 Keterangan <i>header</i> pada file <i>dataset</i> .....	42
Tabel 4.2 Sampel Pasangan Pertanyaan.....	43
Gambar 4.18 Visualisasi hasil <i>training</i> dengan 6 <i>epoch</i> .....	52
Gambar 4.19 Visualisasi hasil <i>training</i> dengan 10 <i>epoch</i> .....	52
Tabel 4.3 <i>Output</i> Confusion Matrix.....	61

## DAFTAR GAMBAR

Gambar 2.1 Ilustrasi <i>Deep Learning</i> .....	17
Gambar 2.2 Ilustrasi <i>web scraping</i> .....	21
Gambar 2.3 Visualisasi sederhana kinerja <i>encoder</i> dan <i>decoder</i> .....	24
Gambar 2.4 : Ilustrasi proses <i>pre-training</i> pada BERT .....	25
Gambar 2.5 Representasi input pada BERT.....	26
Gambar 2.6 Visualisasi proses <i>fine-tuning</i> pada BERT.....	28
Gambar 3.1 <i>Flowchart</i> tahapan penelitian .....	33
Gambar 4.1 Proses <i>scraping</i> melalui <i>terminal</i> .....	40
Gambar 4.2 File <i>output</i> hasil <i>scraping</i> .....	40
Gambar 4.3 Isi file <i>output</i> hasil <i>scraping</i> .....	41
Gambar 4.4 File <i>dataset</i> .....	42
Gambar 4.5 Perintah untuk memasang library <i>transformers</i> .....	44
Gambar 4.6 Perintah untuk <i>clone repository</i> indonlu.....	44
Gambar 4.7 Kode untuk mengimpor <i>library</i> dan <i>function</i> .....	45
Gambar 4.8 Kode untuk mempersiapkan <i>data pre-train</i> .....	45
Gambar 4.9 Kode untuk memeriksa kinerja <i>class tokenizer</i> .....	46
Gambar 4.10 Hasil keluaran proses <i>tokenizing</i> .....	46
Gambar 4.11 Deklarasi masing-masing file <i>dataset</i> .....	47
Gambar 4.12 Tahap <i>preprocessing</i> pada <i>dataset</i> .....	48
Gambar 4.13 Memberi konfigurasi tambahan pada <i>dataset</i> .....	48
Gambar 4.14 <i>Optimizer</i> Adam dan pengiriman <i>model</i> ke GPU.....	49
Gambar 4.15 Kode untuk menjalankan <i>training</i> .....	50
Gambar 4.16 Kode untuk evaluasi <i>dev set</i> .....	50
Gambar 4.17 <i>Output</i> yang ditampilkan ke layar pada saat <i>training</i> .....	51
Gambar 4.17 Visualisasi proses <i>training</i> dan evaluasi .....	51
Gambar 4.20 Kode untuk melakukan <i>testing</i> .....	53
Gambar 4.21 Fungsi <i>forward_sequence_classification</i> .....	54

Gambar 4.22 Isi variabel <i>logits</i> .....	55
Gambar 4.23 Output fungsi <i>torch.topk()</i> .....	55
Gambar 4.24 Skema pengujian dengan <i>typo</i> pada kalimat .....	56
Gambar 4.25 Menyalin label dari file <i>pred.txt</i> ke dalam <i>array</i> .....	57
Gambar 4.26 Perhitungan akurasi, <i>confusion matrix</i> , dan <i>classification report</i> ...	58
Gambar 4.27 Keluaran akurasi, <i>classification report</i> , dan grafik <i>confusion matrix</i> .....	59
Gambar 4.28 Kode perbandingan label file asli dengan <i>labelling BERT</i> .....	60
Gambar 4.29 Grafik perbandingan <i>labelling</i> .....	60



## INTISARI

Quora adalah situs forum *online* yang digunakan oleh banyak pengguna dari seluruh dunia untuk melakukan tanya jawab mengenai topik tertentu. Mengingat banyaknya pengguna yang menggunakan layanan ini dan aktif dalam melakukan tanya jawab menyebabkan setiap pengguna menanyakan pertanyaan dengan makna yang sama namun dengan kata yang berbeda, inilah yang menyebabkan adanya masalah duplikasi pertanyaan. Oleh sebab itu diperlukan metode yang dapat mengidentifikasi adanya duplikasi dari pertanyaan yang di-*inputkan* pengguna.

Penelitian ini menggunakan sebuah metode dalam *deep learning* yaitu *Bidirectional Encoder Representation from Transformer* atau BERT sebagai metode untuk membuat *word embedding*. Untuk memprediksi probabilitas kemiripan pertanyaan digunakan metode *logistic regression*. Pertanyaan yang diidentifikasi akan memiliki 2 label yaitu *duplicate* dan *not\_duplicate*.

Proses permodelan dilakukan sebanyak 8 *epochs*. Data yang digunakan merupakan judul pertanyaan yang diambil dari situs Quora dengan teknik *scraping*. Total pasangan judul pertanyaan adalah sebanyak 10.000 pasang dengan pembagian data latih sebanyak 70%, data evaluasi sebanyak 10%, dan data uji sebanyak 20%. Dari hasil pengujian prediksi probabilitas kemiripan didapatkan nilai akurasi sebesar 99% .

**Kata Kunci** : Quora, *Deep Learning*, BERT, *Word Embedding*, *Logistic Regression*.

## ABSTRACT

*Quora is an online forum used by many users from all over the world to ask and answer questions about certain topics. Given the large number who use this service and are active in questioning and answering, it causes each user to ask the same question(s) but in different way, this then causes the duplication of questions problem. Therefore, we need a method that can identify the duplication of user-inputted questions.*

*This research use a method in Deep Learning namely Bidirectional Encoder Representation from Transformers or BERT as a method to create word embedding. Logistic regression is then used to predict the probability of the similarity of the questions. The identified questions will have 2 labels namely duplicate and not\_duplicate.*

*The modelling process is carried out in 8 epochs. The data used is the title of the questions taken from Quora site using scraping technique. The total number of pairs of questions are 10.000 pairs, which then splitted into training data as much as 70%, evaluation data as much as 10%, and training data as much as 20%. From the testing of probability of the similarity prediction obtained accuracy in the amount of 99%.*

**Keywords:** *Quora, Deep Learning, BERT, Word Embedding, Logistic Regression.*