

**MESIN PENERJEMAH BAHASA INGGRIS – INDONESIA
BERBASIS JARINGAN SARAF TIRUAN DENGAN MEKANISME
ATTENTION MENGGUNAKAN ARSITEKTUR TRANSFORMER**



SKRIPSI

**Disusun Sebagai Salah Satu Syarat
Untuk Memperoleh Gelar Sarjana Komputer
pada Departemen Ilmu Komputer/ Informatika**

Disusun Oleh :

PRASASTO ADI WISMOYO

24010311130026

**DEPARTEMEN ILMU KOMPUTER/ INFORMATIKA
FAKULTAS SAINS DAN MATEMATIKA
UNIVERSITAS DIPONEGORO**

2018

HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Saya yang bertanda tangan di bawah ini:

Nama : Prasasto Adi Wismoyo

NIM : 24010311130026

Judul : Mesin Penerjemah Bahasa Inggris – Indonesia Berbasis Jaringan Saraf Tiruan
dengan Mekanisme Attention Menggunakan Arsitektur Transformer

Dengan ini saya menyatakan bahwa dalam tugas akhir/ skripsi ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar kesarjanaan di suatu Perguruan Tinggi, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan di dalam daftar pustaka.

Semarang, 8 Juni 2018



Prasasto Adi Wismoyo
24010311130026

HALAMAN PENGESAHAN

Judul : Mesin Penerjemah Bahasa Inggris – Indonesia Berbasis Jaringan Saraf Tiruan
dengan Mekanisme Attention Menggunakan Arsitektur Transformer

Nama : Prasasto Adi Wismoyo

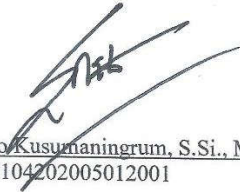
NIM : 24010311130026

Telah diujikan pada sidang tugas akhir pada tanggal 8 Juni 2018 dan dinyatakan lulus
pada tanggal **8 Juni 2018**.

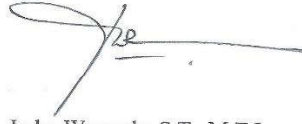
Semarang, 5 Juli 2018

Mengetahui,

Ketua Departemen Ilmu Komputer/Informatika
FSM UNDIP


Dr. Retno Kusumaningrum, S.Si., M.Kom.
NIP. 198104202005012001

Panitia Penguji Tugas Akhir
Ketua,


Indra Waspada, S.T., M.T.I.
NIP. 197902122008121002

HALAMAN PENGESAHAN

Judul : Mesin Penerjemah Bahasa Inggris – Indonesia Berbasis Jaringan Saraf Tiruan
dengan Mekanisme Attention Menggunakan Arsitektur Transformer

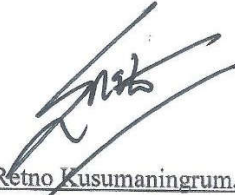
Nama : Prasasto Adi Wismoyo

NIM : 24010311130026

Telah diujikan pada sidang tugas akhir pada tanggal 8 Juni 2018.

Semarang, 5 Juli 2018

Pembimbing



Dr. Retno Kusumaningrum, S.Si., M.Kom.
NIP. 198104202005012001

ABSTRAK

Bahasa Inggris sebagai bahasa pengantar internasional dan bahasa Indonesia sebagai bahasa pengantar nasional memiliki peran yang sangat penting bagi bangsa Indonesia. Sehingga mesin penerjemah (*machine translation*) dirasa sangat membantu dalam proses alih bahasa antara keduanya. Beberapa tahun terakhir mesin penerjemah mengalami perkembangan yang sangat pesat melalui pendekatan mesin penerjemah berbasis jaringan saraf tiruan (*neural machine translation*). Transformer merupakan salah satu arsitektur mesin penerjemah berbasis jaringan saraf tiruan yang mendapatkan hasil *state-of-the-art* dalam beberapa penelitian. Arsitektur ini menggunakan mekanisme *attention* tanpa adanya lapisan *recurrent* maupun *convolution*. Penulis melakukan penelitian untuk mengetahui kualitas terjemahan bahasa Inggris – Indonesia yang dihasilkan oleh Transformer berdasarkan metrik penilaian BLEU. Hasil pengujian menggunakan korpus uji *dev2010*, *tst2010*, *tst2017plus*, dan *test.go.id* menunjukkan bahwa model mesin penerjemah mendapatkan skor BLEU 24,83% (*dev2010*) 24,14% (*tst2010*), 27,31% (*tst2017plus*), dan 34,18% (*test.go.id*) untuk arah terjemahan Inggris – Indonesia. Adapun untuk arah terjemahan Indonesia – Inggris, model mendapatkan skor BLEU 90,13% (*dev2010*), 88,68% (*tst2010*), 90,68% (*tst2010*), dan 69,14% (*test.go.id*).

Kata kunci: mesin penerjemah, jaringan saraf tiruan

ABSTRACT

English as an international lingua franca and Bahasa Indonesia as the national lingua franca have significant role for Indonesian people. Therefore machine translation is considered to be very helpful in translating the two languages. In recent years, machine translation is advancing rapidly through neural machine translation approach. Transformer is one of neural machine translation architectures that achieved state-of-the-art results in several researches. That architecture uses attention mechanism without recurrent nor convolution layer. This research was conducted to discover the quality of English-Indonesian translation generated by Transformer based on BLEU metrics. The test results using test corpus *dev2010*, *tst2010*, *tst2017plus*, and *test.go.id* show that the machine translation model achieved BLEU scores 24.83% (*dev2010*) 24.14% (*tst2010*), 27.31% (*tst2017plus*), and 34, 18% (*test.go.id*) for the English-Indonesian translation direction. As for the Indonesian-English translation direction, the model achieved BLEU scores 90.13% (*dev2010*), 88.68% (*tst2010*), 90.68% (*tst2010*), and 69.14% (*test.go.id*).

Keywords: machine translation, artificial neural network

KATA PENGANTAR

Segala puji syukur bagi Tuhan Yang Maha Esa atas karunia-Nya yang diberikan kepada penulis sehingga dapat menyelesaikan penulisan laporan tugas akhir yang berjudul “Mesin Penerjemah Bahasa Inggris – Indonesia Berbasis Jaringan Saraf Tiruan dengan Mekanisme Attention Menggunakan Arsitektur Transformer”. Laporan tugas akhir ini disusun sebagai salah satu syarat untuk memperoleh gelar sarjana strata satu pada Jurusan Ilmu Komputer/ Informatika Fakultas Sains dan Matematika Universitas Diponegoro Semarang.

Dalam penyusunan laporan ini tentulah banyak mendapat bimbingan dan bantuan dari berbagai pihak. Untuk itu, pada kesempatan ini penulis mengucapkan rasa hormat dan terima kasih kepada:

1. Ibu Dr. Retno Kusumaningrum, S.Si., M.Kom., selaku Ketua Departemen Ilmu Komputer/ Informatika Fakultas Sains dan Matematika Universitas Diponegoro Semarang, sekaligus menjadi Dosen Pembimbing yang telah membantu dalam proses bimbingan hingga selesainya laporan tugas akhir ini.
2. Bapak Helmie Arif Wibawa, S.Si., M.Cs., selaku Koordinator Tugas Akhir sekaligus menjadi Penguji Tugas Akhir.
3. Bapak Indra Waspada, S.T., M.T.I., selaku Ketua Penguji Tugas Akhir.
4. Bapak Sudarwanto, S.H. dan Ibu Adi Sulistyomurti, selaku orang tua yang terus memberikan dukungan dan do'a dalam menyelesaikan tugas akhir ini.
5. Semua pihak yang telah membantu kelancaran dalam pelaksanaan tugas akhir yang tidak dapat penulis sebutkan satu persatu.

Penulis menyadari bahwa dalam laporan ini masih banyak kekurangan baik dari segi materi ataupun dalam penyajiannya karena keterbatasan kemampuan dan pengetahuan penulis. Oleh karena itu, kritik dan saran sangat penulis harapkan. Semoga laporan ini dapat bermanfaat bagi pembaca pada umumnya dan penulis pada khususnya.

Semarang, 8 Juni 2018

Penulis,

Prasasto Adi Wismoyo

24010311130026

DAFTAR ISI

HALAMAN PERNYATAAN KEASLIAN SKRIPSI.....	Error! Bookmark not defined.
HALAMAN PENGESAHAN	ii
HALAMAN PENGESAHAN	iii
ABSTRAK.....	iv
ABSTRACT	vi
KATA PENGANTAR.....	vii
DAFTAR ISI	viii
DAFTAR GAMBAR.....	xi
DAFTAR TABEL	1
BAB I PENDAHULUAN	2
1.1. Latar Belakang.....	2
1.2. Rumusan Masalah.....	3
1.3. Tujuan dan Manfaat	4
1.4. Ruang Lingkup.....	4
1.5. Sistematika Penulisan	4
BAB II TINJAUAN PUSTAKA	6
2.1. Perkembangan Penelitian Mesin Penerjemah Bahasa Inggris - Indonesia ..	6
2.2. Jaringan Saraf Tiruan.....	6
2.3. Machine Translation	10
2.4. Neural Machine Translation	11
2.5. Korpus.....	15
2.6. Tokenisasi	17
2.7. <i>Subword</i>	18
2.8. <i>Embedding</i>	21
2.9. <i>Positional Encoding</i>	23
2.10. Mekanisme <i>Attention</i>	24
2.11. BLEU	27
BAB III METODOLOGI PENELITIAN.....	32

3.1.	Pengumpulan Data	32
3.1.1.	TED Talks	32
3.1.2.	PANL-BPPT	33
3.1.3.	OpenSubtitles2018	33
3.1.4.	Korpus test.go.id.....	34
3.2.	Prapengolahan.....	35
3.3.	Arsitektur	36
3.3.1.	Input.....	36
3.3.2.	Encoder.....	37
3.3.3.	Decoder.....	38
3.3.4.	Feed-Forward Neural Network.....	39
3.3.5.	Output.....	39
3.4.	Pelatihan.....	40
3.5.	Evaluasi.....	40
BAB IV HASIL DAN PEMBAHASAN.....		42
4.1	Implementasi Model	42
4.2	Pengujian <i>Hyperparameter</i>	42
4.2.1	Pencarian Iterasi Mendekati Titik Konvergen	42
4.2.2	Percobaan Learning Rate Warmup Steps.....	45
4.2.3	Percobaan Batch Size	46
4.2.4	Hasil dan Analisis Pencarian Hyperparameter.....	47
4.3	Pengujian Korpus Latih	48
4.3.1	Hasil dan Analisis Pengujian Korpus Latih	49
4.4	Pengujian Model Akhir.....	49
4.4.1	Hasil dan Analisis Pengujian Model Akhir.....	50
BAB V KESIMPULAN.....		52
5.1	Kesimpulan	52
5.2	Saran	53
DAFTAR PUSTAKA.....		54
LAMPIRAN-LAMPIRAN		57
Lampiran 1: Sampel dev2010.en.....		58
Lampiran 2: Sampel dev2010.id		59

Lampiran 3: Sampel tst2010.en.....	61
Lampiran 4: Sampel tst2010.id	63
Lampiran 5: Sampel tst2017plus.en.....	65
Lampiran 6: Sampel tst2017plus.id.....	67
Lampiran 7: Sampel test.go.id.en.....	69
Lampiran 8: Sampel test.go.id.id	71
Lampiran 9: Sampel Hasil Terjemahan Indonesia – Inggris test.go.id	73
Lampiran 10: Sampel Hasil Terjemahan Inggris – Indonesia test.go.id.....	75

DAFTAR GAMBAR

Gambar 2.1. Contoh neuron dengan tiga <i>input</i>	6
Gambar 2.2. <i>Multilayer neural network</i>	8
Gambar 2.3. Model <i>encoder – decoder</i> untuk mesin penerjemah.....	13
Gambar 2.4. Arsitektur RNN (Luong 2016).....	13
Gambar 2.5. Mekanisme <i>attention</i> pada <i>bidirectinoal RNN</i>	14
Gambar 2.6. Visualisasi <i>attention</i> en-fr (Olah and Carter 2016).....	15
Gambar 2.7. Visualisasi lengkap <i>attention</i> en-fr (Olah and Carter 2016).....	15
Gambar 2.8. Ilustrasi pemetaan <i>word embedding</i>	22
Gambar 2.9. Operasi <i>dot product</i> pada <i>embedding</i>	22
Gambar 2.10. RNN.....	24
Gambar 2.11. RNN dengan mekanisme <i>attention</i>	25
Gambar 2.12. <i>Scaled dot-product attention</i>	26
Gambar 2.13. <i>Multi-head attention</i>	26
Gambar 2.14. Perbandingan n-gram.....	29
Gambar 2.15: Multireferensi.	30
Gambar 3.1. Arsitektur <i>Transformer</i> (Vaswani et al. 2017).	36
Gambar 3.2. <i>Input embedding</i>	37
Gambar 3.3. <i>Encoder</i>	37
Gambar 3.4. <i>Encoder</i> paling awal.	38
Gambar 3.5. <i>Decoder</i>	38
Gambar 3.6. Operasi <i>Softmax</i>	39
Gambar 4.1. Kurva pembelajaran pencarian iterasi.	44
Gambar 4.2. Kurva Eksperimen LRWS 1 dan 2	45
Gambar 4.3. Kurva Eksperimen LRWS 3 dan 4.	46
Gambar 4.4. Percobaan <i>batch size</i>	47

DAFTAR TABEL

Tabel 2.1. Perkembangan Penelitian.....	6
Tabel 2.2. Contoh kalimat sebelum tokenisasi.....	18
Tabel 2.3. Contoh kalimat setelah tokenisasi.....	18
Tabel 2.4. Contoh kamus dengan N-kata.....	19
Tabel 2.5. Contoh kamus menggunakan <i>subword</i>	20
Tabel 2.6. Hasil penerapan <i>subword</i> pada kalimat.....	20
Tabel 2.7. Contoh representasi <i>one-hot encoding</i>	21
Tabel 2.8. Indeks <i>subword</i>	23
Tabel 2.9. Penilaian <i>Fluency</i> dan <i>Adequacy</i>	27
Tabel 2.10. Contoh kalimat sumber dan target untuk perhitungan BLEU.....	29
Tabel 2.11. Contoh kalimat prediksi sistem A dan sistem B.....	29
Tabel 2.12. Presisi n-gram.....	29
Tabel 2.13. Metrik penilaian.....	31
Tabel 3.1. IWSLT17.....	32
Tabel 3.2. PANL-BPPT.....	33
Tabel 3.3. OpenSubtitles2018.....	33
Tabel 3.4. Test.go.id.....	34
Tabel 4.1. Pencarian Iterasi.....	43
Tabel 4.2. Hasil Pencarian Iterasi.....	43
Tabel 4.3. Eksperimen LRWS 1 dan 2.....	45
Tabel 4.4. Eksperimen LRWS 3 dan 4.....	45
Tabel 4.5. Percobaan <i>batch size</i>	46
Tabel 4.6. <i>Hyperparameters</i>	48
Tabel 4.7. Hasil uji terhadap <i>dev2010</i>	49
Tabel 4.8. Hasil Pengujian Model Akhir.....	50
Tabel 4.9. Perbandingan Hasil dengan Layanan <i>Online</i>	51

BAB I

PENDAHULUAN

1.1. Latar Belakang

Keberagaman suku dan budaya di Indonesia memberikan pengaruh terhadap bahasa yang digunakan sebagai sarana komunikasi dalam kehidupan bermasyarakat. Untuk menjembatani perbedaan bahasa antarsuku, antardaerah, dan antarbudaya, ditetapkanlah bahasa Indonesia sebagai bahasa pengantar (*lingua franca*) demi kelancaran interaksi sosial dan integrasi bangsa. Asal mula bahasa Indonesia berpangkal dari bahasa Melayu, bahasa yang ketika itu menjadi bahasa perdagangan bagi para saudagar di wilayah Nusantara (Mohamed and Harahap 2013).

Selain bahasa Indonesia dan bahasa daerah, bahasa asing juga digunakan dan diajarkan di Indonesia untuk mendukung komunikasi antarbangsa dalam ruang lingkup internasional. Bahasa asing yang banyak digunakan dan diajarkan di Indonesia adalah bahasa Inggris dan bahasa Arab. Bahasa Inggris merupakan bahasa pengantar dalam komunikasi beragam tema pembicaraan di kancah internasional. Munculnya bahasa Inggris sebagai bahasa pengantar internasional tak lepas dari peran kekuasaan politik dan perdagangan Kerajaan Inggris pada abad 19 dan 20 (Kloss 1967). Sementara itu, Bahasa Arab umumnya digunakan sebagai bahasa pengantar dalam komunikasi antar umat Islam, khususnya komunikasi yang berkaitan dengan tema keagamaan (Arifin and Yundiafi 2014; Kloss 1967).

Merupakan suatu kelaziman bahwa bahasa pengantar digunakan dalam pertukaran informasi di segala bidang seperti perdagangan, pendidikan, teknologi, diplomatik, dan lain sebagainya. Sehingga penggunaan bahasa Indonesia sebagai bahasa pengantar nasional serta bahasa Inggris sebagai bahasa pengantar internasional berimplikasi terhadap tingginya kebutuhan alih bahasa di antara keduanya. Kamus Inggris – Indonesia menjadi acuan primer untuk menemukan kata atau terminologi yang belum dimengerti ketika memahami atau menerjemahkan teks informasi. Penggunaan kamus cetak pada era digital sudah semakin tergantikan oleh kamus digital yang memiliki beberapa keunggulan seperti portabilitas tinggi, kemudahan pencarian lema dan keterbaruan entri pada basis data. Alat bantu penerjemahan selain kamus dwibahasa dapat berupa thesaurus, memori penerjemahan (*translation memory*), pemeriksa ejaan (*spell checker*), pemeriksa tata bahasa (*grammar checker*), maupun mesin penerjemah (*machine translation*).

Mesin penerjemah merupakan perangkat lunak penerjemah bahasa yang proses penerjemahannya dilakukan oleh mesin (komputer) secara otomatis dan hasilnya dapat langsung disajikan kepada pengguna. Opsi penggunaan mesin penerjemah lebih umum dipilih karena cepat dan praktis dalam melakukan penerjemahan teks, berbanding terbalik dengan penerjemahan oleh manusia yang memerlukan waktu lama. Meskipun penerjemah manusia yang berkompeten akan menghasilkan penerjemahan yang lebih baik dan dalam beberapa keadaan tak tergantikan oleh penerjemah otomatis, namun ada banyak kasus di mana cukup untuk diterjemahkan menggunakan mesin penerjemah. Selain itu, hasil penerjemahan oleh mesin pun tidak serta merta digunakan begitu saja, namun disunting kembali oleh manusia. Sehingga penerjemah manusia berkualifikasi tinggi pun juga menggunakan mesin penerjemah untuk mempercepat proses penerjemahan.

Penelitian terkait mesin penerjemah telah banyak dilakukan dan termasuk bidang penelitian yang sangat aktif. Namun penelitian mesin penerjemah bahasa Inggris – Indonesia tergolong masih sedikit ditemui terutama metode-metode berbasis jaringan saraf tiruan yang merupakan *state-of-the-art* untuk mesin penerjemah. *Neural machine translation* merupakan istilah yang digunakan untuk menyebut metode mesin penerjemah dengan arsitektur berbasis jaringan saraf tiruan. Metode-metode yang paling umum digunakan adalah RNN (*Recurrent Neural Network*), CNN (*Convolutional Neural Network*), dan *attention*. Pada tugas akhir ini, *attention* dipilih karena model dasarnya sudah cukup untuk menghasilkan terjemahan yang baik serta kompleksitas waktu komputasi dan ruang memori yang relatif lebih rendah dibandingkan metode lainnya (Chen et al. 2018; Vaswani et al. 2017). Arsitektur *Transformer* adalah arsitektur mesin penerjemah yang menerapkan mekanisme *attention* tanpa adanya tambahan RNN ataupun CNN (Vaswani et al. 2017).

1.2. Rumusan Masalah

Berdasarkan penjelasan latar belakang yang telah disampaikan, dapat disusun rumusan masalah yaitu bagaimana menerapkan metode jaringan saraf tiruan dengan mekanisme *attention* menggunakan arsitektur *Transformer* pada mesin penerjemah Bahasa Inggris – Indonesia.

1.3. Tujuan dan Manfaat

Tujuan yang ingin dicapai dari penelitian tugas akhir ini adalah mengetahui kualitas model penerjemahan Bahasa Inggris – Indonesia berbasis jaringan saraf tiruan dengan mekanisme *attention* menggunakan arsitektur *Transformer* berdasarkan matrik skor BLEU.

Adapun manfaat yang diharapkan dari penelitian tugas akhir ini adalah dapat meningkatkan peran mesin penerjemah dalam kegiatan masyarakat, penelitian, maupun industri.

1.4. Ruang Lingkup

Ruang lingkup penelitian tugas akhir ini adalah sebagai berikut :

1. Korpus latih berupa korpus paralel bahasa Inggris – Indonesia dari PANL-BPPT, IWSLT17, dan OpenSubtitles2018.
2. Korpus uji berupa korpus paralel dari IWSLT17 dan korpus uji yang disusun oleh penulis.
3. Metode penelitian berbasis jaringan saraf tiruan dengan mekanisme *attention* menggunakan arsitektur *Transformer*.
4. Evaluasi kualitas model penerjemahan menggunakan matrik skor BLEU.
5. Output hasil terjemahan berupa *plain text* dengan pemisah kalimat hasil terjemahan berupa baris baru.

1.5. Sistematika Penulisan

Sistematika penulisan yang digunakan dalam laporan tugas akhir ini terbagi menjadi beberapa pokok bahasan, yaitu:

BAB I PENDAHULUAN

Bab ini membahas latar belakang, rumusan masalah, tujuan dan manfaat, ruang lingkup, serta sistematika penulisan laporan.

BAB II TINJAUAN PUSTAKA

Bab ini menjelaskan mengenai teori-teori serta istilah-istilah yang digunakan dalam penelitian tugas akhir meliputi perkembangan penelitian mesin penerjemah bahasa Inggris - Indonesia ,jaringan saraf tiruan, *machine translation*, *neural machine translation*, korpus, tokenisasi,

subword, embedding, positional encoding, mekanisme attention, dan BLEU.

BAB III METODOLOGI PENELITIAN

Bab ini menjelaskan tahapan yang dilakukan dalam penelitian tugas akhir. Tahapan tersebut meliputi pengumpulan data, prapengolahan, arsitektur, pelatihan, dan evaluasi.

BAB IV HASIL DAN PEMBAHASAN

Bab ini menjelaskan hasil eksperimen yang telah dilakukan berupa implementasi model, pengujian *hyperparameter*, serta pengujian model akhir.

BAB V PENUTUP

Bab ini merupakan kesimpulan dari bab-bab yang dibahas sebelumnya dan saran bagi penelitian selanjutnya.