

***NAMED ENTITY RECOGNITION* BERITA BAHASA INDONESIA  
DENGAN POLYGLOT**

**SKRIPSI**



disusun oleh

**Barep Setiyadi**

**19.21.1330**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2021**

**NAMED ENTITY RECOGNITION BERITA BAHASA  
INDONESIA DENGAN POLYGLOT**

**SKRIPSI**

untuk memenuhi sebagian persyaratan  
mencapai gelar Sarjana  
pada Program Studi Informatika



disusun oleh

**Barep Setiyadi**

**19.21.1330**

**PROGRAM SARJANA  
PROGRAM STUDI INFORMATIKA  
FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2021**

# **PERSETUJUAN**

## **SKRIPSI**

### ***NAMED ENTITY RECOGNITION* BERITA BAHASA INDONESIA DENGAN POLYGLOT**

yang dipersiapkan dan disusun oleh

**Barep Setiyadi**

**19.21.1330**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 5 Agustus 2021

**Dosen Pembimbing,**

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

**NIK. 190302108**

# PENGESAHAN

## SKRIPSI

### *NAMED ENTITY RECOGNITION* BERITA BAHASA INDONESIA DENGAN POLYGLOT

yang dipersiapkan dan disusun oleh

**Barep Setiyadi**

**19.21.1330**

telah dipertahankan di depan Dewan Penguji

pada tanggal 19 Juli 2021

**Susunan Dewan Penguji**

**Nama Penguji**

**Tanda Tangan**

**Anggit Dwi Hartanto, M.Kom.**

**NIK. 190302163**

**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 19 Juli 2021

**DEKAN FAKULTAS ILMU KOMPUTER**

**Hanif Al Fatta, M.Kom.**

**NIK. 190302096**


## PERNYATAAN

### PERNYATAAN

Saya yang bertanda tangan di bawah ini menyatakan bahwa, skripsi ini merupakan karya saya sendiri (ASLI), dan isi dalam skripsi ini tidak terdapat karya yang pernah diajukan 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 naskah dan karya yang telah dibuat adalah menjadi tanggung jawab saya pribadi.

Tegal, 20 Agustus 2021

  
Barep Setiyadi

NIM. 19.21.1330



## **MOTTO**

”Melakukan sesuatu selalu dengan alasan”



## PERSEMBAHAN

Segala puji bagi Allah SWT yang telah mencurahkan rahmat dan karuniaNya kepada makhluk-makhlukNya. Sholawat serta salam tidak lupa kita curahkan kepada junjungan nabi besar kita Nabi Muhammad SAW yang kita nantikan syafaatnya di hari kiamat kelak.

Alhamdulillah, penulis ucapkan syukur kehadiran Allah SWT karena atas kehendakNya-lah penulis dapat menyelesaikan laporan skripsi yang berjudul **“NAMED ENTITY RECOGNITION BERITA BAHASA INDONESIA DENGAN POLYGLOT”**. Tidak lupa penulis persembahkan karya tulis ini untuk:

1. Kedua orang tua tercinta, yang senantiasa memberikan kasih sayang dan juga doa yang tak ada batasnya dan yang selalu mendidik tanpa bosannya, semoga selalu dalam keadaan sehat dan selalu berada dalam lindungan-Nya.
2. Ibu Mardhiya Hayati. yang telah membimbing hingga skripsi ini dapat diselesaikan.
3. Tuan R, Tuan E, Tuan F, Tuan D yang telah membimbing selama ini.
4. Rekan-rekan yang telah mendukung.

## KATA PENGANTAR

Puji syukur kepada Allah SWT atas rahmat dan hidayah-Nya, penulis dapat menyelesaikan skripsi yang berjudul **“NAMED ENTITY RECOGNITION BERITA BAHASA INDONESIA DENGAN POLYGLOT”** dengan lancar.

Laporan ini disusun sebagai salah satu syarat kelulusan program S1 Informatika Universitas Amikom Yogyakarta. Dalam penyusunan laporan ini penulis mendapat bantuan dari berbagai pihak. Penulis ingin mengucapkan terima kasih kepada para pihak yang telah membantuk dalam penulisan laporan skripsi ini. Maka dari itu penulis mengucapkan terima kasih kepada:

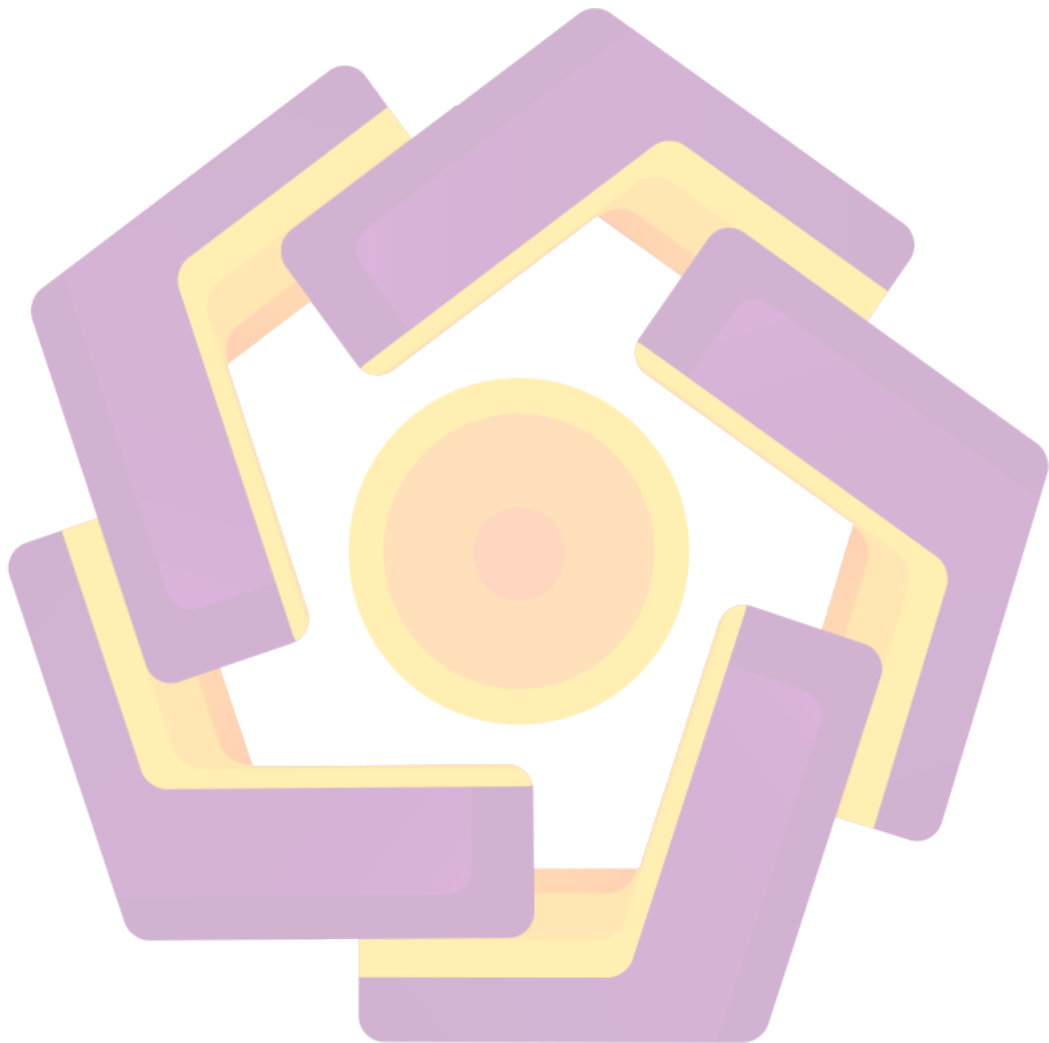
1. Bapak Prof. Dr. M. Suyanto, M.M. selaku Rektor Universitas Amikom Yogyakarta.
2. Bapak Hanif Al Fatta, M.Kom. selaku Dekan Fakultas Ilmu Komputer Universitas Amikom Yogyakarta.
3. Ibu Mardhiya Hayaty, S.T., M.Kom. selaku dosen pembimbing yang telah memberikan bimbingan dan arahan sehingga skripsi ini selesai.
4. Dewan Penguji dan segenap Dosen Universitas Amikom Yogyakarta yang telah berbagi ilmu dan pengalamannya.
5. Kedua orang tua yang selalu mendoakan, memberikan semangat dan dukungan moril.
6. Penulis sumber bacaan, jurnal dan makalah yang penulis jadikan referensi dalam penulisan laporan skripsi ini.

Penulis menyadari bahwa masih ada banyak kekurangan di dalam laporan ini. Namun penulis berharap laporan skripsi ini dapat memberikan manfaat pada para pembaca sekalian.

Tegal, 20 Agustus 2021

Barep Setiyadi  
19.21.1330





## DAFTAR ISI

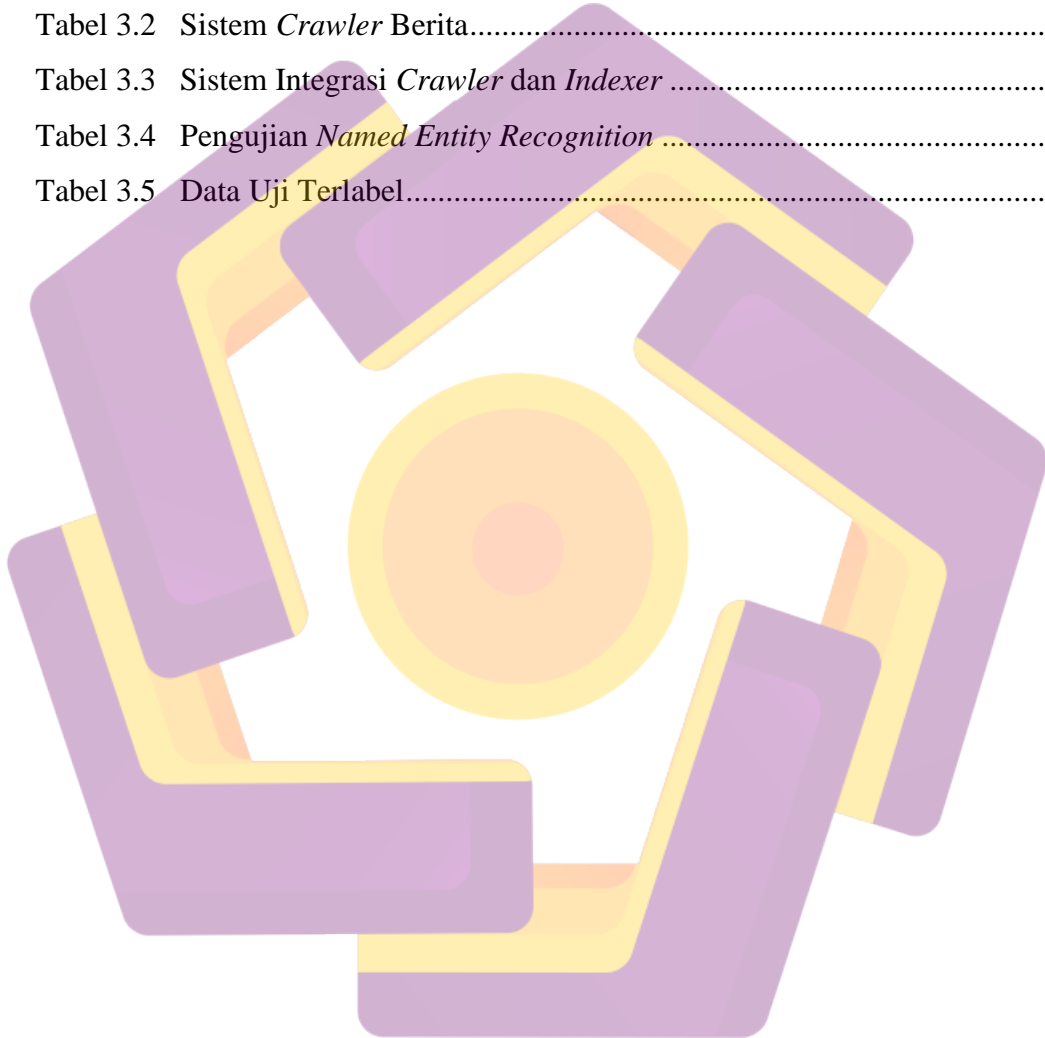
<b>JUDUL .....</b>	<b>i</b>
<b>PERSETUJUAN.....</b>	<b>ii</b>
<b>PENGESAHAN.....</b>	<b>ii</b>
<b>PERNYATAAN.....</b>	<b>iii</b>
<b>MOTTO .....</b>	<b>iv</b>
<b>PERSEMBAHAN.....</b>	<b>v</b>
<b>KATA PENGANTAR.....</b>	<b>vi</b>
<b>DAFTAR ISI.....</b>	<b>viii</b>
<b>DAFTAR TABEL .....</b>	<b>xi</b>
<b>DAFTAR GAMBAR.....</b>	<b>xii</b>
<b>INTISARI .....</b>	<b>xiv</b>
<b>ABSTRAK .....</b>	<b>xv</b>
<b>BAB 1 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 .....	2
1.5 Metode Penelitian.....	3
1.6 Sistematika Penulisan.....	4
<b>BAB II LANDASAN TEORI .....</b>	<b>5</b>
2.1 Tinjauan Pustaka .....	5
2.2 Dasar Teori.....	10

2.2.1	<i>Natural Language Processing</i> .....	10
2.2.2	<i>Named Entity Recognition</i> .....	10
2.2.3	<i>Text Mining</i> .....	14
2.2.4	<i>Crawler</i> .....	15
<b>BAB III METODE PENELITIAN .....</b>		<b>17</b>
3.1	Metode Penelitian.....	17
3.2	Metode Pengembangan Sistem .....	18
<b>BAB IV HASIL PENELITIAN DAN PEMBAHASAN.....</b>		<b>31</b>
4.1	Pengumpulan Data .....	31
4.2	Pemrosesan Data .....	32
4.2.1	Ekstraksi Berita.....	32
4.2.2	<i>Ingesting</i> Berita.....	34
4.2.3	<i>Indexing</i> Berita.....	35
4.3	Rancangan Aplikasi.....	36
4.3.1.	Polyglot NER .....	36
4.3.2.	NER REST API .....	40
4.4	Pengujian .....	42
4.4.1	Pengujian <i>Crawler</i> .....	42
4.4.2	Pengujian <i>Indexer</i> .....	43
4.4.3	Pengujian <i>Polyglot NER</i> .....	44
4.4.4	Pengujian <i>REST API</i> .....	45
<b>BAB V PENUTUP.....</b>		<b>46</b>
5.1	KESIMPULAN .....	46
5.2	SARAN .....	46



## DAFTAR TABEL

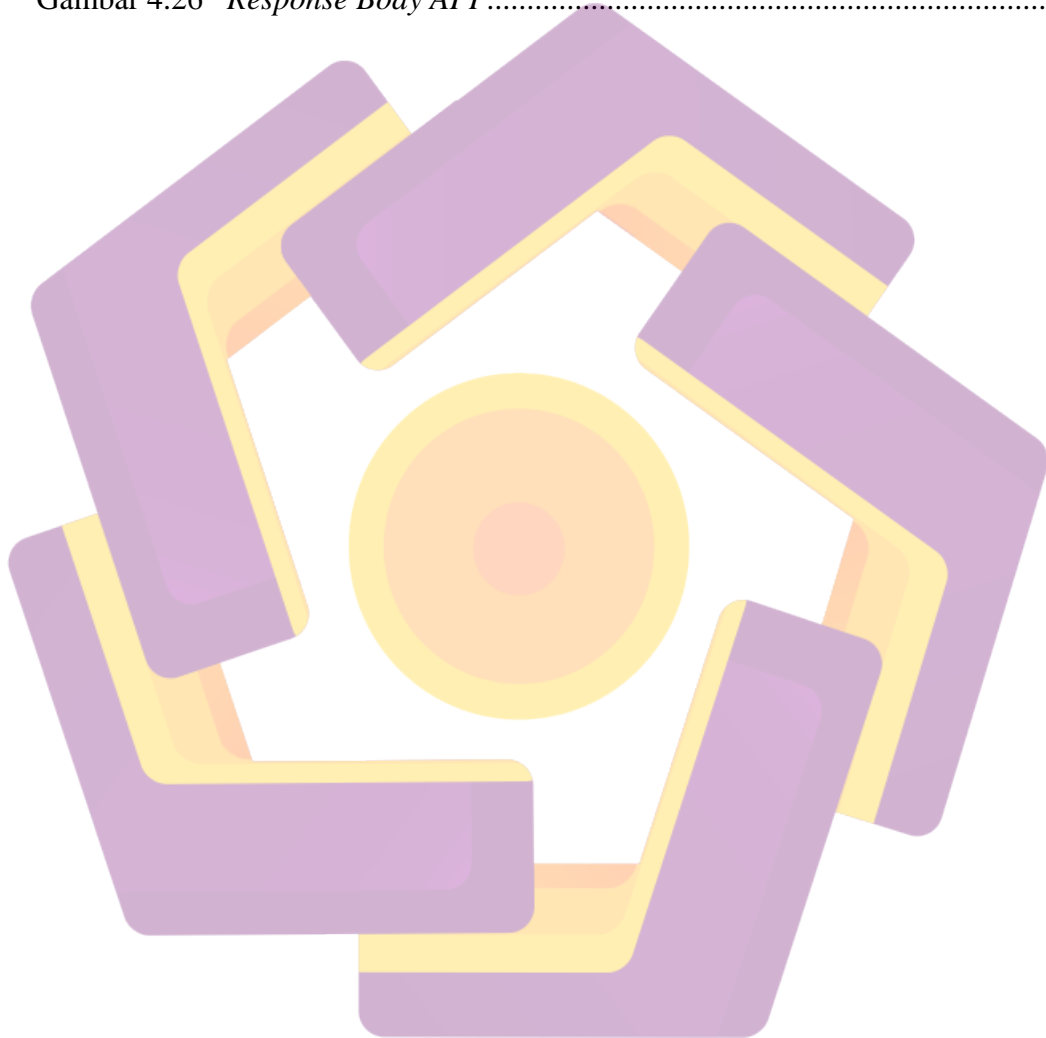
Tabel 2.1	Matrik Perbandingan Tinjauan Pustaka .....	7
Tabel 3.1	Kebutuhan Perangkat .....	21
Tabel 3.2	Sistem <i>Crawler</i> Berita.....	25
Tabel 3.3	Sistem Integrasi <i>Crawler</i> dan <i>Indexer</i> .....	26
Tabel 3.4	Pengujian <i>Named Entity Recognition</i> .....	27
Tabel 3.5	Data Uji Terlabel.....	28



## DAFTAR GAMBAR

Gambar 3.1 Tahapan Penelitian .....	17
Gambar 3.2 <i>Pipeline</i> untuk <i>web crawling</i> .....	18
Gambar 3.3 Alur Sistem.....	22
Gambar 3.4 Halaman <i>Entity Extraction</i> .....	23
Gambar 3.5 <i>Language Detection</i> .....	23
Gambar 3.6 Article Extraction .....	24
Gambar 3.7 Rancangan REST API.....	24
Gambar 3.8 Data Uji Belum Terlabel .....	28
Gambar 4.1 Kode <i>Crawler</i> .....	32
Gambar 4.2 Hasil <i>Crawling</i> berita dengan Format HTML.....	32
Gambar 4.3 Blok Fungsi untuk Penghapusan Tag.....	33
Gambar 4.4 Blok Fungsi untuk Penghapusan Spasi .....	33
Gambar 4.5 Blok Kode untuk Menyimpan ke dalam Format JSON .....	33
Gambar 4.6 Hasil Berita dengan Format JSON .....	34
Gambar 4.7 Blok Kode Koneksi Crawler dan Elasticsearch .....	34
Gambar 4.8 Pemetaan <i>Index</i> .....	35
Gambar 4.9 Daftar <i>Index</i> .....	35
Gambar 4.10 Data Berita Terindeks.....	36
Gambar 4.11 Perintah Pemasangan <i>Numpy libicu-dev</i> .....	36
Gambar 4.12 Paket dan Model Bahasa Indonesia.....	37
Gambar 4.13 Blok Kode Tokenisasi .....	38
Gambar 4.14 Teks Berita Sebelum Tokenisasi .....	38
Gambar 4.15 Teks Berita Setelah Tokenisasi .....	38
Gambar 4.16 Hasil <i>Entity Recognition</i> .....	39
Gambar 4.17 Halaman <i>Entity Extraction</i> .....	40
Gambar 4.18 Hasil Ekstraksi dari <i>API</i> .....	41
Gambar 4.19 Halaman <i>Language Detection</i> .....	41

Gambar 4.20	Halaman <i>Article Extraction</i> .....	42
Gambar 4.21	Hasil <i>Article Extraction</i> .....	42
Gambar 4.23	<i>Index Berita</i> .....	43
Gambar 4.24	Hasil Ekstraksi Berita.....	44
Gambar 4.25	Pengujian <i>REST API</i> .....	45
Gambar 4.26	<i>Response Body API</i> .....	45



## INTISARI

Dalam *Natural Language Processing(NLP)*, *Named Entity Recognition* merupakan sub-bahasan yang cukup banyak digunakan untuk penelitian. Tugas utama dari *Named Entity Recognition* yaitu membantu mengidentifikasi dan mendeteksi nama entitas dari suatu kata yang terdapat dalam kalimat.

Sumber data yang digunakan yaitu berita Bahasa Indonesia yang berasal dari 3 media *online* faktual, yaitu detik.com, kompas.com, dan tribunews.com. Kata yang terdapat pada berita Bahasa Indonesia dapat merujuk nama entitas orang, lokasi, atau organisasi. Dengan mengintegrasikan antara *crawler* berita, *Elasticsearch* sebagai alat untuk menyimpan data berita yang sudah diambil dan *Polyglot-NER* menggunakan *REST API*, maka dimungkinkan untuk membuka ruang kolaborasi antara pengembang sistem analisa media lainnya.

Saran yang disampaikan penulis adalah menerapkan *REST API* sebagai jembatan komunikasi antar sistem, diharapkan mampu mengintegrasikan antar sistem crawler berita dengan sistem *NER* Bahasa Indonesia dengan tingkat akurasi diatas 80%.

**Kata Kunci:** *Natural Language Processing, Named Entity Recognition, REST-API, Polyglot, Elasticsearch.*



## ABSTRAK

*In Natural Language Processing (NLP), Named Entity Recognition is a sub-discussion that is quite widely used for research. The main task of Named Entity Recognition is to help identify and detect the name of the entity from a word contained in a sentence.*

*The data sources used are Indonesian news originating from 3 factual online media, namely detik.com, kompas.com, and tribunews.com. Words in Indonesian news can refer to the name of the entity, person, location, or organization. By integrating news crawlers, Elasticsearch as a tool to store retrieved news data and Polyglot-NER using the REST API, it is possible to open up space for collaboration between developers of other media analysis systems.*

*The suggestion submitted by the author is to apply the REST API as a communication bridge between systems, which is expected to be able to integrate news crawler systems with the Indonesian NER system with an accuracy rate above 80%.*

**Keywords:** *Natural Language Processing, Named Entity Recognition, REST-API, Polyglot, Elasticsearch.*

