

UNIVERSITÄT BIELEFELD

BACHELORARBEIT

---

# Entity Linking through Indexing : Implementierung und Evaluierung

---

*Author:*

Tim PONTZEN

*Supervisor:*

Prof. Dr. Philipp CIMIANO

*A thesis submitted in fulfilment of the requirements  
for the degree of Bachelor of Science*

*in the*

Semantic Computing Group  
Universität Bielefeld

12. Mai 2014

# Declaration of Authorship

I, Tim PONTZEN, declare that this thesis titled, 'Entity Linking through Indexing : Implementierung und Evaluierung' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

---

Date:

---

# Abkürzungen

<b>BI</b>	<b>B</b> lock <b>I</b> ndex
<b>ERU</b>	<b>E</b> ntity- <b>R</b> ecognition- <b>U</b> nit
<b>GUI</b>	<b>G</b> enerell <b>U</b> ser <b>I</b> nterface
<b>IBEL</b>	<b>I</b> ndex <b>B</b> ased <b>E</b> ntity <b>L</b> inker
<b>IBELU</b>	<b>I</b> ndex <b>B</b> ased <b>E</b> ntity <b>L</b> inker <b>U</b> tility
<b>idf</b>	<b>i</b> nverse <b>d</b> ocument <b>f</b> requency
<b>JAR</b>	<b>J</b> ava <b>A</b> Rchieve
<b>NER</b>	<b>N</b> amed <b>E</b> ntity <b>R</b> ecognizer
<b>tf</b>	<b>t</b> erm <b>f</b> requency
<b>URI</b>	<b>U</b> niform <b>R</b> esource <b>I</b> dentifier

# Abbildungsverzeichnis

2.1	Systemübersicht . . . . .	2
2.2	GUI . . . . .	3
2.3	Entity-Extraction-Unit . . . . .	3
2.4	Linking Unit . . . . .	4
2.5	Programmablauf . . . . .	4
5.1	Index Based Entity Linker Utility . . . . .	7

# Inhaltsverzeichnis

<b>Declaration of Authorship</b>	<b>i</b>
<b>Abkürzungen</b>	<b>ii</b>
<b>List of Figures</b>	<b>iii</b>
<b>Inhaltsverzeichnis</b>	<b>iv</b>
<b>1 Einleitung</b>	<b>1</b>
<b>2 Systemübersicht</b>	<b>2</b>
2.1 GUI . . . . .	3
2.2 Entity-Recognition-Unit . . . . .	3
2.3 Linking-Unit und Index . . . . .	4
2.4 Informationsfluss . . . . .	4
<b>3 Implementierung</b>	<b>5</b>
3.1 Index . . . . .	5
3.1.1 BlockIndex . . . . .	5
<b>4 Evaluierung</b>	<b>6</b>
<b>5 Deployment und Erweiterung auf andere Sprachen</b>	<b>7</b>
5.1 Deployment . . . . .	7
5.1.1 Generierung der Indexe mit Utility-1.0 . . . . .	7
5.1.2 Eigenständige Generierung von Dateien . . . . .	9
5.2 Erweiterung um andere Sprachen . . . . .	10
5.2.1 Indexe . . . . .	10
5.2.2 Entity Recognition . . . . .	10
<b>6 Zusammenfassung</b>	<b>12</b>

# Kapitel 1

## Einleitung

## Kapitel 2

# Systemübersicht

In diesem Kapitel wird die Systemarchitektur und dessen Informationsfluss behandelt. Insgesamt besteht das System aus 4 Teilen:

- GUI
- Entity-Reconition-Unit(ERU)
- Linking Unit
- Index

Das nachfolgende Bild bietet einen Überblick über die Architektur und das Zusammenspiel der einzelnen Komponenten. Deren Funktion wird im Anschluss näher erläutert.

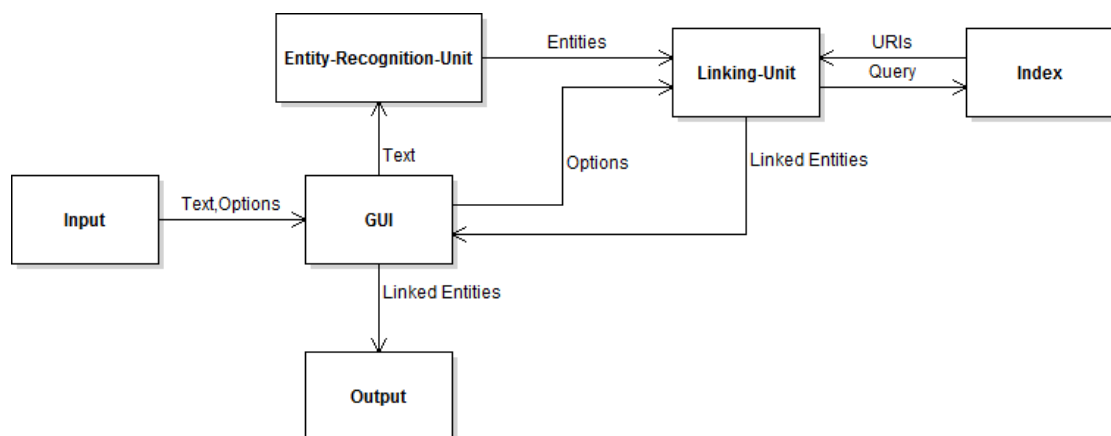


ABBILDUNG 2.1: Systemübersicht

## 2.1 GUI

Die GUI stellt die Schnittstelle zwischen Benutzer und Programm dar. Diese erlaubt ihm seinen zu anotierenden Text zu übergeben und zwischen verschiedenen Optionen zu wählen. Der Input wird entgegengenommen und dessen Text an die Entity-Recognition-Unit weitergeleitet. Am Ende bekommt es die verlinkten Entities von der Linking-Unit zurück und gibt diese an den Benutzer weiter.

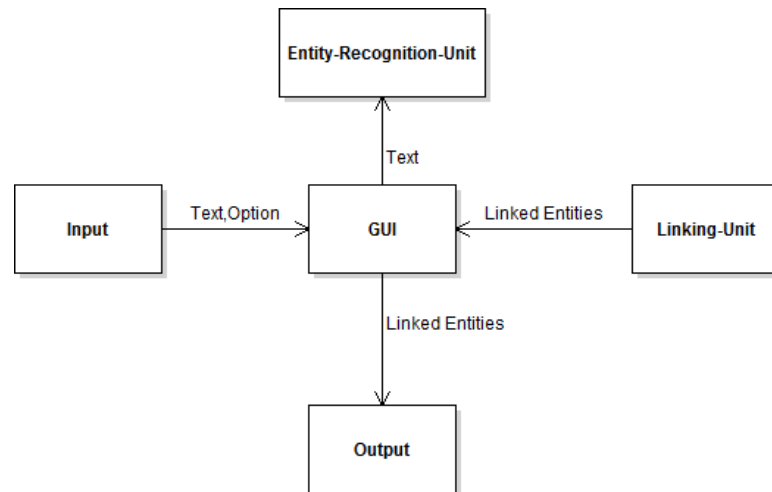


ABBILDUNG 2.2: GUI-Verbindungen

## 2.2 Entity-Recognition-Unit

Diese Einheit ist für die Erkennung von Named Entities verantwortlich. Das bedeutet, dass sie den ihr übergebenen Text analysiert und alle sich darin befindlichen Entitäten extrahiert (z.B. Personen oder Orte). Das Ergebnis dieses Prozesses wird dann an die Linking-Unit weitergeleitet.

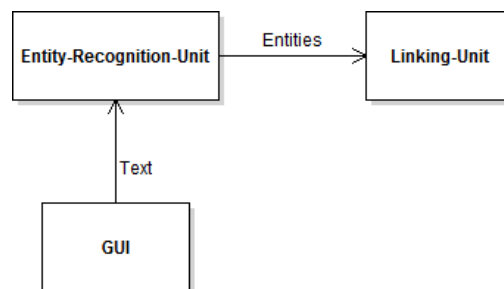


ABBILDUNG 2.3: ERU



## 2.3 Linking-Unit und Index

Die Linking-Unit versucht mit Hilfe des Indexes für jede an sie weitergeleitete Entität eine passende URI zu finden. Dazu erstellt sie, unter Berücksichtigung der vom Benutzer eingestellten Optionen, eine Query und lässt den Index diese ausführen. Dessen Rückgabewert wird dann mit den Named Entities verknüpft und an die GUI weitergereicht.

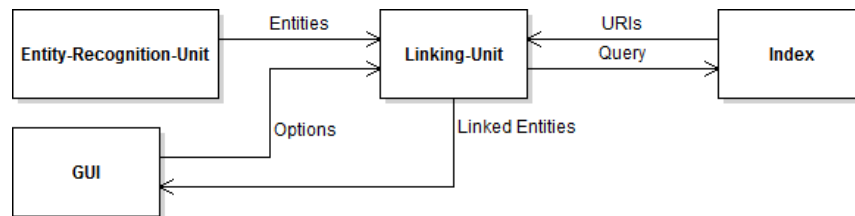


ABBILDUNG 2.4: Linking-Unit

## 2.4 Informationsfluss

bla

bla

bla

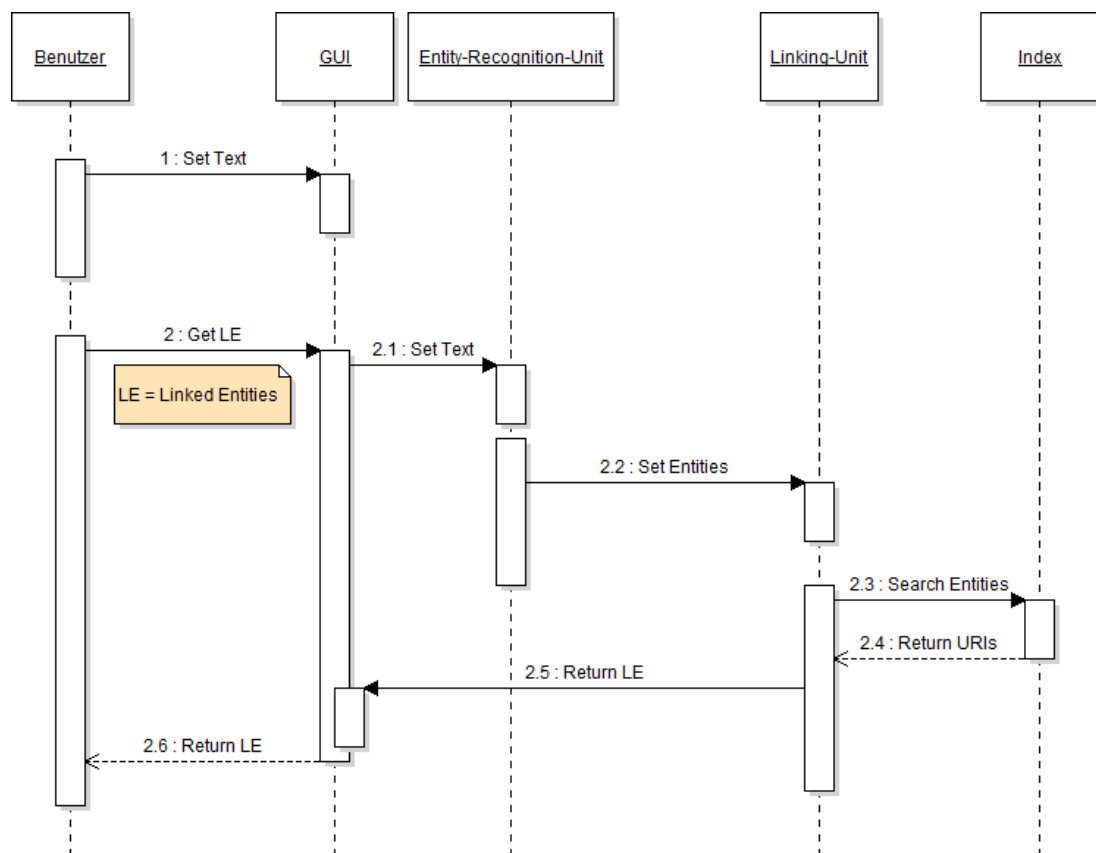


ABBILDUNG 2.5: Programmablauf

# Kapitel 3

## Implementierung

bal bla bla

### 3.1 Index

#### 3.1.1 BlockIndex

---

`<child11>`

`...`

`<childn1>`

`<parent1>`

`...`

`<childm1>`

`...`

`<child1m>`

`...`

`<childnm>`

`<parentm>`

---

`<anchorN11>`

`...`

`<anchorNn1>`

`<URI1,anchor11 ... anchori1,title1,type>`

---

## Kapitel 4

# Evaluierung

## Kapitel 5

# Deployment und Erweiterung auf andere Sprachen

### 5.1 Deployment

Um das Programm zu deployen muss lediglich die JAR mit compilierten Dependencies(falls das Programm aus dem Sourcecode compiliert wurde, ist dies die „IndexBasedEntityLinker-1.0-jar-with-dependencies.jar“), sowie ein passender Entity- und Abstract-Index in ein gleiches Verzeichnis kopiert werden.

#### 5.1.1 Generierung der Indexe mit Utility-1.0

Sollten keine gültigen Indexe vorhanden sein, so können diese entweder manuell oder über die IndexBasedEntityLinker\_Utility(IBELU) generiert werden.

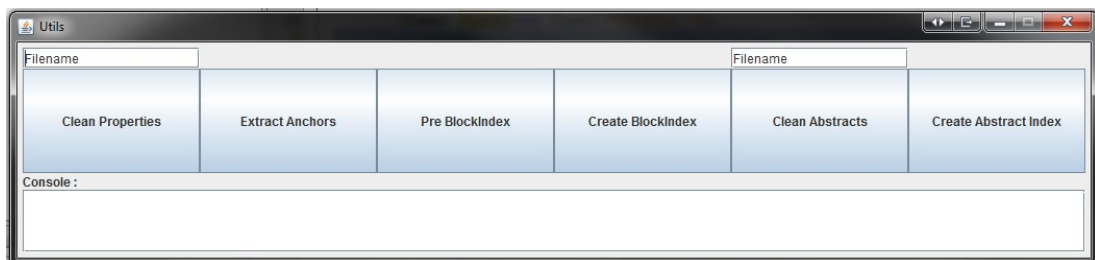


ABBILDUNG 5.1: Index Based Entity Linker Utility

Dafür werden folgende Dateien benötigt:

- Der DBpedia Anchor-Index der Computer Semantic Group der Universität Bielefeld
- Mapping-based Properties (en) von DBpedia<sup>1</sup>
- Extended Abstracts (en) von DBpedia<sup>2</sup>

Die Abfolge der auszuführenden Operationen ist von links nach rechts sortiert und sollte auch in dieser Reihenfolge abgewickelt werden.

Zur Erstellung des Entity-Indexes werden der Anchor-Index und die Mapping-based Properties benötigt. Der Ablauf ist wie folgt:

1. Clean Properties: In dem zugehörigen Textfeld oberhalb des Buttons muss der Dateipfad (lokal oder absolut) zu der Mapping-Based Properties Datei angegeben werden. Diese wird dann von allen irrelevanten Daten gesäubert.  
Output: cleaned\_properties.txt, cleaned\_properties\_neighborToEntity.txt und entities.txt.
2. Extract Anchors: Liest den Anchor-Index aus und schreibt alle <URI,anchor>-Paare in eine Datei.  
Output: anchors.txt
3. Pre BlockIndex: Erstellt für die Schnittmenge der URIs aus entities.txt und anchors.txt Dateien zur BlockIndex Generierung.  
Output: combined.txt, entity\_anchors.txt, entity\_neighbor\_anchorsN.txt, neighbor\_entity\_anchorsE.txt
4. Create BlockIndex Erzeugt den BlockIndex(EntityIndex).  
Output : Blockindex

Zu Generierung des Abstract-Indexes werden die entities.txt aus dem ersten Teil und die Long-Abstracts benötigt:

---

<sup>1</sup>[http://downloads.dbpedia.org/3.9/en/mappingbased\\_properties\\_en.nt.bz2](http://downloads.dbpedia.org/3.9/en/mappingbased_properties_en.nt.bz2)

<sup>2</sup>[http://downloads.dbpedia.org/3.9/en/long\\_abstracts\\_en.nt.bz2](http://downloads.dbpedia.org/3.9/en/long_abstracts_en.nt.bz2)

1. Clean Abstracts: Nach Angabe des Dateipfades (lokal oder absolut) der long-abstracts-Datei werden für alle URIs die ein Abstract besitzen und in entities.txt vorkommen <URI, abstract>-Paare erstellt und gespeichert.  
Output: abstract\_clean.txt
2. Create Abstract Index: Erstellt den Abstract-Index.  
Output: Abstract Index

Über Erfolg oder Misserfolg (Fehlermeldungen) der ausgeführten Operationen wird der Benutzer der IBELU über ein Konsolen-Feld (siehe Abbildung 5.1) informiert.

**Anmerkung:** Diese Operationen sind teilweise sehr Speicherintensiv. Es sollten mindestens 12GB RAM zur Verfügung gestellt werden und sichergestellt sein, dass die JVM diesen auch nutzen darf (VM Parameter: -Xmx12g).

### 5.1.2 Eigenständige Generierung von Dateien

Die IBELU kann auch nur zur reinen Indexerstellung genutzt werden, während die dafür benötigten Dateien anderweitig erstellt werden. Im Nachfolgenden werden die Dateien und ihre Formate näher erläutert.

#### Dateien zur Erstellung des Entity-Indexes

Um den Index über die IBELU zu erstellen, sind 2 Dateien erforderlich:

- combined.txt : In dieser Datei werden alle URIs auf ihre Nachbarn und deren Anchors gemappt. Ein Nachbar einer URI A ist in diesem Fall eine URI B, deren Graphknoten eine Kante zu dem Knoten von A besitzt.

Format:

---

```
URI1 | Neighbor1 | anchor1,1; ... ; anchor1,i
...
URI1 | Neighborm1 | anchorm1,1; ... ; anchorm1,j
...
...
URIn | Neighbormn | anchormn,1; ... ; anchormn,k
```

---

- `entity_anchors.txt` : Diese Datei mappt alle URIs auf ihre Anchors.

Format:

---

```
URI1 | anchor1,1; ...; anchor1,i
...
...
...
URIn | anchorn,1; ...; anchorn,j
```

---

Diese Dateien sollten lexikographisch sortiert sein, um den resultierenden Index performanter zu machen.

**Anmerkung:** Natürlich können auch die Indexe eigenständig erstellt werden. Aufbau und Funktionsweise werden in dem Kapitel „Implementierung“ ausführlich behandelt.

## 5.2 Erweiterung um andere Sprachen

Um IBEL mit anderen Sprachen nutzen zu können, müssen neue Indexe erstellt und die existierende Entity-Recognition-Unit angepasst beziehungsweise ersetzt werden.

### 5.2.1 Indexe

Für jede zu unterstützende Sprache muss ein neuer Entity-Index und ein neuer Abstract-Index erstellt werden. In Kapitel 5.1 wurde dies bereits ausführlich behandelt. Die dafür benötigten Mapping-Based Properties sowie die long-Abstracts werden von DBpedia in 119 verschiedenen Sprachen zur Verfügung gestellt.

Die entsprechenden Anchors für diese Sprache müssen allerdings entweder eigenständig generiert oder aus einer anderweitigen Quelle bezogen werden.

### 5.2.2 Entity Recognition

Für die Entity-Recognition-Unit kann, durch einen Austausch der Klassifizierer, weiterhin der Stanford-NER verwendet werden. Auf der Downloadseite des Stanford-NER<sup>3</sup> finden sich Klassifizierer für die deutsche und chinesische Sprache. Für andere Sprachen

---

<sup>3</sup><http://nlp.stanford.edu/software/CRF-NER.shtml>



müssen eigene Klassifizierer trainiert werden. Das nötige Vorgehen dafür ist auf der FAQ-Seite<sup>4</sup> beschrieben. Um den neuen Klassifizierer nutzen zu können, muss der Ladepfad in der CLASSIFIER\_PATH-Variable in der Klasse NER-Handler entsprechend angepasst werden.

Es wird auch jede andere Implementation von Named-Entity-Regocnition unterstützt. Dafür muss lediglich ein Adapeter<sup>5</sup> geschrieben werden, welcher das EntityExtractor-Interface implementiert. Dieser muss dann in der Klasse App.java der TestingGui anstelle des NER\_Handlers im Konsturktoraufruf übergeben werden.

**Anmerkung** : Da das Programm im Nachhinein in eine Serverapplikation für die Semantic Computer Group umgebaut wird, wird die TestingGui voraussichtlich nicht mehr verwendet werden. Daher muss der Adapter dann den NER\_Hanlder an einer anderen Stelle ersetzen.

---

<sup>4</sup><http://nlp.stanford.edu/software/crf-faq.shtml#a>

<sup>5</sup>[http://de.wikipedia.org/wiki/Adapter\\_\(Entwurfsmuster\)](http://de.wikipedia.org/wiki/Adapter_(Entwurfsmuster))

## Kapitel 6

# Zusammenfassung