

---

# **Optimizing a Natural Language Processing pipeline for the automatic creation of RDF data**

**Dambowy, Nils**

---



**Bachelor Thesis**

Institute for Computer Science  
Goethe University Frankfurt

supervised by:  
Dr. Karsten Tolle

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Problem description . . . . .	3
1.2	Thesis strucutre . . . . .	4
1.3	Related work . . . . .	5
<b>2</b>	<b>Background</b>	<b>6</b>
2.1	Corpus Nummorum . . . . .	6
2.2	D4N4 . . . . .	7
2.3	Natural Language Processing . . . . .	7
2.3.1	Named Entity Recognition . . . . .	7
2.3.2	Relationship Extraction . . . . .	9
2.4	Resource Description Framework . . . . .	10
2.5	RDFLib . . . . .	11
2.6	D2RQ . . . . .	12
2.6.1	D2R Mapping Language . . . . .	13
2.6.2	D2R Server . . . . .	13
<b>3</b>	<b>Assignment</b>	<b>14</b>
3.1	Overview . . . . .	14
3.2	Current state of the pipeline . . . . .	14
3.3	Implementation of the revised pipeline . . . . .	15
<b>4</b>	<b>Results</b>	<b>15</b>
4.1	Results . . . . .	15
4.2	Comparison to the previous pipeline . . . . .	15
<b>5</b>	<b>Conclusion and outlook</b>	<b>16</b>
<b>6</b>	<b>List of figures</b>	<b>17</b>
<b>7</b>	<b>Literature</b>	<b>18</b>

## **Erklärung zur Abschlussarbeit**

**gemäß § 25, Abs. 11 der Ordnung für den Bachelorstudiengang Informatik vom 06. Dezember 2010:**

Hiermit erkläre ich, dass ich die vorliegende Arbeit selbstständig und ohne Benutzung anderer als der angegebenen Quellen und Hilfsmittel verfasst habe.

Frankfurt am Main, der March 9, 2023

Dambowy, Nils

## Abbreviations

CN	Corpus Nummorum
RDF	Resource Description Framework
NLP	Natural Language Processing
CN	Corpus Nummorum
CN	Corpus Nummorum
CN	Corpus Nummorum

# 1 Introduction

Numismatics is the academic discipline focussing on the study of different forms of currency, including e.g. coins, paper money or medals. These objects can be as old as 200 years BC, modern numismatics usually begins at around 1750. Numismatists made it their task to collect, preserve and categorize these archaeological findings and thereby play an important role in the enlightenment of our history. For example when it comes to determining the date of an archaeological site, coins can be of great use since coinage could indicate a time span or the distance from its place of discovery to its origin could tell us about the mobility during these times. In our current time most numismatic institutions have to go through the challenge of finding an adequate way of managing their data and mostly<sup>1</sup> really on databases. In this thesis I am focussing on the **CN** database, which will be further explained in a later chapter. Currently working with this database is the **D4N**<sup>4</sup> (Data quality for Numismatics based on Natural language processing and Neural Networks)<sup>2</sup> project which goal it is to improve the assignment of existing data(descriptions or images) to different research portals. The assignment is done in two different ways. Firstly the assigning is done with **Natural Language Processing** which extracts conditions and description and lastly image recognition is used to recognize the coin. Furthermore the project aims to continue the development and implementation of tools useful for numismatic research portals.

## 1.1 Problem description

In this thesis I am optimizing a currently existing **Natural Language Processing** pipeline, which will be described in more detail in chapter 5.1. The pipeline consists of two parts. First **Natural Language Processing** is used to extract a name entities (see **Name Entity Relationship**) and word relationships (see **Relationship Extraction**) out of the **CN** database and afterwards written back again. Lastly, with new data and the programm **D2RQ** RDF data is created with the aim to be published later. Currently the whole process, from the execution of the notebooks to the setup of the D2RQ program has to be done manually. To improve this, the goal of this thesis is to refine the state of current pipeline by automating the process of the NLP and creation of RDF data without having to rely on the tools D2RQ offers. At the end the whole execution of the pipeline should require as little as possible human input. When given a coin description as an input it should return the results in the RDF format.

---

<sup>1</sup><http://nomisma.org/datasets>

<sup>2</sup><http://www.bigdata.uni-frankfurt.de/d4n4/>

## 1.2 Thesis strucutre

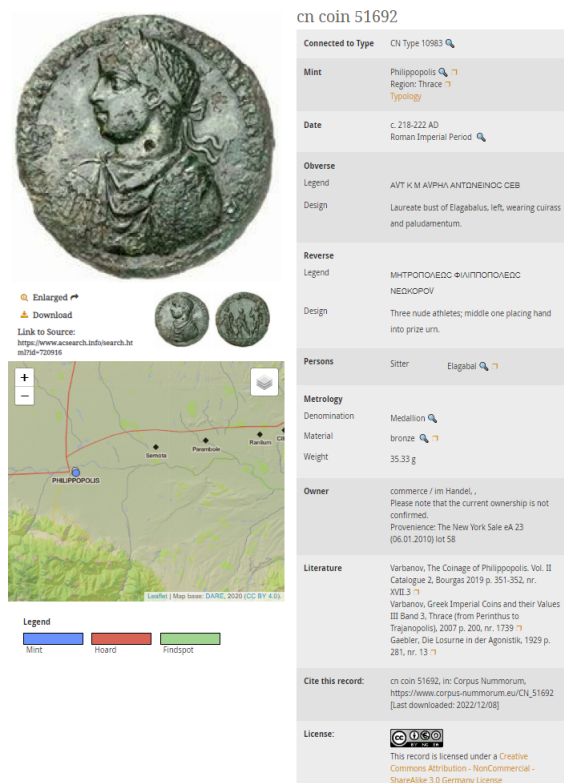
In the following chapter I am going go to discuss the CN project and its successor, the D4N4 project which is the database the pipeline is working with. Furthermore in **chapter 2** technical background information about the different technologies used is given to understand the functionality of the new and old pipeline. Afterwards, in chapter 4, the state of current pipeline and the implementation of the revised one is discussed. Both pipelines will also be compared in this chapter. Chapter 5 leaves room for the conclusion and the outlook.

### **1.3 Related work**


## 2 Background

### 2.1 Corpus Nummorum

The Corpus Nummorum (CN) is a research database, which is the result of the joint work of the Münzkabinett Berlin, Berlin-Brandenburg Academy of Sciences and Humanities (BBAW) and the Big Data Lab of Goethe University, with the motivation to offer ancient Greek coinage for research purposes.<sup>1</sup> It contains information about coins with origin in the regions of Lower Moesia, Thrace, Mysia, and the Troad. Added together, the database contains information about approx. 27,500 coins, 14,000 coming from the area of Thrace and 14,500 from the remaining regions. It is also possible to contribute coins to the database yourself.



cn coin 51692

Connected to Type	CN Type 10983
Mint	Philippopolis Region: Thrace Typology
Date	c. 218-222 AD Roman Imperial Period
Obverse	Legend: ΑΥΤ Κ Μ ΑΥΡΗΑΝΤΩΝΕΙΝΟC ΕΒ Design: Laureate bust of Elagabalus, left, wearing cuirass and paludamentum.
Reverse	Legend: ΜΗΤΡΟΠΟΛΕΩC ΦΙΛΙΠΠΟΠΟΛΕΩC ΝΕΚΡΟΠΟΛΩC Design: Three nude athletes; middle one placing hand into prize urn.
Persons	Sitter: Elagabal
Metrology	Denomination: Medallion Material: bronze Weight: 35.33 g
Owner	commerce / im Handel. Please note that the current ownership is not confirmed. Provenience: The New York Sale eA 23 (06.01.2010) lot 58
Literature	Varbanov, The Coinage of Philippopolis, Vol. II Catalogue 2, Bourgas 2019 p. 351-352, nr. XVII.3 Varbanov, Greek Imperial Coins and their Values III Band 3, Thrace (from Perinthus to Trajanopolis), 2007 p. 200, nr. 1739 Gaebler, Die Losurme in der Agonistik, 1929 p. 281, nr. 13
Cite this record:	cn coin 51692, in: Corpus Nummorum, <a href="https://www.corpus-nummorum.eu/CN/51692">https://www.corpus-nummorum.eu/CN/51692</a> [last downloaded: 2022/12/08]
License:	 This record is licensed under a Creative Commons Attribution - NonCommercial - ShareAlike 3.0 Germany License

An example of a coin in the database<sup>2</sup>

All of this information is accessible through the CN Online Website by using the provided search tool, which lets you filter coins by e.g. epoch, tribe, weight or material. Beside the mentioned characteristics, the CN offers text descriptions of the depicted images on the front- and back side of the coins. To allow a better comparison between different coins, all properties have to be entered in accord to a standardized scheme.<sup>3</sup> This not only guarantees the previous mentioned upsides but also make entering coins more accessible since e.g. volunteers do not have to worry about having to come up with a scheme themselves. The scheme offers guidelines on how to describe the obverse and reverse of the coin, figures or architecture, portraits, scenes or some general information on how properly describe a coin. After a submitting a coin, it has to reviewed before it is published.

<sup>1</sup><https://www.corpus-nummorum.eu/about>

<sup>2</sup><https://www.corpus-nummorum.eu/coins/51692>

<sup>3</sup><https://www.corpus-nummorum.eu/pdf/ExternalCoinEntry.pdf>



## 2.2 D4N4

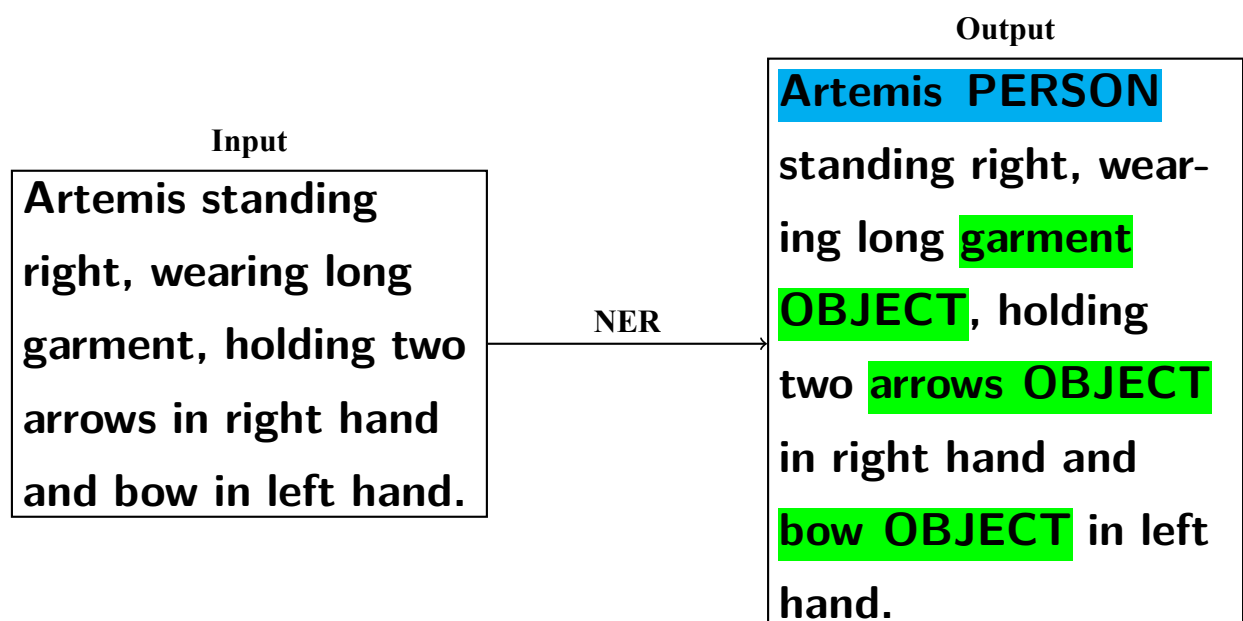
The **D4N<sup>4</sup>** is research project and a successor of the **Corpus Nummorum** project. The Name of the project is abbreviated from "Data quality for Numismatics based on Natural language processing and Neural Networks". The project aims to improve the usage of the many available images and descriptions of coins by classifying and assigning them to the different databases. Officially starting in the July of 2021, the project, much like its predecessor is the collaborative effort of the Münzkabinett Berlin, Berlin-Brandenburg Academy of Sciences and Humanities (BBAW) and the Big Data Lab of Goethe University. It achieves that by continuing and driving the development of tools forward which are needed for numismatics.

## 2.3 Natural Language Processing

NLP deals with the processing of *natural* Language. The term natural referring to the creation of the language which was natural e.g. from through human conversation unlike artificially created languages like programming languages. The language is given as input in text form for analyzation. Two different forms are being used in the pipeline and are going to be further discussed.w

### 2.3.1 Named Entity Recognition

In Name Entity Recognition natural text is being analysed for named entities. Named entities are word which give reference to a real world object using a proper noun. For Example "Goethe" or "Germany", which identify only one object. On the opposite there are words like teacher which does not refer to a specific person. In our case Named Entity Recognition could look like this:



As shown in the picture *named entity recognition* take a description of a coin and labels the named entities. Altogether there are four different labels, which are being used: **PERSON**, **OBJECT**, **ANIMAL** or **PLANT**.

### 2.3.2 Relationship Extraction

In text, words are connected over different relationships. Example: Car is a vehicle or car has wheels. With relation ship it is possible to extract this relationship. The relationship is usually between the subject and object of the sentence, thus relationship extraction is focussing on this relation.

#### Input

("Artemis standing right, wearing long garment,holding two arrows in right hand and bow in left hand.",  
"Artemis", "garment"),  
("Artemis standing right, wearing long garment,holding two arrows in right hand and bow in left hand.",  
"Artemis", "arrows"),  
("Artemis standing right, wearing long garment,holding two arrows in right hand and bow in left hand.",  
"Artemis", "bow")

#### Output

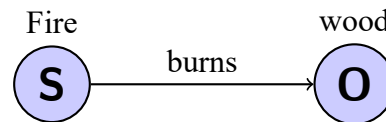
("Artemis", "wearing", "garment"),  
("Artemis", "holding", "arrows"),  
("Artemis", "holding", "bow")

## 2.4 Resource Description Framework

A **Resource Description Framework**, often abbreviated as **RDF**, is a standardized model which is most commonly used in the context of the Semantic Web. It is used to describe or exchange graph data. In an **RDF** model, data is represented as a directed graph, which consists out of triple statements.<sup>4</sup>

A triple graph statement is made out of the following components:

- A node for the *subject*
- A node for the *object*
- A *predicate* connecting the two nodes



For each of the triple statements exists such graph relationship, and altogether they make up the **RDF** model. Statements inside the **RDF** model that refer to the same subject or object are connected and form a semantic network.

Internally the components of the statement can be represented in different ways. A subject can be either a URI reference or a blank node, an object can be either a URI reference, blank node or a literal and a predicate can exist only as URI reference.<sup>4</sup> A **Uniform Resource Identifier reference**, often shortened to **URI ref**, is a Unicode string that only is made up from characters out of the ASCII Alphabet<sup>5</sup>. It is constructed out of five elements: **scheme**, **authority**, **path**, **query** and **fragment**.

Altogether, the generic URI syntax looks like this:

**foo://example.com:8042/over/there?name=ferret#nose**<sup>6</sup>

This shows the similarity to the more popular **URL**, which is a special case of an URI. A **blank node** is a node that is neither a URI ref nor a literal, that's why it is also called an anonymous resource. It contains no information. A **literal** is used to represent values such as numbers, text strings or dates.

---

<sup>4</sup><https://www.w3.org/TR/rdf-concepts/>

<sup>5</sup><https://www.ietf.org/rfc/rfc2396.txt>

<sup>6</sup>Example taken from: [https://en.wikipedia.org/wiki/Uniform\\_Resource\\_Identifier](https://en.wikipedia.org/wiki/Uniform_Resource_Identifier)

## 2.5 RDFLib

RDFLib is an open-source<sup>3</sup> python package created for working with RDF. It allows the user to conveniently add information to a RDF graph and offers to different parsers/serializers to work with. It was initially created in 2002 and is still being maintained/updated with the latest major version being published in 2021.

### Example:<sup>4</sup>

```
1 from rdflib import Graph
2 from rdflib.namespace._XSD import XSD
3
4 g = Graph()
5
6 g.add((
7     URIRef("http://example.com/person/nick"),
8     FOAF.givenName,
9     Literal("Nick", datatype=XSD.string)
10 ))
11
12 g.serialize(format="turtle")
```

In this example the structure *Graph* is being introduced, it functions as the primary interface/container when working with RDFLib. It contains all of our triples and allows to perform common set-operations(like `add()`, to add triples). In this example we are adding one triple to the graph. The triple is made out of a *URI-reference*, the property *FOAF.givenName* and a *Literal*. Afterwards the data can be serialized using the different serializers. In the example turtle was chosen as serializer.

### Output:

```
1 @prefix ns1: <http://xmlns.com/foaf/0.1/> .
2 @prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
3
4 <http://example.com/person/nick> ns1:givenName "Nick"^^xsd:string .
```

In the script, used in our pipeline, only the functions of the example above were used, however N-triples was used as a serializer.

---

<sup>3</sup><https://github.com/RDFLib/rdflib>

<sup>4</sup>Example taken from the RDFLib documentation

## 2.6 D2RQ

To create RDF graphs from the **D4N<sup>4</sup>** database used in the pipeline, D2RQ is used. D2RQ is an abbreviation for *Database to RDF Query* and is a open source mapping tool which allows you to access your data as a RDF graph without having to store it in an RDF format in your database. It consists out of:

- the D2RQ Mapping Language,
- the D2RQ Engine,
- and D2R Server.

However in the current pipeline only the D2RQ Mapping Language and D2R Server are used. In order to create RDF graphs with D2RQ, firstly a mapping file has to be created.<sup>5</sup> The D2RQ Mapping Language will be further discussed in the following chapter. With the mapping file created, it is now possible for us to view the data via d2r-server.



With the d2r-server running locally, we are now able to view our data via either HTML or using a semantic web browser. More importantly we are able to generate RDF dumps via the command line.

### An example of a coin in the database<sup>6</sup>

```
1
2 # Generate a mapping file
3 generate-mapping -o mapping.ttl -d driver.class.name -u db-user -p db-
  password jdbc:url:...
4
5 # Start D2R Server
```

<sup>5</sup><http://d2rq.org/getting-started>

```

6 d2r-server mapping.ttl
7
8 # Generate an RDF dump
9 dump-rdf mapping.ttl -o dump.nt

```

### 2.6.1 D2R Mapping Language

The D2RQ mapping language is a declarative language used to map the data from the **D4N4** database to RDFS vocabularies. It is responsible for how the virtual RDF graph looks.

#### Example:

```

1 map:coins a d2rq:ClassMap;
2   d2rq:dataStorage map:database;
3   d2rq:uriPattern "https://www.corpus-nummorum.eu/coins/@@data_coins.id@@";
4   d2rq:class nmo:NumismaticObject;
5   d2rq:condition "data_coins.publication_state = 1";
6   .

```

In the example we are creating the 'd2rq:ClassMap' *coins*, which represents one or multiple classes in the D2RQ Mapping Language. Instances of the class are being handled in accord to the definition in 'd2rq:ClassMap' by the D2RQ Mapping Language. In the next line reference to where the data is stored is given by mapping the database to 'd2rq:dataStorage'. In this case database is a variable defined earlier which represents a connection to the database. After this a reference to the resource is given by passing an URI pattern and an RDFS Class is being created. Every instance created from the ClassMap will also be an instance of this class. Lastly we are limiting the mapped coins to those with the value set accordingly, it functions as a SQL WHERE statement.

### 2.6.2 D2R Server

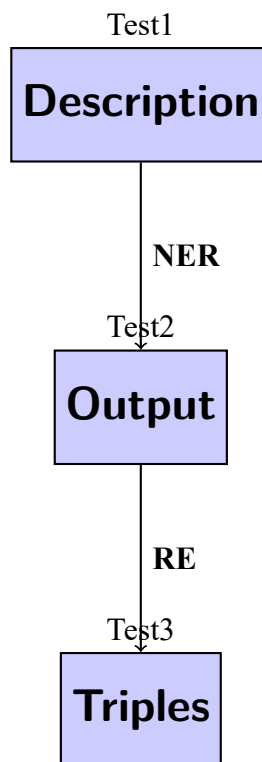
Ihr Text sollte unter Einbindung von Grafiken und Tabellen in Absätze gegliedert werden. Dabei ist zu beachten, dass ein Absatz einen thematischen Gedanken erfasst, wobei am Anfang des Absatzes im Regelfall die Kernaussage zu finden ist und von dieser ausgehend durch weitere Erörterungen innerhalb des Absatzes gegliedert wird

### 3 Assignment

#### 3.1 Overview

#### 3.2 Current state of the pipeline

At the current state the steps of executing the pipeline are associated with a lot of manual work, starting from the manual execution of the NLP notebooks to the running from the D2RQ application. At the beginning of the pipeline the **NER** is being executed. The results of the **NER** builds the foundation on that the **RE** is build. In the notebook **NER** is being applied to every design of the database and the results are then uploaded again to the database. Afterwards, **RE** follows. As previously mentioned, **RE** takes as input a design, a subject and a object and tries to create a relationship between the subject and the object. With completion of the **RE** the created triples are uploaded to the database. The last step of the pipeline is the one, which requires the most amount of manual work and thereby offers the most room of improvement. In order to created the **RDF** data, the program **D2RQ**-server is launched. As mentioned, **D2RQ**-server builds the connection between the database and the **D2RQ** program. It allows us to apply a mapping, which will map the contents of our database to **RDF**. Lastly the execution of the mapping programm has to be done manually via the command line and the results is our data represented in RDF. Summarized the pipeline looks like this:





### **3.3 Implementation of the revised pipeline**

## **4 Results**

### **4.1 Results**

### **4.2 Comparison to the previous pipeline**

## **5 Conclusion and outlook**

## **6 List of figures**

## 7 Literature

### References

Sarkar, Dipanjan (2016). *Text Analytics with Python A Practical Real-World Approach to Gaining Actionable Insights from Your Data*. Apress.