

Universitatea "Alexandru Ioan Cuza" din Iași
Facultatea de Informatică



Named Entity Recognition

For Romanian Language

Drumea Alexandru-Daniel

Sesiunea: Iulie, 2019

Coordonator științific: Prof.Dr. Cristea Dan

Contents

1	Introduction	5
1.1	Abstract	5
1.2	Motivation	7
2	Theoretical Considerations, State-of-the-Art Technologies and Data Formats	8
2.1	Data Formats Used in the Current Thesis	8
2.2	Deploying Web Applications Using Python	10
2.3	The Stanford Named Entity Recognition Software (Charles Sutton <i>et al</i> , 2015)	11
2.3.1	The Conditional Random Field Algorithm, Theoretical Considerations of the Stanford Software	11
2.3.2	The Operation of the Stanford Named Entity Recognition Software	11
2.4	The SpaCy Entity Recognition Module	12
2.5	The Theoretical Foundations of the Gazetteer Method in Named Entity Recognition	13
2.5.1	The Brute Force Method	13
2.5.2	Improving the Brute Force Method	14
2.6	Issues of the Current State-of-the-Art Methods with Respect to The Romanian Language	14
3	The Romanian Named Entity Recognition Software	15
3.1	Data Preprocessing	15
3.2	The Stanford Named Entity Recognition Method	17
3.3	The Practical Application of the Gazetteer Method	21
3.3.1	Creation of the List of Entities	21
3.3.2	The Operation of the Gazetteer Method	23
3.4	The SpaCy Named Entity Recogniser using the Multi-language Model	25
3.5	Statistics	26
3.6	The Voting Function	29
3.7	The Web Application	30
4	Conclusions	32
5	Bibliography	34

Declarație privind originalitate și respectarea drepturilor de autor

Prin prezenta declar că Lucrarea de licență cu titlul "*Titlul complet al lucrării*" este scrisă de mine și nu a mai fost prezentată niciodată la o altă facultate sau instituție de învățământ superior din țară sau străinătate. De asemenea, declar că toate sursele utilizate, inclusiv cele preluate de pe Internet, sunt indicate în lucrare, cu respectarea regulilor de evitare a plagiatului:

- toate fragmentele de text reproduse exact, chiar și în traducere proprie din altă limbă, sunt scrise între ghilimele și dețin referința precisă a sursei;
- reformularea în cuvinte proprii a textelor scrise de către alți autori deține referința precisă;
- codul sursă, imagini etc. preluate din proiecte *open-source* sau alte surse sunt utilizate cu respectarea drepturilor de autor și dețin referințe precise;
- rezumarea ideilor altor autori precizează referința precisă la textul original.

Iași, *data*

Declarație de consimțământ

Prin prezenta declar că sunt de acord ca Lucrarea de licență cu titlul "*Titlul complet al lucrării*", codul sursă al programelor și celelalte continuturi (grafice, multimedia, date de test etc.) care însoțesc această lucrare să fie utilizate în cadrul Facultății de Informatică. De asemenea, sunt de acord ca Facultatea de Informatică de la Universitatea Alexandru Ioan Cuza Iași să utilizeze, modifice, reproducă și să distribuie în scopuri necomerciale programele-calculator, format executabil și sursă, realizate de mine în cadrul prezentei lucrări de licență.

Iași, *data*

Absolvent *Prenume Nume*

(semnătura în original)

Acord privind proprietatea dreptului de autor

Facultatea de Informatică este de acord ca drepturile de autor asupra programele-calculator, format executabil și sursă, să aparțină autorului prezentei lucrări, *Prenume Nume*.

Încheierea acestui acord este necesară din următoarele motive:

[Se explică de ce este necesar un acord, se descriu originile resurselor utilizate în realizarea produsului-program (personal, tehnologii, fonduri) și aportul adus de fiecare resursă.]

lași, *data*

Decan *Prenume Nume*

Absolvent *Prenume Nume*

(semnătura în original)

(semnătura în original)

1 Introduction

1.1 Abstract

Computer science is the study of processes that interact with large volumes of data, that can be represented in the form of programs. It enables the use of algorithms to manipulate, compute, store and communicate digital information. The above mentioned field is relatively new in the scientific community, as it employs the use of computers in order to ease human labour and aid in completing activities which otherwise would take significantly longer time to complete.

Such activities include natural language processing, a topic of research which is concerned with the interactions between computers and human (natural) language, in particular how to program computers to process and analyze large quantities of natural language data. One of the challenges in this field is “named entity recognition”, a topic of great research in modern times.

Named entity recognition is a subtask of information extraction that seeks to locate and classify mentions of known named entities in unstructured texts into pre-defined categories, such as person names, organization names, geographical locations and miscellaneous such as date and time, geo-political relations etc.

Presently, conferences such as CoNLL address the issues and challenges of the task of classifying named entities and award prizes for papers that have most contributed to perfecting human language analysis in general and named entity recognition in particular.

The scientific community in this field has currently reached a level of great precision and recall for highly structured languages such as English and German, which is a much appreciated baseline for other multi-language platforms on which to build custom systems. The aforementioned structured languages have been studied to great extent in order to achieve accurate representations of named entities.

Techopedia, the leading online technical dictionary, explains “Named-entity recognition is a state-of-the-art intelligence system that works with nearly the efficiency of a human brain. The system is structured in such a way that it is capable of finding entity elements from raw data and can determine the category in which the element belongs. The system reads the sentence and highlights the important entity elements in the text. NER might be given separate sensitive entities depending on the project. This means that the NER system designed for one project may not be reused for another task. Similarly, NER faces many challenges which include the

extraction of correct information for specific but closely related categories”. In addition to those previously stated, the problem of recognizing named entities may be broken down, conceptually, in two problems, which are the detection of names and classification of said names by the type of entity they refer to (e.g. person, organization, location and miscellaneous).

In order to ~~understand~~ the results of such systems, one must define the concepts of precision, recall and F-measure functions, as defined by academic conferences such as CoNLL.

Precision is the number of predicted entity name spans that line up exactly with spans in the gold standard evaluation data. I.e. when [Person Hans] [Person Blick] is predicted but [Person Hans Blick] was required, precision for the predicted name is zero. Precision is then averaged over all predicted entity names, ~~recall is similarly~~ the number of names in the gold standard that appear at exactly the same location in the predictions, while the F-measure ~~of F1 functions~~ is the harmonic mean of these two.

The technology currently available for this type of information extraction is comprised of Python libraries and Java applications, which offer support for the aforementioned languages (German and English), with great F-measure scores. To describe the available resources for said languages, one needs to understand the Conditional Random Fields Classifiers, which are described Sutton and McCallum (2010) paper, “An Introduction to Conditional Random Fields” as follows “A solution(n.r Named ent~~ity~~ recognition) to this problem is to model the conditional distribution $p(\mathbf{y}|\mathbf{x})$ directly, which is all that is needed for classification. This is a conditional random field (CRF). CRFs are essentially a way of combining the advantages of classification and graphical modeling, combining the ability to compactly model multivariate data with the ability to leverage a large number of input features for prediction. The advantage to a conditional model is that dependencies that involve only variables in \mathbf{x} play no role in the conditional model, so that an accurate conditional model can have much simpler structure than a joint model”. **In addition, one must define the meaning of the word “corpus” in the context of NER. A corpus is a large data set consisting of bodies of text, contextually linked, labeled in a way that is understandable for a NER system.**

Stanford NER is a Java implementation of a Named Entity Recognizer. Stanford NER is also known as ~~CRF~~ Classifier. The software provides a general implementation of the linear chain Conditional Random Field (~~Charles~~ Sutton *et al*, 2011) sequence models. That is, by training ones own models on labeled data, one can actually use this code to build sequence models for NER or any other task. The English corpus is included in the Stanford NER package and will correctly tag above 90% of named entities ~~correctly~~.

NLTK (Natural Language Toolkit) is a Python package providing a series of natural language corpora and APIs of wide varieties of NLP algorithms and it works in three stages: word tokenization, part-of-speech tagging and named entity recognition.

SpaCy is a natural language processing library in Python, known for its industrial-strength classification capabilities. SpaCy supports numerous entity types as well as multi language support to create models for other languages which are not yet implemented.

All of the technological solutions above are benchmarks for either the independent, hobby researcher or for the international community restlessly working to improve the processing of human language with the use of computers. However, all of the aforementioned APIs are designed primarily for the languages listed above, despite the fact that some offer multi-language and language-independent support.

1.2 Motivation

The writer of this paper is highly motivated by the problem of creating a model NER for the Romanian language, as it may serve as a template for other languages with Latin elements, such as French, Spanish, Portuguese etc. The steps taken to achieve the task at hand include the use of the aforementioned technologies concurrently in order to create a web application which provides one with statistics of the training process, pre-trained corpora as well as a classifier for unlabeled texts.

NER models are used in a wide range of applications in the field of Natural Language processing and Information Retrieval. The most eloquent examples are as follows: automatically summarizing resumes, which consists of summarizing a resume(or CV) enabling departments such as Human Resources to compile large volumes of data at a quick glance, optimizing search engine algorithms by, instead of searching for an entered query across the millions of articles and websites online, running a NER model on the articles once and storing the entities associated with them permanently. Another application of NER is simplifying customer support as it can be used in recognizing relevant entities in customer complaints and feedback such as Product specifications, department or company branch details, so that the feedback is classified accordingly and forwarded to the appropriate department responsible for the identified product.


2 Theoretical Considerations, State-of-the-Art Technologies and Data Formats

2.1 Data Formats Used in the Current Thesis

This chapter describes a series of theoretical notations and the meta-data used in the process of named entity recognition. The files used contain specialized annotations used by the compiler in the categorization of the entities and prediction of their type.

The CoNLL-U Plus Format can encode any kind of annotation using a combination of the sentence-level comments and the MISC attributes, this format can have any number of non-zero columns, as well as sentence-level comment (depicted by the use of # at the beginning of the line). This extends the CoNLL-U format which is widely used in Natural Language Processing, but is impractical in Named Entity Recognition due to its fixed number of columns (ten).

QuoVadis is a corpus displaying semantic relations in free text. As stated in the paper, "Developing a lexical-semantic knowledge base to be used in Natural Language Processing (NLP) applications is the main goal of the research described in this chapter (n.r. paper). Such resources are not available for many languages, mainly because of the high cost of their construction. The knowledge base is built in such a way as to facilitate the training of programs aiming to automatically recognize in text entities and semantic relations.[...] it includes annotations for the spans of text that display mentions of entities and relations, for the arguments (poles) of the relations, as well as for the relevant words or expressions that signal relations" (D. Cristea *et al*, 2014). The data that comprises this corpus will be used in the training of the NER tagger this paper intends on building. The layers of annotation in the aforementioned corpus are as follows: segmentation at text level (marks the sentence boundaries in the raw book text), tokenization (demarcates words or word compounds, but also numbers, punctuation marks, abbreviations etc), part-of-speech tagging, lemmatization (determines lemmas of words) and noun-phrase chunking, noun phrase chunking (explores the previously generated data and adds information regarding noun phrase boundaries and their head words). In the scope of training the NER, this application uses the part-of-speech tagging provided by the QuoVadis corpus, as well as the noun phrase chunking to create a tab separated values(.tsv) file to be trained by the Stanford CRF classifier. The data of the aforementioned corpus, which is initially comprized in an .xml type file, will be extracted and each "entity" tag will be added to the .tsv file, as highlighted in Figure 1.




```

<S id="3" offset="82">
  <W lemma="in" MSD="Sp" POS="ADPOSITION" deprel="o.c.t." head="3" id="1" offset="0">in</W>
  <W Case="direct" Definiteness="no" Gender="masculine" lemma="ajun" MSD="Noun" Number="singular" POS="NOUN" Type="common" deprel="prep." head="1" id="2" offset="3">ajun</W>
  <W EXTRA="interactiv" lemma="fi" MSD="Vml3s" Mood="indicative" Number="singular" POS="Verb" Person="third" Tense="long" Type="predicative" deprel="o.c.t." head="1" id="3" offset="8">fusesce</W>
  <W lemma="la" MSD="Sp" POS="ADPOSITION" deprel="o.c.t." head="3" id="4" offset="15">la</W>
  <W lemma="petrecere" MSD="Noun" Number="singular" POS="NOUN" Type="common" deprel="prep." head="4" id="5" offset="18">petrecere</W>
  <ENTITY id="R000300001" TYPE="PERSON">
    <W Case="oblique" Definiteness="no" EXTRA="NotIndict" Gender="feminine" lemma="Nero" MSD="Noun" Number="plural" POS="NOUN" Type="proper" deprel="prep." head="4" id="6" offset="18">Nero</W>
  </ENTITY>
  <W lemma="la" MSD="Sp" POS="ADPOSITION" deprel="o.c.t." head="3" id="6" offset="23">la</W>
  <W Case="direct" Gender="feminine" lemma="un" MSD="Tifer" Number="singular" POS="ARTICLE" Type="indefinite" deprel="det." head="6" id="7" offset="27">o</W>
  <CLOSE CONTINUE="1" is="CLAUZE">
    <W Case="direct" Definiteness="no" Gender="feminine" lemma="petrecere" MSD="Noun" Number="singular" POS="NOUN" Type="common" deprel="prep." head="6" id="8" offset="29">petrecere</W>
  </CLOSE>

```

Figure 1: Data of the QuoVadis corpus



RONEC(Stefan Dumitrescu *et al*, 2018) is a corpus of annotated data consisting of 5127 sentences, classified in 16 classes with a total of 26376 entities. It is used as a template for the training files, as well as for the named entity data it contains. From this file, the classifier will use the template of the file (CoNLL-UP format), as well as the data classified as person, organization, geographical location, or miscellaneous. [The RONEC.conllup file is depicted in figure 2.](#)


```

# sent_id = 5
# text = Flankerul selectionatei de rugby a României și al echipei franceze din Grenoble, Florin Corodeanu, și-a reluat antrenamen!
1 Flankerul Flankerul PROPON Ncmary 19 nsubj 1:PERSON
2 selectionatei selectionată NOUN Ncfsay Case=Dat,Gen|Definite=Def|Gender=Fem|Number=Sing 1 nmod _ _ *
3 de de ADP Spas AdpType=Prep|Case=Acc 4 case _ _ _
4 rugby rugby NOUN Ncms-n Definite=Ind|Gender=Masc|Number=Sing 2 nmod _ _ *
5 a al DET Tsfs Gender=Fem|Number=Sing|Poss=Yes|PronType=Prs 6 det _ _ _
6 României României PROPON Npfsay Case=Dat,Gen|Definite=Def|Gender=Fem|Number=Sing 4 nmod _ _ 2:GPE
7 și și CONJ Crssp Polarity=Pos 9 cc _ _ _
8 al al DET Tsms Gender=Masc|Number=Sing|Poss=Yes|PronType=Prs 9 det _ _ _
9 echipei echipă NOUN Ncfsay Case=Dat,Gen|Definite=Def|Gender=Fem|Number=Sing 2 conj _ _ _
10 franceze francez ADJ Afpfson Case=Dat,Gen|Definite=Ind|Degree=Pos|Gender=Fem|Number=Sing 9 amod _ _ 3:NAT_REL_POL
11 din din ADP Spas AdpType=Prep|Case=Acc 12 case _ _ _
12 Grenoble Grenoble PROPON Np 9 nmod _ _ SpaceAfter=No 4:GPE
13 , , PUNCT COMMA 14 punct _ _ _
14 Florin Florin PROPON Np 1 appos _ _ 5:PERSON

```

Figure 2: Data of the RONEC corpus

After the processing of the data in the two aforementioned corpora, the custom file to be used in the process of training, will consist of roughly 195000 entries, with their named entity recognition label, which is detailed below.



```

intalnirea O
Bucuresti GPE
Tudor PERSON
Argezi PERSON
Vardar ORGANIZATION
Skopje ORGANIZATION

```

Figure 3: Data of the processed corpus

In Figure 3, one can observe the form of the processed corpus is of the form [entity, type separated by a tab](#). The types of entities are: person, organization, o (others, which describes an entity which is not a named one),

GPE which depicts countries, cities, states, as well as MISC, which will cover geo-politic entities, date and time, facilities and works of art.

2.2 Deploying Web Applications Using Python

As the result of the application of the theoretical considerations and algorithms presented in this thesis will be comprised in a web application, it is very important to note some of the frameworks and libraries to be used.

One of the applications which provide one with the best ease of use is called Flask.

Flask is a microframework for Python, that comes with built-in development, server and debugger, integrated unit testing support, RESTful request dispatching, is Unicode based, and is extensively documented.

A use-case of this framework is depicted in the figure below.

```
from flask import Flask
from flask import Flask, request, render_template
import os
import Vote

app = Flask(__name__)
PEOPLE_FOLDER = os.path.join('static', 'images')
app.config['UPLOAD_FOLDER'] = PEOPLE_FOLDER

@app.route('/')
def home():
    menu = os.path.join(app.config['UPLOAD_FOLDER'], 'menu.png')
    index = os.path.join(app.config['UPLOAD_FOLDER'], 'index.png')
    sect1 = os.path.join(app.config['UPLOAD_FOLDER'], 'section1.png')
    footer = os.path.join(app.config['UPLOAD_FOLDER'], 'footer.png')
    submit = os.path.join(app.config['UPLOAD_FOLDER'], 'submit.png')
    if len(request.args)!=0:
        string = "<p>" + str(Vote.vote(request.args['sentence'])) + "</p>"
        return render_template('index.html', menu_img = menu, index_img = index, sect1_img = sect1, footer_img = footer, submit_img = submit, data = string)
    return render_template('index.html', menu_img = menu, index_img = index, sect1_img = sect1, footer_img = footer, submit_img = submit)

if __name__ == '__main__':
    app.run(debug=True)
```

The aforementioned use-case is as follows: a page that consists of three elements, a navigation bar, contents and the footer respectively. To be noted is that the page will present in the content section of its body a text-box in which the user will insert a sentence that will be sent to the server, the server catching that will execute a Python function using the parameters of the request. Finally, the server will use the results of said Python function and will alter the front-end accordingly.

2.3 The Stanford Named Entity Recognition Software (~~Charles~~ Sutton *et al*, 2015)

2.3.1 The Conditional Random Field Algorithm, Theoretical Considerations of the Stanford Software

It is denoted as $\mathbf{x} = (x_1, \dots, x_n)$ as the input sequence i.e. the words of a sentence and $\mathbf{s} = (s_1, \dots, s_n)$ as the sequence of output states i.e. the named entity tags. In conditional random fields, the conditional probability $p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{s}_1, \dots, \mathbf{s}_n)$ is modeled. This feat is accomplished by defining a feature map function $\Phi(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}_1, \dots, \mathbf{s}_n) \in \mathbb{R}^d$, this maps the input sequence \mathbf{x} paired with an entire state sequence \mathbf{s} to some d-dimensional vector. Then one can model the probability as a log-linear model with the parameter vector $\omega \in \mathbb{R}^d$:

$$p(s|x; \omega) = \frac{\exp(\omega \cdot \Phi(x, s))}{\sum_{s'} \exp(\omega \cdot \Phi(x, s'))},$$

where s' ranges over all possible output sequences. For the estimation of ω , it is assumed there is a set of n labeled examples (x^i, s^i) , where i ranges between 1 and n . Now, the log-likelihood function is computed:

$$\sum_{i=1}^n \log(p(s^i | x^i, \omega))$$

The parameter vector ω^* is then estimated as:

$$\omega^* = \operatorname{argmax}_{\omega \in \mathbb{R}^d} L(\omega)$$

Now, the most likely tagging of a sentence x may be computed as:

$$s^* = \operatorname{argmax}_s p(s|x, \omega^*)$$

Having understood the aforementioned algorithm, one may apply said algorithm both for part-of-speech tagging, as well as named entity recognition. ~~Having said all of the above,~~ this algorithm is most prominently used by the Stanford NLP team as the nucleus for their named entity recognition software, on which this thesis has based its research into NER tagging in the Romanian language.

2.3.2 The Operation of the Stanford Named Entity Recognition Software

This model uses the Conditional Random Field Algorithm in order to predict the features of an unknown word. Of high importance is the fact

that the features of a word are more important than the model itself, as the latter may be obtained by combining the words and their features.

All word features are stored in a window, that being comprised of the current word, previous word, next word, as well as orthographic features such as capitalization and overall form of a word (Jenny is generically depicted as Xxxxx, while AH1-N1 as XX#-X#). Additionally, features of high importance are also stored, features such as prefixes as suffixes, which are also stored in a window (Jenny is depicted as <Jen, <...ny> >), as well as label sequences as feature conjunctions.

This method of recognition also employs the use of the Distributional Similarity Features on a large unannotated corpus, which will provide context of named entities, which will then be clustered on how similar their distributions are within the corpus in order to combat sparsity.

The Stanford Named Entity Recognition Software is highly used in the academical due to its multi-language and language-independent support, as well as to its speed and its open-source availability.

2.4 The SpaCy Entity Recognition Module

SpaCy is a free, open source library that has made Natural Language Processing (NLP) much simpler in Python.

It provides users with a statistical system which is highly efficient for named entity recognition, which can assign labels to groups of tokens which are contiguous. It can accurately recognize a wide variety of named entities or numerical entities, as well as enabling the addition of arbitrary classes to the model, by training it to enable bespoke examples.

Due to the fact that the model is designed to provide a mixture of speed as well as good F-measure scores, the company developing it is yet to provide the scientific community with a detailed architecture of the model, the only technical detail being that it uses convolutional neural networks (CNN).

A convolutional neural network is a class of deep neural networks, which are used in a wide variety of tasks in deep learning, ranging from analyzing visual imagery, to sentence classification and named entity recognition.

"Convolutional neural networks (CNN) utilize layers with convolving filters that are applied to local features" (Nal Kalchbrenner *et al*, 2014). Originally invented for computer vision, CNN models have subsequently been shown to be effective for NLP and have achieved excellent results in semantic parsing (Wen-tau Yih *et al*, 2014), sentence modeling (Xiaodong Liu *et al*, 2015), and other traditional NLP tasks.

Convolutional Neural Networks are used in Natural Language Processing in general and Named Entity Recognition in particular because of their ability to extract features. These networks are applied to embedding vectors of a given sentence with the aim of extracting useful features such as relationships between words that are closer together in a sentence, or the relation between sentences. Moreover, Convolutional Neural Networks will capture the relationships and general meanings of a sentence, which also represents an important advantage for good classification.

As previously stated in this thesis, SpaCy is an industrial-strength classifier, meaning that it is currently used by companies in order to achieve different tasks of a large scope, tasks which include product management, human resource management as well or search engine tuning. According to the lead-developer of SpaCy, companies such as Airbnb (which is an accommodation booking company), Quora (an online technology forum) or Allen Institute for Artificial Intelligence currently employ this software to extend and improve their services.

2.5 The Theoretical Foundations of the Gazetteer Method in Named Entity Recognition

2.5.1 The Brute Force Method

The first step in implementing this method implies the creation of a different list of items for each entity one desires to recognize, on which search operations will be applied in a later step in order to classify names.

Secondly, a large corpus of data needs to be appended to each list in order to achieve a great level of precision and recall for each of the named entities.

This method is quite restrictive due to the difficulties of comprising a corpus large enough to accomodate the varieties of named entities one may encounter, as well as the difficulties in keeping it up to date automatically. Moreover, the ambiguity resolution factor (the interpretation of metaphors, entities of the same name etc) is rather low, a fact which also leads to low precision.

Despite its disadvantages, the Gazetteer Method is a highly powerful method in named entity recognition, due to its speed, ease of manual improvement and relative ease of use.

2.5.2 Improving the Brute Force Method

In order to improve the method presented in the previous subsection, one must use a part-of-speech tagger in order to extract only said words that are proper nouns, such as to optimize the speed of the search algorithm and improve the overall performance. After having taken these steps, the search algorithm may be applied on each proper noun which is an output of the part-of-speech tagger.

To achieve those previously stated, the application presented in this thesis uses a part of the Stanford NER model, the Stanford POS Tagger. This tool applies the Conditional Random Field algorithm on the set of untagged data for obtaining the Part-Of-Speech Tagging necessary. The input data for the POS Tagger is comprised of a Romanian language model available on the Stanford platform, said model containing over 125000 tagged parts of speech.

2.6 ~~Issues of the Current~~ State-of-the-Art Methods with Respect to The Romanian Language

The aforementioned methods of tagging named entities are widely used by companies and researchers alike, being highly appreciated for their reliability and overall correctness. However, these methods have been created for certain internationally-spoken languages, achieving the best F-measure score in languages which are highly restrictive with respect to the topic of the sentence.

For example, in the English language, one may represent a noun-phrase simply by the following regular expression:

$$pattern = ' NP : < DT >? < JJ > * < NN > ',$$

in which NP depicts noun-phrase, DT is a determiner, JJ an adjective and NN a noun. The previously stated regular expression translates as: a noun phrase, NP, should be formed whenever the chunker finds an optional determiner, DT, followed by any number of adjectives, JJ, and then a noun, NN. One needs a noun phrase in order to more efficiently classify words as chunks, as well as to get a better understanding of the context and reduce the number of entities to be classified. For example in the following sentence: "Google, the giant American company, has been fined \$5.1 billion dollars.", the chunk "the giant American company" is a noun phrase which will offer the named entity recognition software the context needed to classify Google as an organization.

Meanwhile, in the Romanian language, the above stated sentence may be translated as "*Google, giganta companie americana, a fost amendata cu 5.1 miliarde de dolari.*", which is a natural context having a natural topic of the sentence. One may observe that the order in which the adjectives are placed is highly flexible, ranging from the above stated $\langle DT \rangle? \langle JJ \rangle * \langle NN \rangle \langle JJ \rangle *$ to $\langle DT \rangle? \langle NN \rangle \langle JJ \rangle *$ ("compania giganta americana"), a fact which is much harder to represent using regular expressions because there exist ~~much~~ more variations in the topic of a sentence ~~rather~~ than in English.

In addition to the added flexibility of the topic and of the more nuanced variety of expressions, the Romanian language, being a much less known language than English benefits from a much smaller international research community to create models, gazetteers, training data, or named entity recognition software. The previous argument is to be added to the author's motivation to create a template for the aforementioned research community in the field of Natural Language Processing in general and Named Entity Recognition in particular.

3 The Romanian Named Entity Recognition Software

The object of the current thesis is to ~~comprise~~ the three methods described in the previous chapter ~~and to implement~~ a reliable, fast Named Entity Recognition Software for the Romanian language in the form of a web application.

The resources that will be used are to be detailed and exemplified in the current chapter, as well as the application logic and the results of the processing.

3.1 Data Preprocessing

In order to obtain fast, reliable data, the corpora described in the ~~the~~ first section of the previous chapter in the current thesis has to be preprocessed in order to obtain a ConLL-UP formatted file, the one that is needed for the training of the customly-built named entity tagger described in this thesis.

Firstly, one has to extract the relevant data from the RONEC file, in order to comprise the training file for the Stanford NER software. This feat is to be accomplished by making use of the Python script in Figure 4. After data has been processed, the ronec.conllu is presented as in Figure 5.


```

1  import os
2
3  f = open("../Resources/ronec.conllu", 'r', encoding='utf8').read()
4
5  preprocessed = ""
6
7  for line in f.split('\n'):
8      if len(line) > 1 and '#' not in line:
9          tab = line.split('\t')
10         if len(tab) >= 10:
11             preprocessed += tab[1] + '\t' + tab[10][2:] + '\n'
12

```

Figure 4: Preprocessing the RONEC corpus

```

Grupul ORGANIZATION
Şcolar
"
M.Eminescu
"
Jimbolia
(
205 NUMERIC_VALUE
)
,
Grupul ORGANIZATION
Şcolar
Sănnicolau
Mare

```

Figure 5: Output of the Python Script in Figure 4.

One may notice that the shape of the file to be used for training is not ideal at this point, as there are words without tagging, mainly those that are parts of noun phrases. In the interest of simplicity, the following Python script will be used in order to complete the file with the missing tags. One is inclined to see that if between two tagged words there are words which are not tagged, there is a high probability that the words between the two tags are of the same tag as the former. Errors are insignificant to the training, as the corpus contains 120000 entries.

This script will loop 500 times to ensure that the **dummy** Romanian corpus is complete and contains no white spaces in the column assigned to

the named entity type.

```
def loop_forever():
    f = open("../Resources/ronec.collu", 'r', encoding='utf8')

    data = f.read()
    lines = data.split("\n")
    f.close()

    # print(lines[672].split("\t"))
    to_write = ""
    for idx in range(0, len(lines)):

        if len(lines[idx].split('\t')) == 1:
            lines[idx] += '\t'
        if lines[idx].split("\t")[1] == '':
            to_write += lines[idx] + lines[idx-1].split("\t")[1] + '\n'
        else:
            to_write += lines[idx] + '\n'

    g = open("D:/Anul_3/Lic/stanford/train/dummy-romanian-corpus.tsv", 'w', encoding='utf8')
    g.write(to_write)

for i in range(0, 500):
    loop_forever()
```

Figure 6: Completion Python script.

One may proceed with the extraction of data from the QuoVadis corpus, with the mention that a python xml parser will be needed to extract the words containing the tag <entity>< /entity>, as well as to extract the "Type" attribute and the text of the entry in the .xml file.

In the aftermath of those above, a complete and comprehensive training file is comprised and ready ~~ready~~ to serve both as a gazetteer and an input for the Conditional Random Field classification algorithm.

3.2 ~~The~~ Stanford Named Entity Recognition Method

Having comprised the training data, one is ready to apply the Stanford NER, by far the most reliable Natural Language Processing library on the market. The aforementioned package is written in Java, so one will need a proper Java Development Kit installed in order to run it.

Having prepared the environment needed, one needs to load the Stanford NER engine, available online for free for academic purposes.

```

în O
special O
Ucrainei GPE
, O
Georgiei GPE
și O
Republicii GPE
Moldova GPE
, O
unde O
Moscova GPE
și- O
a O
pierdut O
influența O
după O
revoluțiile O
oranj O
. O
Se O
asigură O
condiții O
deosebite O
pentru O
asimilarea O
cunoștințelor O
, O
C.P.P.P. ORGANIZATION
..

```

Figure 7: The output of the Python script in Figure 6.

```

<ENTITY id="8000200000" type="PERSON">
  <S Case="oblique" Definiteness="no" EXTRA="NotIndict" Gender="feminine" LDRG="Petronius" MSD="Npfpn" Number="plural" POS="NOUN" Type="proper" deprel="obj" head="7" id="5" offset="20">Petronius</S>
</ENTITY>

```

Figure 8: The exemplification of an "entity" type entry in the QuoVadis corpus.

The classification file compiled at the previous step (Subsection 3.1 Data Preprocessing) will now serve as a training file in the context of the program itself, as well as the two vectors (**x** - word vector, **s** - state vector) presented in Subsection 2.1.1 - "The Conditional Random Field Algorithm, Theoretical Considerations of the Stanford Software", the respective file is to be placed in the folder `./train` of the home folder of the Stanford classifier software. Moreover, the classifier software needs to be given a list of proprieties, most of which are fixed for the Named Entity Recognition, some of which may be changed. Said proprieties are presented in Figure 9.

The `trainFile` field represents the location of the `.tsv` (tab separated variables) file presented in Section 2.0.1 and built at the previous section, `serializeTo` field is the location to which the serialized, trained file will be written, while the `map` field is the exact layout of the columns of the `.tsv` file, in this case the numbering of the columns will start from 0, the first column being named as "word", whilst the second column (1) will be named as "answer", meaning the type of the named entity. It is important to note

that the first column will be the vector of x (words) in the Conditional Random Field Algorithm (Subsection 2.1.1), while the second will be the vector of s (states) in the same algorithm.

```
trainFile = D:/Anul_3/Lic/stanford/train/dummy-romanian-corpus.tsv
serializeTo = D:/Anul_3/Lic/stanford/train/dummy-ner-model-romanian.ser.gz
map = word=0,answer=1

useClassFeature=true
useWord=true
useNGrams=true
noMidNGrams=true
maxNGramLeng=6
usePrev=true
useNext=true
useSequences=true
usePrevSequences=true
maxLeft=1
useTypeSeqs=true
useTypeSeqs2=true
useTypepySequences=true
wordShape=chris2useLC
useDisjunctive=true
```

Figure 9: The list of proprieties of the Stanford NER.

Having set the proprieties, one has to run the command in Figure 10 in the home folder of his Stanford NER software.

```
PS D:\Anul_3\Lic\stanford> java -cp "*" -mx4g edu.stanford.nlp.ie.crf.CRFClassifier -prop train/prop.txt
```

Figure 10: The command to run the Stanford Trainer.

Having run the command, the training will output the serialized, trained file, having applied the Conditional Random Field Classifier on all the windows it produced for over 100 iterations. The training process is time and memory consuming. For the latter, one may observe that the fourth parameter of the command in Figure 10 is the memory allocated for the training process (in the current case 4 GB of memory have been allocated).

In Figure 11, the training time for 5000 data entries is presented, which is around 6 minutes, the training having been done after 131 iterations and will produce satisfying results at tagging level. However, for a higher F-measure score, this thesis will use all of the pretrained data that has been comprised at the previous subsection, adding up to 195000 entries.

```

Iter 120 evals 141 <D> [M 1.000E0] 9.456E2 353.80s {5.058E0} {2.421E-4} 2.491E-4 -
Iter 121 evals 142 <D> [M 1.000E0] 9.453E2 356.50s {1.222E1} {5.851E-4} 2.353E-4 -
Iter 122 evals 143 <D> [M 1.000E0] 9.451E2 359.17s {7.131E0} {3.413E-4} 2.241E-4 -
Iter 123 evals 144 <D> [M 1.000E0] 9.449E2 361.84s {4.657E0} {2.229E-4} 2.113E-4 -
Iter 124 evals 145 <D> [M 1.000E0] 9.448E2 364.50s {1.229E1} {5.883E-4} 1.866E-4 -
Iter 125 evals 146 <D> [M 2.810E-1] 9.447E2 368.31s {8.443E0} {4.042E-4} 1.777E-4 -
Iter 126 evals 148 <D> [M 1.000E0] 9.445E2 371.13s {4.352E0} {2.083E-4} 1.729E-4 -
Iter 127 evals 149 <D> [M 1.000E0] 9.444E2 374.08s {3.023E0} {1.447E-4} 1.549E-4 -
Iter 128 evals 150 <D> [M 1.000E0] 9.443E2 376.94s {4.487E0} {2.148E-4} 1.498E-4 -
Iter 129 evals 151 <D> [M 1.000E0] 9.443E2 380.33s {1.053E1} {5.041E-4} 1.363E-4 -
Iter 130 evals 152 <D> [M 1.000E0] 9.442E2 384.24s {5.090E0} {2.437E-4} 1.247E-4 -
Iter 131 evals 153 <D> [M 1.000E0] 9.441E2 387.83s {3.260E0} {1.560E-4} 1.033E-4 -
QNNMinimizer terminated due to average improvement: | newest_val - previous_val | / |newestVal| < TOL
Total time spent in optimization: 391.42s
CRFClassifier training ... done [395.8 sec].
Serializing classifier to D:/Anu1_3/Lic/stanford/train/dummy-ner-model-romanian.ser.gz... done.

```

Figure 11: Result for 5000 entries.

In order to run the classifier on all 195000 entries, various tweaks had to be made to the training file, such as the removal of some irrelevant punctuation (such as "():,"), because otherwise the Java stack would be exceeded.

```

Iter 236 evals 268 <D> [M 1.000E0] 7.432E3 624.17s {4.904E1} {2.306E-4} 1.142E-4 -
Iter 237 evals 269 <D> [M 1.000E0] 7.432E3 626.54s {2.377E1} {1.118E-4} 1.064E-4 -
Iter 238 evals 270 <D> [M 1.000E0] 7.431E3 628.89s {7.996E1} {3.760E-4} 1.003E-4 -
QNNMinimizer terminated due to average improvement: | newest_val - previous_val | / |newestVal| < TOL
Total time spent in optimization: 631.27s
CRFClassifier training ... done [643.3 sec].
Serializing classifier to D:/Anu1_3/Lic/stanford/train/dummy-ner-model-romanian.ser.gz... done.

```

Figure 12: Result for all 195000 entries.

In order to run the now-trained classifier on raw, unannotated data, one must implement a Python script, using the StanfordNERTagger class of the NLTK (Natural Language Toolkit) Python library this thesis has mentioned in the Introduction section.

An example of such Python script is presented in Figure 13, this exact script will be used in the web application this thesis aims to implement. Moreover, in Figure 14, one notices the results of a tagging example for the purpose of visualizing the results. These results are to be used as data for the Statistics section of this thesis, where one will be introduced to the numerical facts of this classification, compared to the other two that will be presented in the following subsections.

One may notice that the output of the classification in this case is a list of tuples having the following format [(word_text1, NER_tag1),(word_text2, NER_tag2), ..., (word_textn, NER_tagn)]. This format is the one to be used in the manual training of the test data used in the Statistics section of the application, as well as for the other classifiers.

```

1 import nltk
2 from nltk.tag.stanford import StanfordNERTagger
3
4 def run(sentence):
5     import os
6     java_path = "C:/Program Files/Java/jdk-12.0.1/bin"
7     os.environ['JAVAHOME'] = java_path
8
9     jar = 'D:/Anul_3/Lic/stanford/stanford-ner.jar'
10    model = 'D:/Anul_3/Lic/stanford/train/dummy-ner-model-romanian.ser.gz'
11
12    ner_tagger = StanfordNERTagger(model, jar, encoding='utf8')
13
14    words = nltk.word_tokenize(sentence)
15    raw_data = ner_tagger.tag(words)
16    to_return = []
17    for tuplu in raw_data:
18        if tuplu[1] != 'O':
19            to_return.append(tuplu)
20
21    return to_return

```

Figure 13: Example of Python script to run the trained classifier.

```

PS D:\Anul_3\Lic\Training Methods> python .\ner_romanian.py
[('Iudea', 'PERSON'), ('Arghezi', 'PERSON'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('Academia', 'PERSON'), ('Romana', 'PERSON'), ('21', 'MISC'), ('mai', 'MISC'), ('1880', 'MISC'), ('14', 'MISC'), ('Iulie', 'MISC'), ('1967', 'MISC'), ('scriitor', 'PERSON'), ('roman', 'MISC')]

```

Figure 14: Output of the Classifier.

3.3 The Practical Application of the Gazetteer Method

3.3.1 Creation of the List of Entities

The creation of this corpus is highly time consuming. The internet is the best resource as to gain as much information as possible in order to complete the list necessary for a satisfactory classification of some more complex sentences.

Data from the Wikipedia article dedicated to Romanian surnames have been extracted, a list of the most common 2000 first names in the language, as well as data from the previously presented RONEC and QuoVadis corpora have been extracted in order to comprise the "PERSON.txt" file for named entities of this particular types. The same procedure has been applied for the other 3 categories of named entities (ORGANIZATION, MISC, GPE), with the mention that for the geographical entities, data has also been extracted from the gazetteer comprised in the study "A Mixed Approach in

Recognising Geographical Entities in Texts” (D.Cristea *et al*, 2015).

The format of the above mentioned files is one-per-line in a simple text file, for easier access and reduced complexity of the search algorithm. Figure 15 depicts a snippet of the ”PERSON.txt” file, comprising of 43727 entries.

43718	Țăroi
43719	Țăroiu
43720	Țăruș
43721	Țăruși
43722	Țărână
43723	Țărîndă
43724	Țărângoiu
43725	Țărăpănel
43726	Țărăscu
43727	Țărău

Figure 15: Snippet of the PERSON.txt file.

3.3.2 The Operation of the Gazetteer Method

First and foremost, one must run the Stanford POS Tagger in order to appropriately tag the parts of speech in the given test text. The code of said tagging is depicted in Figure 16, while the output in Figure 17. The test data to which the tagger is applied is the sentence: "Tudor Arghezi, pseudonimul lui Ion Nae Theodorescu, Academia Romana, nascut la 21 mai 1880, Bucuresti decedat in 14 iulie 1967 a fost un scriitor roman, cunoscut pentru contribuția sa la dezvoltarea liricii românești sub influența baudelaireanismului. Opera sa poetică, de o originalitate exemplară, reprezintă o altă vârstă marcantă a literaturii române. A scris, între altele, teatru, proză (notabile fiind romanele Cimitirul Buna Vestire și Ochii Maicii Domnului), pamflete, precum și literatură pentru copii. A fost printre autorii cei mai contestați din întreaga literatură română."

```
def extract_pos(doc):
    result = []
    for sent in doc.sentences:
        for wrd in sent.words:
            result.append((wrd.text, wrd.upos))
    #return a dataframe of pos and text
    return result
```

Figure 16: The function for the Part of Speech Tagger.

```
[('Tudor', 'PERSON'), ('Arghezi', 'PERSON'), (',', 'PUNCT'), ('pseudonimul', 'NOUN'), ('lui', 'DET'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('.', 'PUNCT'), ('Academia', 'NOUN'), ('Romana', 'PERSON'), (',', 'PUNCT'), ('nascut', 'VERB'), ('la', 'ADP'), ('21', 'NUM'), ('mai', 'NOUN'), ('1880', 'NUM'), (',', 'PUNCT'), ('Bucuresti', 'PERSON'), ('decedat', 'VERB'), ('in', 'ADP'), ('14', 'NUM'), ('iulie', 'NOUN'), ('1967', 'NUM'), ('a', 'AUX'), ('fost', 'AUX'), ('un', 'DET'), ('scriitor', 'NOUN'), ('roman', 'ADJ'), (',', 'PUNCT'), ('cunoscut', 'VERB'), ('pentru', 'ADP'), ('contribuția', 'NOUN'), ('sa', 'DET'), ('la', 'ADP'), ('dezvoltarea', 'NOUN'), ('liricii', 'NOUN'), ('românești', 'ADJ'), ('sub', 'ADP'), ('influența', 'NOUN'), ('baudelaireanismului', 'NOUN'), (',', 'PUNCT'), ('Opera', 'NOUN'), ('sa', 'DET'), ('poetică', 'ADJ'), (',', 'PUNCT'), ('de', 'ADP'), ('o', 'DET'), ('originalitate', 'NOUN'), ('exemplară', 'ADJ'), (',', 'PUNCT'), ('reprezintă', 'VERB'), ('o', 'DET'), ('altă', 'DET'), ('vârstă', 'NOUN'), ('marcantă', 'ADJ'), ('a', 'ADP'), ('literaturii', 'NOUN'), ('române', 'ADJ'), (',', 'PUNCT'), ('A', 'AUX'), ('scris', 'VERB'), (',', 'PUNCT'), ('între', 'ADP'), ('altale', 'PERSON'), (',', 'PUNCT'), ('teatru', 'NOUN'), (',', 'PUNCT'), ('proză', 'NOUN'), (',', 'PUNCT'), ('notabile', 'ADJ'), ('fiind', 'AUX'), ('romanele', 'NOUN'), ('Cimitirul', 'NOUN'), ('Buna', 'ADJ'), ('Vestire', 'NOUN'), ('și', 'CCONJ'), ('Ochii', 'NOUN'), ('Maicii', 'NOUN'), ('Domnului', 'NOUN'), (',', 'PUNCT'), (',', 'PUNCT'), ('pamflete', 'NOUN'), (',', 'PUNCT'), ('precum', 'ADP'), ('și', 'CCONJ'), ('literatură', 'NOUN'), ('pentru', 'ADP'), ('copii', 'NOUN'), (',', 'PUNCT'), ('A', 'AUX'), ('fost', 'AUX'), ('printre', 'ADP'), ('autorii', 'NOUN'), ('celi', 'DET'), ('mai', 'AUX'), ('contestați', 'VERB'), ('din', 'ADP'), ('întreaga', 'ADJ'), ('literatură', 'NOUN'), ('română', 'ADJ'), (',', 'PUNCT'), ('.', 'PUNCT'), ('Tudor', 'PERSON'), ('Arghezi', 'PERSON'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('Romana', 'PERSON'), ('Bucuresti', 'PERSON')]
```

Figure 17: Output run on test data.

As the gazetteer method this thesis applies does not take context into consideration, one may filter the results presented above and extract only said words that are nouns and proper nouns, and then proceed to the tagging itself, with the returned result being of the same type as that of the one outputted by the classifier in the previous section. These facts are depicted in Figures 18 and 19.

```
def runGazetteer(sentence):
    f = open(r"D:\Anu1_3\Lic\Resources\PERSON.txt", "r", encoding="utf8").read()
    g = open(r"D:\Anu1_3\Lic\Resources\GPE.txt", "r", encoding="utf8").read()
    h = open(r"D:\Anu1_3\Lic\Resources\MISC.txt", "r", encoding="utf8").read()
    j = open(r"D:\Anu1_3\Lic\Resources\ORGANIZATION.txt", "r", encoding="utf8").read()
    # extract pos
    # print(extract_pos(doc))
    person = []
    gpe = []
    misc = []
    organization = []
    for line in f.split('\n'):
        person.append(line)
    for line in g.split('\n'):
        gpe.append(line)
    for line in h.split('\n'):
        misc.append(line)
    for line in j.split('\n'):
        organization.append(line)
    sent = nlp(sentence)
    pos_tagged = extract_pos(sent)
    ner_list = []
    # print("Arghezi" in person)

    for word in pos_tagged:
        if word[1] == 'NOUN' or word[1] == 'PROPN':
            if word[0] in person:
                ner_list.append((word[0], "PERSON"))
            elif word[0] in gpe:
                ner_list.append((word[0], "GPE"))
            elif word[0] in organization:
                ner_list.append((word[0], "ORGANIZATION"))
            elif word[0] in misc:
                ner_list.append((word[0], "MISC"))

    return ner_list
```

Figure 18: The function for the classification of the proper nouns.

One may see that the proper nouns that describe persons have been correctly classified, as well as those describing geographical location, meanwhile those describing miscellaneous named entities or organizations named

entities, due to the vast scope of their respective entities are not perfectly recognized.

```
[('Istoc', 'PERSON'), ('Arghezi', 'PERSON'), ('Ise', 'PERSON'), ('Iac', 'PERSON'), ('Theodorescu', 'PERSON'), ('Academia', 'ORGANIZATION'), ('Rumania', 'PERSON'), ('mal', 'MISC'), ('Iulie', 'MISC'), ('Ispire', 'ORGANIZATION'), ('Istaitu', 'ORGANIZATION'), ('Moldu', 'MISC'), ('Dumitru', 'MISC'), ('Cantii', 'ORGANIZATION')]
```

Figure 19: Output run on test data.

3.4 The SpaCy Named Entity Recogniser using the Multi-language Model

In order to run the SpaCy software on a language that does not have any pretrained model, one needs to download the multi language model available on the SpaCy website in order to access data extracted from various Wikipedia pages, comprised into a training data set for the convolutional neural network (described in Subsection 2.2). The aforementioned model supports identification of PERSON, GPE, ORGANIZATION and MISC entities, similar to the other two methods described in the previous chapters.

Having downloaded the resources necessary for the execution, one has to import and load the model into the memory, as described in Figure 20.

```
import spacy
import xx_ent_wiki_sm
nlp = xx_ent_wiki_sm.load()
```

Figure 20: SpaCy imports

The SpaCy NER tagger will then work automatically, tagging the instances, the output being formatted to the format used in the previous two sections of the thesis (list of tuples of the form (word, entity), where entity is one of PERSON, ORGANIZATION, GPE or MISC). The execution of the SpaCy tagger is shown in Figure 21, while the output is in Figure 22, the input data being the same sentence as the aforementioned two other methods of classification.

It is to be observed that the classification is correct to a certain degree and becomes satisfying when taking into account that this classifier has no specific model for the Romanian language.

```

def run_english(sent):
    import spacy
    import xx_ent_wiki_sm
    nlp = xx_ent_wiki_sm.load()

    sentence = nlp (sent)

    l = [(X.text, X.label_) for X in sentence.ents]
    # print(l)
    to_return = []

    # print(l)
    for item in l:
        if len(item[0].split(' ')) == 1:
            if item[1] == 'PER':
                new_tuple = (item[0], 'PERSON')
            elif item[1] == 'ORG':
                new_tuple = (item[0], 'ORGANIZATION')
            elif item[1] == 'GPE':
                new_tuple = (item[0], 'GPE')
            else:
                new_tuple = (item[0], 'MISC')
            to_return.append(new_tuple)
        else:
            for tab in item[0].split(' '):
                if item[1] == 'PER':
                    new_tuple = (tab, 'PERSON')
                    to_return.append(new_tuple)
                elif item[1] == 'ORG':
                    new_tuple = (tab, 'ORGANIZATION')
                    to_return.append(new_tuple)
                elif item[1] == 'GPE':
                    new_tuple = (tab, 'GPE')
                    to_return.append(new_tuple)
                else:
                    new_tuple = (tab, 'MISC')
                    to_return.append(new_tuple)

    return to_return

```

Figure 21: SpaCy execution and formatting the data.

3.5 Statistics

In order to achieve comprehensive statistics, one must run the three aforementioned methods of classification on multiple examples from different

```
[('Tudor', 'PERSON'), ('Arghezi', 'PERSON'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('Academia', 'ORGANIZATION'), ('Romana', 'ORGANIZATION'), ('Bucuresti', 'MISC'), ('Gimnaziul', 'PERSON'), ('Buna', 'PERSON'), ('Vestire', 'PERSON'), ('Ochii', 'PERSON'), ('Maicii', 'PERSON'), ('Domului', 'PERSON')]
```

Figure 22: SpaCy results

spheres and contexts. For doing so, a list of 20 pieces of text has been comprised, the contents of which varying from scientific text, journalistic text, social media, academic or online article.

The format chosen for this collection of data is a dictionary with two keys, the 'raw' tab, and the 'tagged' tab, with the raw tab containing the text as it is, unannotated, while the tagged tab contains the Named Entity Recognition done manually by the author. Figure 23 presents a snippet of the file mentioned above.

```
{"sentence": {
  "raw": "Tudor Arghezi, pseudonimul lui Ion Nae Theodorescu, Academia Romana, nascut la 21 mai 1980, a decedat in 14 iulie 1967 a fost un scriitor rom\u00e2n, cunoscut pentru",
  "tagged": "[('Tudor', 'PERSON'), ('Arghezi', 'PERSON'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('21', 'MISC'), ('mai', 'MISC'), ('1980', 'MISC'), ('14', 'MISC'), ('iulie', 'MISC'), ('1967', 'MISC'), ('fost', 'MISC'), ('un', 'MISC'), ('scriitor', 'MISC'), ('rom\u00e2n', 'MISC'), ('cunoscut', 'MISC'), ('pentru', 'MISC')]"
}
```

Figure 23: Dictionary of test data.

In order to present the precision and recall functions, one must define some key terms.

The first of those is the term of true positive. One calls a word a true positive, if the word is a named entity, the classifier tagged the word as a named entity, and the category in which it placed, the word is correct.

Secondly, one calls a false positive a word which is not a named entity, but has been classified as such.

Thirdly, a false negative is a either a named entity which has not been classified at all, or one which has been wrongly classified.

Lastly, a true negative is a word that is neither a named entity, nor has it been classified as such.

One computes the following metrics, which have theoretically addressed in Section 1 of the current thesis, for calculating the performance of each classifier:

1.

$$precision = \frac{true_positives}{true_positives + false_positives}$$

2.

$$recall = \frac{true_positives}{true_positives + false_negatives}$$

3.

$$F - measure = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

Each of these three metrics are relevant to the overall performance of one's Named Entity Recognition software. In this current thesis we address all 3 of them in order to assess the level of correctness each method of classification previously described achieves.

For each of the methods, one will compute a function in order to compute the parameters previously stated, as in the example in Figure 24.

```
elif method.lower() == 'blind':  
    for sentence in sentences:  
        tagged_sentence = ner_english.run_english(sentence['raw'])  
        for item in tagged_sentence:  
            if str(item) in sentence['tagged']:  
                true_positives += 1  
            elif str(item[0]) in sentence['tagged']:  
                false_negatives += 1  
            else:  
                false_positives += 1  
        false_negatives += (all_correct_tagged - true_positives)  
        true_negatives = all_words - all_correct_tagged  
        # print(all_words, all_correct_tagged)
```

Figure 24: The statistics function for the SpaCy method.

It is highly important to mention that variables "all_words" and "all_correct_tagged" have been previously computed, the former being the sum of all words in the test examples (be them named entities or not), while the latter the sum of all manually tagged words (all of the named entities).

The result of these calculations is an excellent indicator to how the system performs, and which of the methods is better on this particular data set.

Figure 25 depicts the output of the precision algorithm presented above, similarly, Figures 26 and 27 present the output of the recall and F-measure algorithms respectively.

One must note that the results are bounded by 1 (it being the maximum precision, recall and F-measure achievable). If we were to scale the result to percentages, the most exact overall classifier (the metric of choice is the F-measure score), would be the trained Stanford Conditional Random

```
The precision of the trained Stanford CRFClassifier is: 0.8390804597701149
The precision of the Gazetteer Classifier is: 0.5858585858585859
The precision of the SpaCy multi-language model Classifier is: 0.75
```

Figure 25: The output of the precision algorithm.

```
The recall of the trained Stanford CRFClassifier is: 0.863905325443787
The recall of the Gazetteer Classifier is: 0.7733333333333333
The recall of the SpaCy multi-language model Classifier is: 0.4067796610169492
```

Figure 26: The output of the recall algorithm.

```
The F-measure result of the trained Stanford CRFClassifier is: 0.8513119533527697
The F-measure result of the Gazetteer Classifier is: 0.6666666666666667
The F-measure result of the SpaCy multi-language model Classifier is: 0.5274725274725275
```

Figure 27: The output of the F-measure algorithm.

Field Classifier, which has been previously presented in Subsection 2.1 from a theoretical standpoint and in Subsection 3.2 from a practical standpoint.

The results show that inevitably the best classifier ~~is~~, with a score of 85.1% [▲] is the trained Stanford classifier, while the least accurate is the SpaCy, an expected result, due to its lack of Romanian language-specific models.

The lack of higher accuracy in the context of the Gazetteer method is due to its lack of context data, as well as due to the lack of sufficient data in order to interpret different ORGANIZATION and MISC entities, despite the relatively high number of instances present in the training data.

3.6 The Voting Function

In order for the application to achieve the best possible classification comprised of the three methods, one must implement a voting system in order to best choose the classification needed for the untagged sentence.

This thesis proposes a voting function which combines the F-measure results computed **at** the previous subsection and **makes them into ratios of the their sum**.

~~This being said,~~ the general formula for each of the weights is:

$$\frac{f_i}{\sum_{i=1}^3 f_i}$$

One will then run each of the classifiers, ~~and~~ presuming the same noun phrase is caught by all classifiers in the same way, will compute the overall ratio of said noun phrase.

All of the weights will be saved in a list as proper numbers. So, the weight of the first classifier is 0.42 (the result of $\frac{0.86}{2.04}$, where 0.86 is the F-measure score of the trained Stanford NER Classifier and 2.04 is the sum of all classifiers), that of the second classifier is 0.32 (the result of $\frac{0.66}{2.04}$, where 0.66 is the F-measure score of the Gazetteer Classifier and 2.04 is the sum of all classifiers), and the weight of the third classifier is 0.25 (the result of $\frac{0.52}{2.04}$, where 0.52 is the F-measure score of the SpaCy Classifier and 2.04 is the sum of all classifiers). One might add that the sum of all these weights is 1, meaning that each classifier has been assigned the weight corresponding to how much it contributes to the sum of F-measure scores.

For example, if the word "Google" is present in the output of the first classifier, its score will be 0.42, if the same word is present in the second classifier, one will add to its score 0.32 and so on.

Having stated those above, one will compare the sum of the weights left in the list of weights (those in which the word is not present in) to the score it has after the aforementioned computations have been performed. If the score is higher than the sum of ratios, it will be present in the final result list.

The aforementioned steps will maximize the chances of the classification to be correct and complete.

3.7 The Web Application

For the user interface, this thesis has adopted the use of a web application with a user-friendly front-end, as well as all the necessary information about the thesis. One can observe that a minimalist approach has been chosen, as depicted in Figure 30.

The operation of the aforementioned application is fairly simple, as the user has to input the sentence he wants to label, and the result of said classification will appear on the right side of the box, just as Figure 31 depicts.

The logic of this application is as follows, the user input is extracted from the front-end, it being passed on to the voting function described previously. The voting function will do the named entity recognition as described in the previous sections, then call the beautify function (Figure 32), which will return the newly formed text, with the words labeled, using the colours described in the legend section of the page.

In order for all of the above to be fulfilled, Flask, the microframework described at Section 2.2 has been used, employing the same approach

```

def vote(sentence):
    result = []
    trained = ner_romanian.run(sentence)
    gazetteer = stanford.runGazetteer(sentence)
    spacy = ner_english.run_english(sentence)

    for data in trained:
        result.append(data)
    for data in gazetteer:
        coefs = [0.85, 0.09, 0.06] #according to each method's F-measure score
        score = 0.09
        coefs.remove(0.09)
        if data in trained and data not in result:
            score += 0.85
            coefs.remove(0.85)
        if data in spacy and data not in result:
            score += 0.06
        if score > sum(coefs):
            result.append(data)

    for data in spacy:
        coefs = [0.85, 0.09, 0.06] #according to each method's F-measure score
        score = 0.06
        coefs.remove(0.06)
        if data in trained and data not in result:
            score += 0.85
            coefs.remove(0.85)
        if data in spacy and data not in result:
            score += 0.09
        if score > sum(coefs):
            result.append(data)

    return result

```

Figure 28: The Python implementation of the voting function.

```

[[('Tudor', 'PERSON'), ('Arghizil', 'PERSON'), ('Ion', 'PERSON'), ('Nae', 'PERSON'), ('Theodorescu', 'PERSON'), ('Academia', 'PERSON'), ('Romana', 'PERSON'), ('21', 'MISC'), ('mail', 'MISC'), ('1888', 'MISC'), ('Bucuresti', 'LOC'), ('14', 'MISC'), ('Julie', 'MISC'), ('187', 'MISC'), ('scriitor', 'PERSON'), ('ramin', 'MISC'), ('ramdeviti', 'MISC'), ('ramdev', 'MISC'), ('ramdevia', 'PERSON'), ('Catinina', 'PERSON'), ('Bun', 'PERSON'), ('vestire', 'PERSON'), ('Ogali', 'PERSON'), ('Miciu', 'PERSON'), ('Domulul', 'PERSON'), ('il', 'PERSON'), ('ramdev', 'MISC')]]

```

Figure 29: The output of the voting function.

described in great detail there.

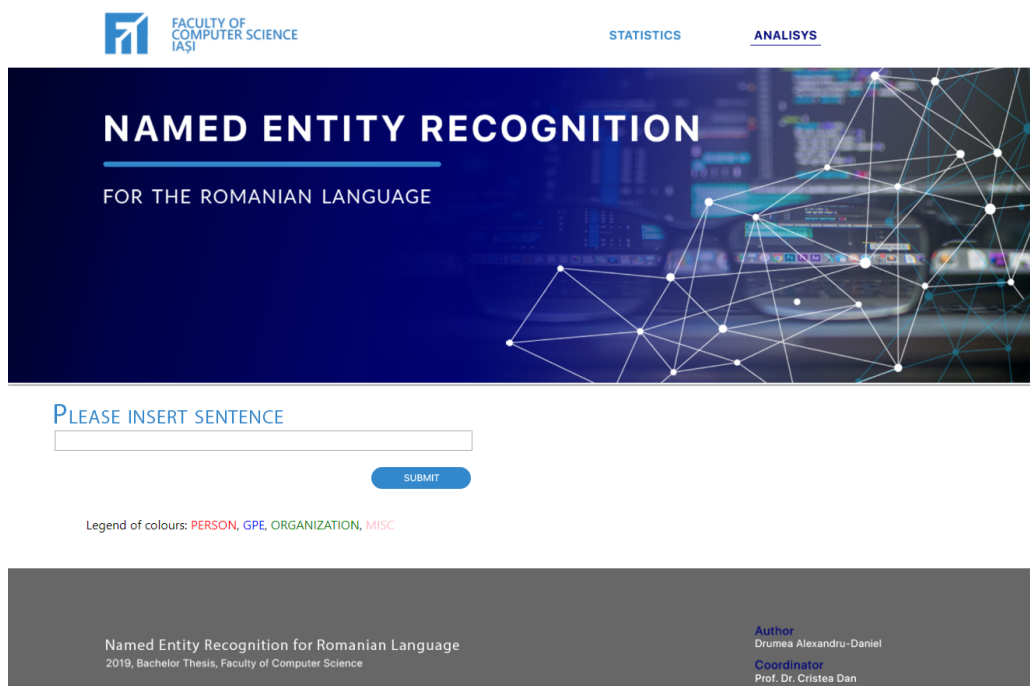


Figure 30: The homepage of the web application.

4 Conclusions

The methods used in the current thesis have returned results that are satisfying as well as comprehensive. One has managed to achieve an F-measure score of 85%, which consistently makes reliable classifications of sentences when it comes to named entities. If on a relatively short training data set the results have been such satisfying, one may ~~only~~ assume that on a much more complex data set, the metrics will ~~only~~ improve, the application ~~reaching~~ a higher level of ~~complexity~~.

As one may have noticed in the course of this brief documentation, the premise of having a named entity recognizer has been achieved in a very comprehensive manner.

Starting from what the market provides as support (SpaCy multi-language models), which not being trained specifically for the Romanian language, achieve an F-measure score of just 52%, despite being not satisfying enough, is not quite such a negligible score.

The Gazetteer method, which in some areas provide one with quite accurate classifications, has its drawbacks, having achieved an F-measure score of 66%, which primarily comes from the difficulty of keeping a gazetteer

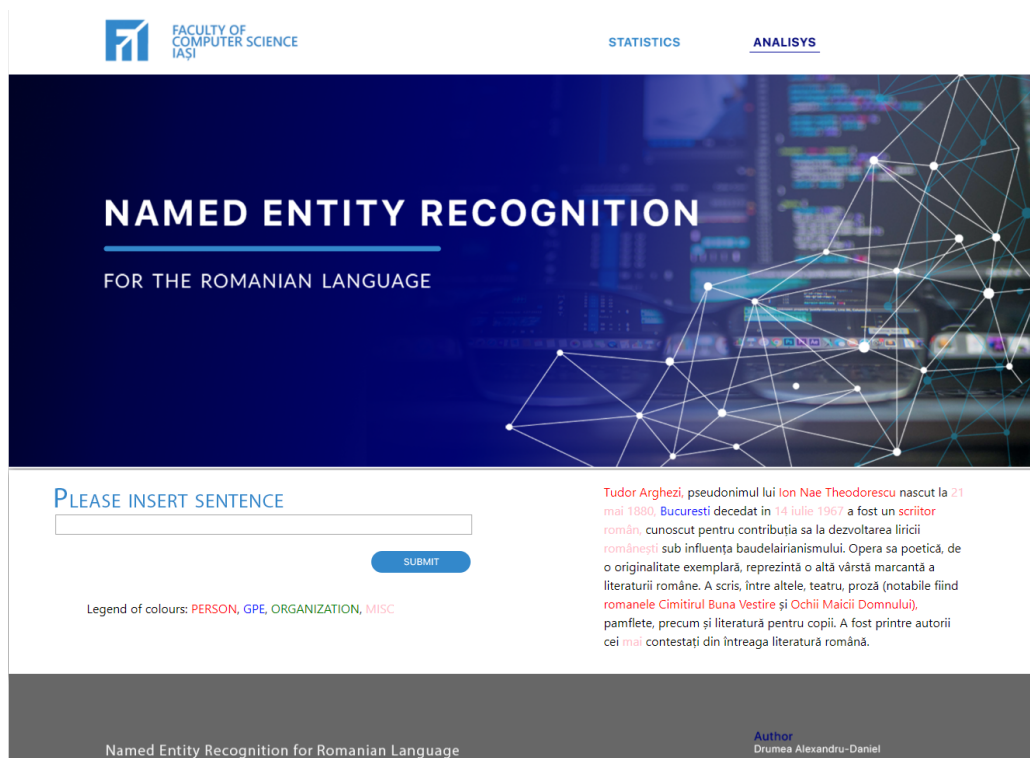


Figure 31: The homepage of the web application with labeled text.

up-to-date, as well as its lack of context awareness, thus making it not such a reliable, long-term tool.

The aim of this thesis has been reached by using the Conditional Random Field Algorithm of the Stanford Named Entity Recognition Software, which has obtained a score of 85% accuracy. Even though it is not quite at the level of the industry standard, it is a fine attempt in solving the problem of named entity recognition for the Romanian language.

Improvements are required in order for the application to become a reliable, powerful, industry standard tool, these improvements being: bigger contexts, more data, more ambiguity applied to said data, as well as much more powerful computers than the one used in the current thesis.

In conclusion, the aim set at the beginning of this paper has been achieved, the knowledge gathered has been documented and the improvements needed have been identified and, therefore, the premises for further studies with more ambitious goals have been set.

```

def beautify(sentence, ner_tags):

    to_return = ""
    for word in sentence.split(' '):
        aux = word
        new_word = ""
        for tupl in ner_tags:
            if tupl[0] == word.strip(',').strip(" ").strip("("):
                # print(word)
                if tupl[1] == 'PERSON':
                    # print(tupl[0])
                    new_word = "<span style=\"color: red\">" + aux + "</span>"
                    break
                elif tupl[1] == 'GPE':
                    new_word = "<span style=\"color: blue\">" + aux + "</span>"
                    break
                elif tupl[1] == 'ORGANIZATION':
                    new_word = "<span style=\"color: green\">" + aux + "</span>"
                    break
                else:
                    new_word = "<span style=\"color: pink\">" + aux + "</span>"
                    break
            if len(new_word) == 0:
                new_word = aux
        to_return += new_word + ' '

    return to_return

```

Figure 32: The beautify function.

5 Bibliography

1. Charles Sutton, University of Edinburgh and Andrew McCallum University of Massachusetts Amherst, "An Introduction to Conditional Random Fields", Vol. 4, No. 4 (2011) 267–373
2. Anca-Diana Bibiri, Mihaela Colhon, Paul Diac, Dan Cristea (2014). Statistics Over A Corpus Of Semantic Links: "QuoVadis". In Mihaela Colhon, Adrian Iftene, Verginica Barbu Mititelu, Dan Cristea, Dan Tufiş (eds.) (2014). Proceedings of the 10th International Conference "Linguistic Resources And Tools For Processing The Romanian Language, Craiova, 18-19 September 2014", "Alexandru Ioan Cuza" University Publishing House, pag. 33-44
3. Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learn-

- ing Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998
4. Wen-tau Yih Xiaodong He Christopher Meek Microsoft Research, Association for Computational Linguistics, "Semantic Parsing for Single-Relation Question Answering", June 2014
 5. Xiaodong Liu, Jianfeng Gao, Xiaodong He, Li Deng, Kevin Duh, Ye-Yi Wang, NAACL, "Representation Learning Using Multi-Task Deep Neural Networks for Semantic Classification and Information Retrieval", May 2015
 6. Nal Kalchbrenner, Edward Grefenstette, Phil Blunsom, Department of Computer Science, "A Convolutional Neural Network for Modelling Sentences", University of Oxford, 2014
 7. Ronan Collobert, Jason Weston, Leon Bottou, Michael Karlen, Koray Kavukcuoglu, Pavel Kuksa, "Natural Language Processing (Almost) from Scratch", Journal of Machine Learning Research 12 (2011) 2493-2537, 2011
 8. Stefan Daniel Dumitrescu, Andrei Avram, Luciana Morogan, Stefan Toma, "Romanian Named Entity Corpus", 2018
 9. <https://nlp.stanford.edu/software/CRF-NER.shtml>
 10. https://ro.wikipedia.org/wiki/Categorie:_Nume_de_familie_romanesti
Cristea D., Gifu D., Pistol I., Sfirnaciuc D., Niculiță M. (2016) A Mixed Approach in Recognising Geographical Entities in Texts. In: Trandabăț D., Gifu D. (eds) Linguistic Linked Open Data. RUMOUR 2015. Communications in Computer and Information Science, vol 588. Springer, Cham
 11. <https://spacy.io/models/xx>
 12. <https://towardsdatascience.com/named-entity-recognition-with-nltk-and-spacy-8c4a7d88e7da>
 13. <https://towardsdatascience.com/precision-vs-recall-386cf9f89488>
 14. Charles Sutton, University of Edinburgh and Andrew McCallum University of Massachusetts Amherst, "An Introduction to Conditional Random Fields", Foundations and Trends in Machine Learning Vol. 4, No. 4, 267-373, 2011

15. Manning, Christopher D., Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pp. 55-60.
16. <https://blog.sicara.com/train-ner-model-with-nltk-stanford-tagger-english-french-german-6d90573a9486>
17. <https://profs.info.uaic.ro/~dcristea/papers/RUMOUR-Cristea%20et%20al.pdf>
18. <https://pythonprogramming.net/named-entity-recognition-nltk-python/>
19. Manning, Christopher D. and Surdeanu, Mihai and Bauer, John and Finkel, Jenny and Bethard, Steven J. and McClosky, "The Stanford CoreNLP Natural Language Processing Toolkit", Association for Computational Linguistics, 55-60, 2014
20. <http://www.aclweb.org/anthology/P/P14/P14-5010>
21. <https://pythonhow.com/how-a-flask-app-works/>
22. <https://profs.info.uaic.ro/~dcristea/cursuri/IA/2018-2019/Curs13-%20Analiza%20limbajului%20natural.ppt.pdf>
23. <https://developers.google.com/machine-learning/crash-course/classification/precision-and-recall>
24. <https://universaldependencies.org/>
25. Doug Cutting and Julian Kupiec and Jan Pedersen and Penelope Sibun, "A Practical Part-of-Speech Tagger", ANLC '92 Proceedings of the third conference on Applied natural language processing, Pages 133-140, 1992
26. <https://medium.com/@jrodthoughts/ambiguity-in-natural-language-processing-15f3e4706a9a>
27. <https://www.quora.com/How-does-CNN-work-with-NLP>
28. https://en.wikipedia.org/wiki/Convolutional_neural_network#Natural_language_processing
29. <http://flask.pocoo.org/>
30. <https://www.techopedia.com/definition/13825/named-entity-recognition-ner>