

# TOWARDS MULTILINGUAL READABILITY ASSESSMENT

by

Ion Madrazo Azpiazu

A thesis

submitted in partial fulfillment

of the requirements for the degree of

Master of Science in Computer Science

Boise State University

December 2014



BOISE STATE UNIVERSITY GRADUATE COLLEGE

**DEFENSE COMMITTEE AND FINAL READING APPROVALS**

of the thesis submitted by

Ion Madrazo Azpiazu

Thesis Title: Towards Multilingual Readability Assessment

Date of Final Oral Examination: 1st December 2014

The following individuals read and discussed the thesis submitted by student Ion Madrazo Azpiazu, and they evaluated the presentation and response to questions during the final oral examination. They found that the student passed the final oral examination.

Sole Pera, Ph.D.

Chair, Supervisory Committee

Tim Andersen, Ph.D.

Member, Supervisory Committee

Nere Lete, Ph.D.

Member, Supervisory Committee

The final reading approval of the thesis was granted by Sole Pera, Ph.D., Chair of the Supervisory Committee. The thesis was approved for the Graduate College by John R. Pelton, Ph.D., Dean of the Graduate College.

## ACKNOWLEDGMENTS

The author wishes to express gratitude to Foo. This work would have been partially supported by some particular grant, if there was one.

## ABSTRACT

An *abstract* is a brief summary of the document. A typical abstract provides a brief introduction, enough to provide context for the document, explains the purpose of the thesis or project, and summarizes the major results and conclusions. Keep in mind that a casual observer is likely to judge the content of the document by the abstract and title alone. (There is an old adage: “in a joke, the punchline comes at the end; in a paper [or thesis], it comes in the abstract.”) A single concise paragraph usually suffices for the abstract. If it spills onto a second page, it is probably too long.

## TABLE OF CONTENTS

<b>ABSTRACT</b> .....	v
<b>LIST OF TABLES</b> .....	viii
<b>LIST OF FIGURES</b> .....	ix
<b>LIST OF ABBREVIATIONS</b> .....	x
<b>LIST OF SYMBOLS</b> .....	xi
<b>1 Introduction</b> .....	1
<b>2 Related work</b> .....	5
<b>3 Method</b> .....	8
3.1 Tools .....	8
3.1.1 NLP Toolkits .....	8
3.1.2 WordNet .....	9
3.2 Text processing .....	10
3.2.1 Tokenization .....	11
3.2.2 Stopping .....	11
3.2.3 Stemming/Lemmatization .....	12
3.2.4 Part of Speech Tagging .....	13
3.2.5 Shallow Parsing .....	14

3.2.6	Dependency Parsing . . . . .	15
3.2.7	Named entity detection . . . . .	15
3.3	Design overview . . . . .	16
3.4	Document type detection . . . . .	18
3.5	Feature extraction . . . . .	18
3.5.1	Shallow features . . . . .	18
3.5.2	Morphological features . . . . .	19
3.5.3	Structural features . . . . .	20
3.5.4	Semantic features . . . . .	20
3.5.5	Social features . . . . .	20
3.5.6	Metadata features . . . . .	21
3.6	Fusioning Strategy . . . . .	21
<b>4</b>	<b>Evaluation . . . . .</b>	<b>22</b>
4.0.1	Datasets . . . . .	22
4.0.2	Metrics . . . . .	22
4.0.3	Overall Assessment . . . . .	23
<b>5</b>	<b>MRAS in Action . . . . .</b>	<b>25</b>
<b>6</b>	<b>Conclusions . . . . .</b>	<b>26</b>
	<b>REFERENCES . . . . .</b>	<b>27</b>

## LIST OF TABLES

3.1	Comparison of existing NLP tools . . . . .	10
3.2	Tokenization of sentences in different languages. . . . .	11
3.3	Example of lemmatization for different languages. . . . .	13
3.4	Part of Speech tagging of the sentence “Did they win the race?” in multiple languages. . . . .	13
3.5	Named entity detection of the sentence “Salvador Dali was born in the city of Figures” . . . . .	16
4.1	Data resources identified for MRAS development and validation . . . . .	23



## LIST OF FIGURES

3.1	Example of a reduced tree generated around the vehicle concept in Wordnet.....	10
3.2	Shallow parsing of the sentence “Did they win the race?” .....	14
3.3	Dependency parsing of the sentence “Did they win the race?” .....	16
3.4	MRAS .....	17

## LIST OF ABBREVIATIONS

**EN** – English

**EU** – Basque

**ES** – Spanish

**IR** – Information Retrieval

**NE** – Named Entity

**NLP** – Natural Language Processing

**PoS** – Part of Speech

## LIST OF SYMBOLS

$\sqrt{2}$  square root of 2

$\lambda$  lambda symbol, normally used in lambda calculus but it sometimes gets used for wavelength as well

## CHAPTER 1

### INTRODUCTION

Reading is an important skill in the academic environment, a competence that can be critical for students' educational opportunities and their careers [49]. As reported by Lennon and Burdick [40] reading for learning takes place when the reader comprehends 75% of a text. This represents an appropriate balance that allows the reader to positively understand the text, while also finding challenges in the reading process that will motivate him to improve his skills [40]. Outside the educational environment, reading generally takes place for comprehension rather than for learning. In this context, it is critical to provide people with texts they can fully understand. For example, patients that properly understand documents disclosed to them before surgery, are known to be less anxious before the operation and obtain more satisfactory results during posterior treatment [48]. However, recent studies[38, 48, 46] show that even medical documents that are supposed to be suited for average readers, tend to be too specialized and even well-educated adults have trouble understanding them. Whether for learning or understanding, the complexity of texts to be read needs to be determined.

Every reader has different reading skills and the levels of difficulty of the texts they need depends also upon their personal objective. Therefore, providing institutions and readers with tools that can measure the complexity of a text so that they can assess

whether it is adequate for a user is imperative. *Readability Assessment (RA)* tools<sup>1</sup> are certainly aimed for handling such a task by providing a mean to determine the degree of ease with which a reader can understand a given text, i.e. the *Readability Score (RS)* of the text.

Historically, teachers have been the main stakeholders of RA formulas, using them to select new materials for their courses and curriculum design. However, lately, more stakeholders have found benefits in using RA tools outside the academic environment. Automatic text simplification[52, 50], summarization for people with reading difficulties[33], book recommendation [47], literacy assessment[54], or legal and medical document complexity assessment[38, 43, 48, 46] are only a few examples of applications that take advantage of the complexity levels generated by existing RA tools. Even in commercial environments, book publishers require professional linguistic services in order to tag their publications with a readability level required for their intended audience, a task that could similarly be completed by an automatic tool.

In estimating the complexity of texts, traditional formulas, such as Flesh [35], became very popular in the late 1940's among educators for manually determining text difficulty. Most of these formulas relied on *shallow features*, which could easily be adapted to multiple languages and provide a simple way of determining text complexity. The multilingualism achieved by traditional formulas offered numerous benefits in contexts where the readability of more than one language was needed, i.e., book translation or learning a second language. However, traditional formulas were known to lack precision. For example, they could classify nonsense text as *simple to read*, just because it contained short and frequently-used words [30]. The insufficient

---

<sup>1</sup>RA tool and RA formula are used interchangeably in this document.

precision encouraged researchers to study and develop better and more sophisticated methods for RA that depended upon more in-depth text analysis [36, 19]. These new formulas continued taking advantage of shallow features, but incorporated more complex features based on the syntax and semantics of text. With the addition of new text complexity indicators, the tools became more precise, but at the same time more constrained regarding their language adaptability [23, 34]. In fact, they used increasingly more language-dependent techniques, which made the systems unadaptable to estimate RS for text in languages other than the one they were designed for. As a result, the multilingualism that was possible in the early stages disappeared.

With multilingualism and precision in mind, we propose to develop **MRAS**, a **Multilingual Readability Assessment System**. This tool should both show results comparable to monolingual state-of-the-art systems and maintain the multilingualism the early tools in the RA field had. For doing so, we will (1) explore features and methods used in literature, (2) design novel features that positively influence the readability level estimating process and (3) analyze how all those features can be adapted to be used in multilingual RA. MRAS will be *open source* and *easily connected* to different applications that require RA as a service, potentially permitting the analysis of all sorts of texts, including text snippets, books, websites and even short and unstructured texts, such as the ones found in social media. In doing so, we will create a system that will adapt itself to the input text language and use an adequate subset of features for the corresponding language for readability prediction, creating, to the best of our knowledge, the first multilingual readability assessment system.

As a byproduct of our research work, we will create a leveled dataset with readability-labeled documents for different languages, which currently is unavailable. In addition, we will create an in-depth report surveying existing strategies for readability predic-

tion.

It is important to note that, for practical purposes, the proposed application will only be tested in three different languages: *English*, for state of the art comparison purposes and as reference of germanic languages. *Spanish*, as a reference for romance languages, and *Basque* as an example of a pre-indoeuropean and minority language.

## CHAPTER 2

### RELATED WORK

From the past six decades, different RA systems have been developed with high diversity in terms of both languages and features [34, 23]. Initial readability formulas, such as Flesh [35], Dale-Chall [26], and Gunning FOG [16] made use of **shallow features**, mostly based on ratios of characters, terms, and sentences. These formulas, were basic enough even to be computed manually, providing a simple way of estimating a text’s complexity, even if the formulas lacked precision in some cases [30]. This simplicity, however, made them easy to be adapted to estimate readability scores in different languages [51].

In recent years, readability formulas have evolved to supervised learning based systems that use a combination of traditional shallow features and new natural language processing based ones, which consider language aspects, such as syntax or semantics of texts. However, incorporating new features has brought a drawback to the area, evidenced by the fact that current systems are too focused in certain languages, making them only functional in the languages they were created for. Current state-of-the-art is composed by methods focused on specific languages, as discussed below:

For **English**, the RA system presented in [19] predicted only two levels of difficulty, simple or complex, using elaborated features, such as ambiguity among the terms in



the texts. Other authors [33], oriented their system for assessing the difficulty level of a text for people with intellectual disabilities by developing features that were intended to detect how well a text was structured. A readability prediction system for financial documents was presented in [25], which was based on features such as the presence of active voice or number of hidden verbs. It is also important to mention two commercial RA tools, Lexile<sup>1</sup> and AR<sup>2</sup>, which are widely used among English speaker academic professionals. Even if their algorithms are not public, they are known to use shallow features showing how common terms of a text are and how long sentences are in average [40]. The literature pertaining to RA for text in English is abundant. For more in-depth discussion on RA formulas refer to [34, 23].

In contrast to English, **Spanish** RA has not seen any significant improvement regarding features in recent years, as most of the existing works are still based on shallow features. Among the well-known RA tools for Spanish, SSR [51] was based on the analysis of sentence length and number of rare words per sentences, whereas LC and SCI [21] were based on density of low frequency words in text. Other systems [53, 32] presented strategies to combine the aforementioned methods to improve RA estimation.

Compared to other languages, **Basque** RA is reduced to only one system. Due to the fact that Basque is considered a minority language and shares little similarity with most spoken languages, limited research has been done in the area. So far, ErreXail [37] is the only system created for Basque RA. ErreXail was developed to predict two different readability values, simple or complex, using features mostly based on ratios of common natural language processing labels, such as Part-of-Speech tags or

---

<sup>1</sup><https://www.lexile.com/>

<sup>2</sup><http://www.renaissance.com/products/accelerated-reader/atos-analyzer>

morphology annotations.

Similar to Basque, the literature for **Arabic** RA is limited as well. Al-Ajlan and Al-Khalifa [15] developed a RA tool based on only two features: average letters per term and average terms per sentence. These features were analyzed using a Support Vector Machine classifier in order to classify text as simple or complex.

Opposed to RA tools for previous languages, structural features do not seem to have such a success for **Chinese** RA. Therefore, most of the research works related to Chinese RA have been focused only on lexical features, such as Tf-Idf of terms [27, 28].

In contrast to the aforementioned techniques, the authors of [31] presented a RA system for **Italian** aimed at assessing readability at sentence level, which combined traditional, lexical, and syntactical methods.

Rather than focusing on the general reader, François and Fairon [36] developed a RA system for **French** with foreign language learners in mind. The objective was to determine which features were more important for a foreign language learner to understand a text. They tested lexical, syntactical and semantic features and showed that semantic ones performed poorly in their case.

Even if the number of RA systems that tackle individual languages is high, they are usually focused on a specific set of features and materials they can analyze. In addition, to best of our knowledge, none the RA systems presented are **multilingual**. MRAS will not only be multilingual, but will also be based on a comprehensive set of existing and novel features which will be general enough to potentially be able to handle all sorts of reading materials. All those characteristics will make MRAS a unique system in the area.

## CHAPTER 3

### METHOD

MRAS is a readability assessment system capable of automatically predicting the readability level of any given document, being able to handle multiple document types, differing in format, length and language. For doing so MRAS takes advantage of several tools described in section 3.1, as well as various text processing strategies described in 3.2. This chapter also provides an overview of the design of MRAS in section 3.3 and further in-depth detail in sections 3.5 and 3.6.

### 3.1 Tools

Whether for text processing or for feature extraction, MRAS takes advantage of several existing tools and techniques, which we describe below.

#### 3.1.1 NLP Toolkits

For developing MRAS we analysed several existing NLP toolkits. The goal of this analysis was to find a suitable NLP toolkit than can work in a widest variety of languages with enough functionalities so that MRAS could be developed.

We narrowed our search to the 6 toolkits presented in table 3.1. NLTK, OpenNLP and CoreNLP are tools that were originally designed for English, with the possibility of extension for other languages. However, even if they can be extended, they require

models specifically designed for them, with very limited languages available currently. On the other side, SyntaxNet adopts an standards based approach, being compatible with the Universal Dependency Treebank [1] which contains model for more than 70 languages. However, given the recency of its development, it does not offer support for some common NLP tasks that MRAS requires yet.

Freeling is a more consolidated project offering a good balance between languages and functionalities. Freeling includes tokenization, PoS tagging, syntactic parsing, dependency parsing and semantic labelling, and supports 14 languages: Asturian (as), Catalan (ca), German (de), English (en), French (fr), Galician (gl), Croatian (hr), Italian (it), Norwegian (nb), Portuguese (pt), Russian (ru), Slovene (sl), Spanish (es), and Welsh (cy).

Katea is the only set of NLP tools developed for Basque. Katea is composed by Morpheus [14] (morpho-syntactic analysis), eustagger [18] (lemmatization and syntactic function indentification), eihera [17] (named entity detection), ixati [14] (shallow parsing) and maltixa [22] (dependency parsing).

Showing support for the three languages analyzed in this thesis (English, Spanish and Basque), SyntaxNet could be the perfect tool for developing it. However, given that it still does not offer some basic features we decided to discard it from MRAS for now. However, this tool might be of interest in the near future. In order to implement MRAS we decided to use two tools that complement each other, Freeling and Katea. Together they fullfill both the language and functionality needs of MRAS.

### 3.1.2 WordNet

Wordnet [42] is a lexical database for English where terms, i.e., nouns, verbs and adjectives, are grouped into sets of synonyms, expressing different concepts. These

Tool	Lang.	Tok.	Stem. \Lemm.	PoS	Shallow Par.	Dependency	NE detection
NLTK [24]	English +	Yes	No	Yes	Yes	In dev.	Poor
CoreNLP [41]	English +	Yes	Some lang.	Yes	Yes	Some lang.	No
OpenNLP [2]	English +	Yes	Yes	Yes	Yes	Yes	Yes
SyntaxNet [3, 20]	70+ lang.	Yes	No	Yes	No	Yes	No
Freeling [44, 45]	14 lang.	Yes	Yes	Yes	Yes	Yes	Yes
Katea [14, 18, 22]	Basque	Yes	Yes	Yes	Yes	Yes	Yes

Table 3.1: Comparison of existing NLP tools

concepts are related by each other by several semantic relationships, such as hyperonymie or hyponimie. An example of this structure can be seen in figure 3.1, where *motor vehicle* has *vehicle* as hypernym and *car* and *truck* as hyponyms. Note that this structure forms a tree, permitting basic operations such as distance or similarity measuring between concepts.

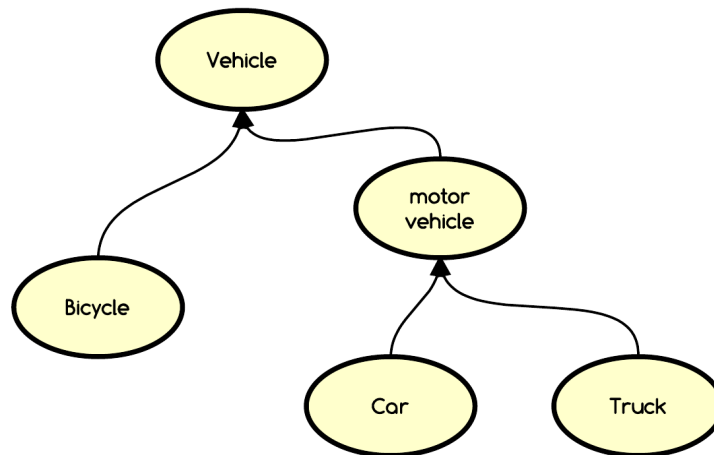


Figure 3.1: Example of a reduced tree generated around the vehicle concept in Wordnet.

## 3.2 Text processing

Natural language text is generally unstructured, making it hard to handle for computers. The aim of text processing is to infer inherent structures in text in order

to ease later processes. This section describes the text processing operation used in MRAS.

### 3.2.1 Tokenization

Tokenization is the process of splitting a text into smaller parts, i.e. tokens. A token represents each sensical part of a text, which usually corresponds to a term, a number, or a punctuation mark.

Even if tokenization may look trivial, the reality proofs that is not. Looking at the examples in Table 3.2, it can be seen how tokenization needs to go further than splitting a sentence by spaces. Both the words *aren't* and *student's* need to be split into two tokens, while the dot following the year 1991 needs to be keep together to show ordinality in Basque numbers.

A precise tokenization is important for most tasks, since most of the processes described later, require of a previous tokenization step.

1991. urtean jaio zen.						
1991.	urtean	jaio	zen	.		
La O.T.A.N. ha firmado el acuerdo.						
la	O.T.A.N.	ha	firmado	el	acuerdo	.
Those aren't student's tables.						
those	are	n't	student	's	tables	.

Table 3.2: Tokenization of sentences in different languages.

### 3.2.2 Stopping

Stopping or stopword removal refer to the process of removing stopwords from a text. A stopword is a term that does not add any important information to the task that is performed, usually adding unnecessary noise that hinders valuable information among

a document. The frequency of a term is a good indicator for stopwords since, usually the terms that most appear in a text are the ones that less information contain. Some examples of general purpose stopwords are *a*, *the* or *is*, however depending on the domain, terms such as *computer* can also be stopwords given its high frequency. The purpose stopping is usually 2-way, speeding up later processes and noise reduction, and is usually performed without the need of any specific tool, just using a stopwords list.

### 3.2.3 Stemming/Lemmatization

The goal of both stemming and lemmatization is to achieve a normalized version of a term. This normalization, is usually helpful for search and comparison task as it reduced the search space among all the terms. As an example, when computing a metric about the verb *play* it might be interesting to compute it for all its word-forms (*play*, *plays*, *played*). This process is usually simpler if all the word-forms are normalized to one canonical form. Stemming and Lemmatization differ in the way the normalized form is obtained. While lemmatization is able to achieve real canonical form (i.e. lemma) of a word (the one appearing in the dictionary), stemming simply chops common prefixes and suffixes to obtain an approximation of it. When both techniques are available MRAS will use lemmatization over stemming, however, some languages do not have any lemmatization technique available. Freeling will be used for lemmatization in Spanish and English, while katea will serve for the same purpose in Basque. Several examples of lemmatization can be seen in Table 3.3.

	English			Spanish			Basque		
<b>Original</b>	plays	are	won	fuiste	vine	comieron	nator	balituzte	dakit
<b>Lemmatized</b>	play	be	win	ir	venir	comer	etorri	edun	jakin

Table 3.3: Example of lemmatization for different languages.

<b>English</b>	Did	they	won	the	race	?
	Verb (Auxiliar)	Pronoun	Verb	Determinant	Noun	Symbol
<b>Spanish</b>	Ganaron	(ellos)	la	carrera		?
	Verb	Pronoun (Elliptic)	Determinant	Noun		Symbol
<b>Basque</b>	Lasterketa	irabazi	al	zuten	(haiek)	?
	Noun	Verb	Particle	Verb (Auxiliar)	Pronoun (Elliptic)	Symbol

Table 3.4: Part of Speech tagging of the sentence “Did they win the race?” in multiple languages.

### 3.2.4 Part of Speech Tagging

Part of Speech (PoS) Tagging is the process of labelling each token with a tag that represent the function each token has in a sentence. PoS tags usually differ from language to language<sup>1</sup>, however, the most predominant tags, such as verb, adjective or noun, exist among all the languages. As it is shown Table 3.5, for the sentence ”Did they won the race?”, most of the PoS tags used in the 3 languages are similar, despite slight differences, such as, the presence of auxiliary verbs or particles.

PoS tags are the building block of structure in text, since they define the first organizational layer on it. They can be used alone for applications that need to consider very basic structure in text or need to better identify terms (a term can have different meaning depending its PoS)**TODO Cannot come up with an example.** However PoS tagging usually constitutes a preliminary step used as part of further text analysis processes.

---

<sup>1</sup>As an example, the Penn Treebank project defines 36 PoS tags for English, which can be seen here [https://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn.treebank.pos.html](https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn.treebank.pos.html)



### 3.2.5 Shallow Parsing

Shallow parsing, also called chunking, refers to the process of grouping tokens into chunks. A chunk usually consists of a small phrase of about 1 to 4 terms. The terms in each chunk are somehow connected to each other and together express a senseful concept. There are two types of chunks, depending if they express a noun or a verb phrase. An example of a shallow parsing of a sentence can be seen in Figure 3.2. The example is only shown for English since this shallow parsing works the same way among most languages.

Shallow parsing is interesting for tasks that require to use bigger chunks of text than just tokens. One simple example could be autocompletion, where full senseful chunks can be suggested to the user instead on just separate tokens. Named entity recognition also benefits from shallow parsing, since for most of the Named Entities are usually exact noun-chunks.

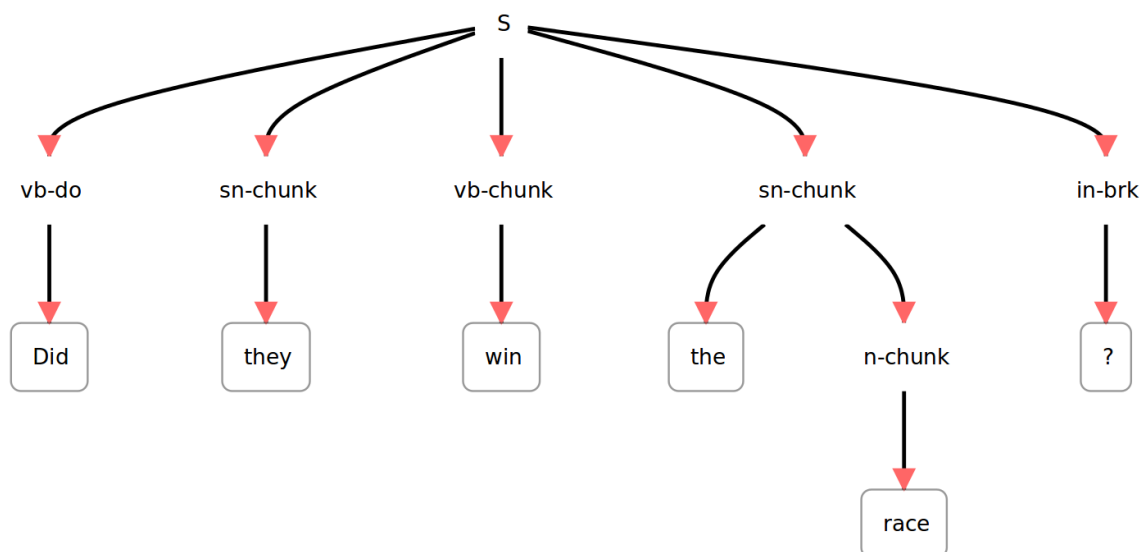


Figure 3.2: Shallow parsing of the sentence “Did they win the race?”

### 3.2.6 Dependency Parsing

Dependency parsing goes further than shallow parsing, determining relationships between tokens rather than just grouping them. Given these relationships, a dependency tree is generated, which usually has a root node representing the main verb of the sentence, which has the subject and objects of the sentence as children. An example of a dependency parsed sentence can be seen in Figure 3.3. Only one language is presented in the example since dependency parsing is very similar for the 3 languages.

One of the benefits of dependency parsing over shallow parsing is the possibility to select a granularity when selecting sensical chunks. Using the tree generated by the dependencies, any cut will cause the sub-trees to be fully sensical. As an example, cutting the tree in the node *did*, would create sensical chunks *they* and *win the race* while cutting the tree in the node *win* would create *the race*. Because of this possibility, dependency parsing is highly used in language transformation applications, e.g. transforming a statement into a question.

### 3.2.7 Named entity detection

A named entity is a token or group of tokens that represent and known entity such as a person, a location, or an organization. Depending on the complexity of the tool that performs this analysis, those entities can also be linked to a knowledge base such as Dbpedia [39] where more structured information about the entity can be found.

A Named Entity is usually a relevant term in a text, something that a user may be interested or looking for. For this reason, it is common give preference to Named Entities, or consider them as key points for a document, when searching or indexing documents.

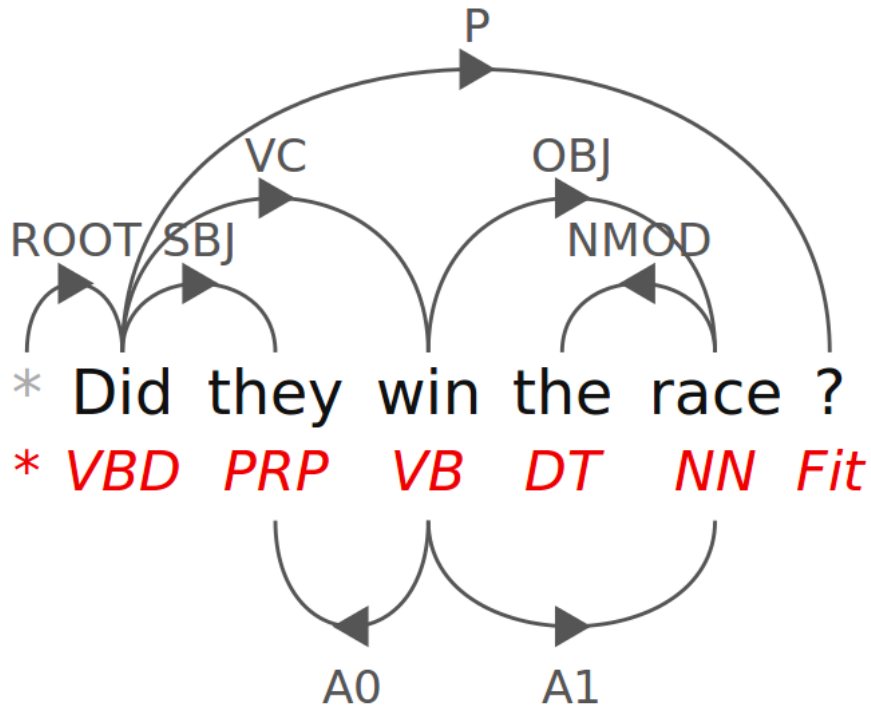


Figure 3.3: Dependency parsing of the sentence “Did they win the race?”

Salvador	Dali	was	born	in	the	city	of	Figures	.
person	person							location	

Table 3.5: Named entity detection of the sentence “Salvador Dali was born in the city of Figures”

### 3.3 Design overview

MRAS is based on a supervised learning approach that relies on knowledge acquired from a leveled corpora. In designing MRAS we followed the steps illustrated in Figure 3.4 and discussed below.

MRAS receives two different inputs: a collection  $C$  of documents for each of which a readability label is already assigned, and a document  $d$  which its readability is unknown for MRAS and, thus, will be predicted. Both inputs are taken through a

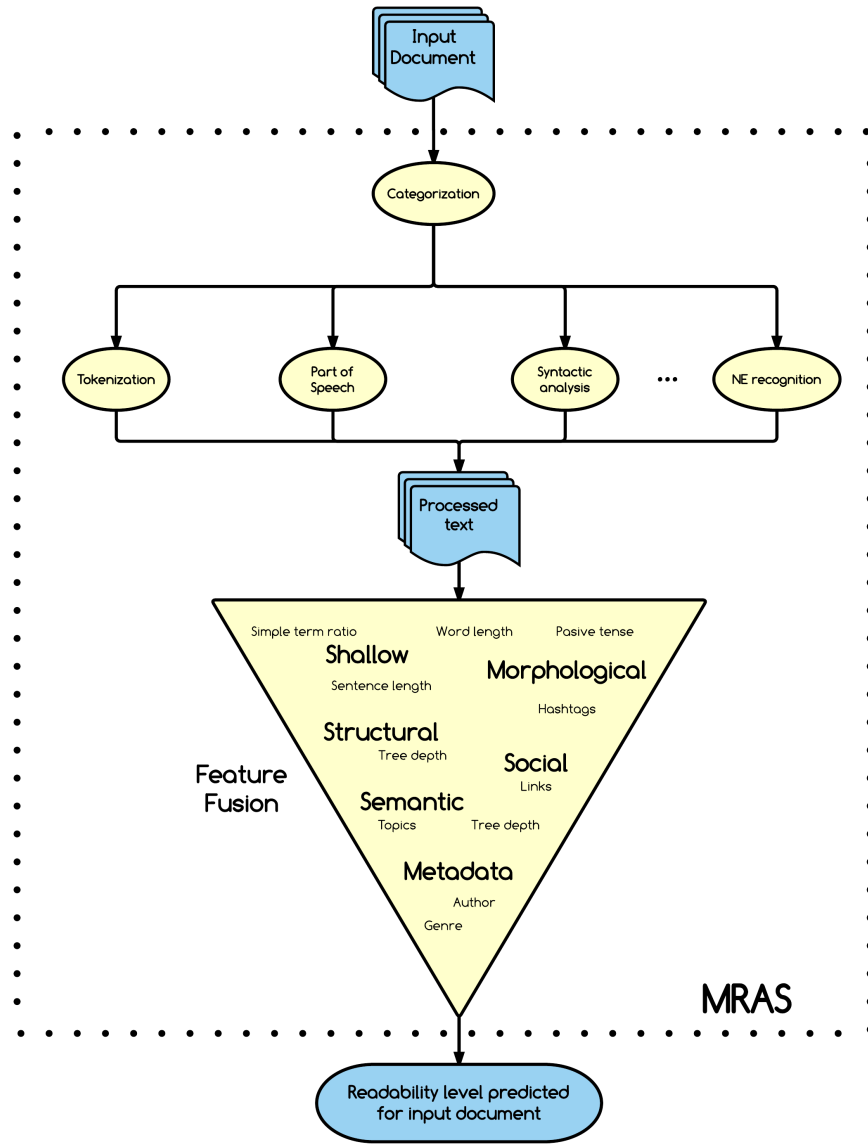


Figure 3.4: MRAS

preprocessing step described in section ??, which cleans, filters and normalizes their content, to prepare them for the feature extraction step described in section ??. These features, serve as a numeric representation for each document. MRAS is capable of learning patterns over the representations extracted from each document in  $C$  and

use these patterns to predict readability scores for new unlabelled documents such as *d*.

### 3.4 Document type detection

One of the main features of MRAS is its versatility, since MRAS is capable to predict readability values for documents of different format, length and language. Consequently, not all documents handled by MRAS can be treated the same way.

I need to figure out how we will do this, before writing further

### 3.5 Feature extraction

Feature engineering is one of the most important aspects of this thesis. A good feature set determines the the quality of a classifier, and therefore the quality of a readability assessment system. A description of each feature included in MRAS is provided below separated in different linguistic levels.

#### 3.5.1 Shallow features

Shallow features [35, 26, 16] have historically shown to be of good use when predicting readability. Even if they sometimes lack precision [30], they serve as a good baseline for readability assessment systems. A description for each shallow feature is provided below.

#### Word length

Everyday terms are usually short in most languages, they are preferred for spoken language over their longer analogues, in order to maintain more fluent conversations.

In the opposite side, difficult terms are longer the more technical or scientific they are. Therefore, short terms the one youngsters first learn and better comprehend. To take advantage of this fact, MRAS bases 4 features on it, average length of terms and average length of their lemmas, considering or not stopwords for both cases.

### **Sentence length**

Sentences oriented people with low understanding skill are usually short. They just focus on simple ideas and facts with very simple argumentations. On the opposite side, documents oriented to more technical and complex aspects contain more argumentation and therefore more subclauses in the sentences, resulting in longer sentences. MRAS benefits of this fact generating two features, average word length of sentence with and without stop words.

### **Simple terms**

#### **3.5.2 Morphological features**

Morphological features capture how terms are formed from their root. Even if this aspect is not relevant in some languages, such as English, it has been shown to be a strong predictor for readability scores in morphology rich languages such as Basque [37]. A description of morphological features included in MRAS is described below.

### **Morphological phenomena frequencies**

#### **Inflection ratio**

difference in size between the wordform and the lemma

### 3.5.3 Structural features

Structural features are the ones that describe how a text is organized. They can both describe structure within the sentence (syntactical structure) or structure between sentences (pragmatical structure). The features that have been implemented in MRAS can be seen below.

#### Structural complexity

based on a neural network and distributional semantics

### 3.5.4 Semantic features

Semantic features go beyond the tokens and structure of the text in order to analyse the concepts laying on it.

#### Semantic closeness

How close is a text from the terms of each group? Using simple cosine similarity and using the cloud of points created using distributional semantics.

#### Concept densisty

#### Followability

#### Ngram frequencies

### 3.5.5 Social features

RA can be used in more than just plain text. Internet is evolving into a new social era and so are text resources. Increasingly more resources contain social information,

such as hashtags, mentions, or links, a type of information that is usually ignored by readability formulas.

### **3.5.6 Metadata features**

Metadata based features can be useful in environments where text access is limited (i.e. copyrighted material). An exploration of this type of features will also be done in order to expand the types of texts MRAS can handle.

## **3.6 Fusioning Strategy**



## CHAPTER 4

### EVALUATION

Even if MRAS is designed to be language independent, for practical purposes the evaluation will only be conducted in three languages that we think can faithfully represent the diversity of existing languages. For this purpose, we have chosen a germanic, a romance, and a pre-indioeuropean language, i.e. English, Spanish, and Basque respectively.

#### 4.0.1 Datasets

The ideal dataset for developing MRAS would be a multilingual leveled dataset that would contain the exact same documents written in different languages, as well as human judgments, in terms of readability scores for each document. However, to the best of our knowledge, such a dataset does not currently exist. Consequently, we have identified various sets of leveled documents for each individual language that can suit MRAS' needs and can be used for evaluation purposes. Details on the datasets considered for evaluation purposes can be seen in Table 4.1.

#### 4.0.2 Metrics

The performance of MRAS will be evaluated by means of (1) common classification evaluation methods, such as absolute error [29], (2) regression evaluation methods

	Dataset	Description
<b>English</b>	Lexile [4]	Contains book titles associated with its readability level
	Stantarized tests [5, 6]	Tests for English level, they contain various texts per test
	Other [7, 8, 9]	News for kids, exercises for learning English
<b>Spanish</b>	Lexile [4]	Contains book titles associated with its readability level
	Learning resources [10, 11, 12]	Various exercises for learning Spanish
<b>Basque</b>	Learning resources [4]	Various exercises for learning Basque
<b>Multilingual</b>	Parallel corpus [13]	Contains same texts translated into Spanish and English

Table 4.1: Data resources identified for MRAS development and validation

such as MSE (Mean Square Error) [29] and (3) methods common in the readability assessment domain, such as adjacent accuracy [36].

#### 4.0.3 Overall Assessment

The study and performance analysis of this thesis will aim at answering the following questions:

- Which learning model performs better for MRAS? Which feature subset?
- Which features add more value in terms of predicting readability? Do they add same value for each language?
- How does MRAS perform compared to baseline shallow feature based formulas? and compared to state of the art systems?
- Would MRAS give the same prediction for a text that is translated manually into another language? and for a text that is automatically translated?
- How efficiently can MRAS predict the readability levels of written text in a language for which it has not been trained? If we train MRAS for two languages can we use it to predict the readability of a text in a third one?

- If we have a really small dataset for one language, would adding more data from another language improve the prediction results of the first one?

## CHAPTER 5

### MRAS IN ACTION

hashtag rec in twitter

## CHAPTER 6

## CONCLUSIONS

## REFERENCES

- [1] <http://universaldependencies.org/>.
- [2] <http://opennlp.apache.org>.
- [3] <https://github.com/tensorflow/models/tree/master/syntaxnet>.
- [4] <http://www.ikasbil.eus>.
- [5] <http://www.flo-joe.co.uk/exams.htm>.
- [6] [http://learnenglishteens.britishcouncil.org/content?field\\_language\\_level\\_tid=50&field\\_section\\_tid=1129&field\\_topics\\_tid=&language=en](http://learnenglishteens.britishcouncil.org/content?field_language_level_tid=50&field_section_tid=1129&field_topics_tid=&language=en).
- [7] <https://www.readinga-z.com/books/leveled-books/>.
- [8] <http://www.breakingnewsenglish.com/news-for-kids.html>.
- [9] <http://www.newsinlevels.com/>.
- [10] <http://cvc.cervantes.es/aula/lecturas/>.
- [11] <http://aprenderespanol.org/lecturas/lecturas-ejercicios.html>.
- [12] [http://www-k6.thinkcentral.com/content/hsp/reading/Senderos/na/common/online\\\_senderos\\\_libros\\\_graduables\\\_para\\\_lectores/senderos\\\_SE/launch.html](http://www-k6.thinkcentral.com/content/hsp/reading/Senderos/na/common/online\_senderos\_libros\_graduables\_para\_lectores/senderos\_SE/launch.html).
- [13] <http://albalearning.com/audiolibros/textosparalelos.html>.
- [14] Itziar Aduriz, Maxux J Aranzabe, Jose Maria Arriola, Arantza Díaz de Ilarraza, Koldo Gojenola, Maite Oronoz, and Larraitz Uria. A cascaded syntactic analyser for basque. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 124–134. Springer, 2004.
- [15] Amani A Al-Ajlan, Hend S Al-Khalifa, and A Al-Salman. Towards the development of an automatic readability measurements for arabic language. In *Digital Information Management, 2008. ICDIM 2008. Third International Conference on*, pages 506–511. IEEE, 2008.

- [16] Judith Albright, Carol de Guzman, Patrick Acebo, Dorothy Paiva, Mary Faulkner, and Janice Swanson. Readability of patient education materials: implications for clinical practice. *Applied Nursing Research*, 9(3):139–143, 1996.
- [17] Inaki Alegria, Olatz Arregi, Irene Balza, Nerea Ezeiza, Izaskun Fernandez, and Ruben Urizar. Design and development of a named entity recognizer for an agglutinative language. In *First International Joint Conference on NLP (IJCNLP-04). Workshop on Named Entity Recognition*, 2004.
- [18] Iñaki Alegria, Xabier Artola, Kepa Sarasola, and Miriam Urkia. Automatic morphological analysis of basque. *Literary and Linguistic Computing*, 11(4):193–203, 1996.
- [19] Sandra Aluisio, Lucia Specia, Caroline Gasperin, and Carolina Scarton. Readability assessment for text simplification. In *Proceedings of the NAACL HLT 2010 Fifth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 1–9. Association for Computational Linguistics, 2010.
- [20] Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta, Kuzman Ganchev, Slav Petrov, and Michael Collins. Globally normalized transition-based neural networks. *arXiv preprint arXiv:1603.06042*, 2016.
- [21] Alberto Anula. Tipos de textos, complejidad lingüística y facilitación lectora. In *Actas del Sexto Congreso de Hispanistas de Asia*, pages 45–61, 2007.
- [22] Kepa Bengoetxea and Koldo Gojenola. Application of different techniques to dependency parsing of basque. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, pages 31–39. Association for Computational Linguistics, 2010.
- [23] Rebekah George Benjamin. Reconstructing readability: Recent developments and recommendations in the analysis of text difficulty. *Educational Psychology Review*, 24(1):63–88, 2012.
- [24] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python*. ” O’Reilly Media, Inc.”, 2009.
- [25] Samuel B Bonsall, Andrew J Leone, and Brian P Miller. A plain english measure of financial reporting readability. *Available at SSRN 2560644*, 2015.
- [26] Jeanne Sternlicht Chall and Edgar Dale. *Readability revisited: The new Dale-Chall Readability Formula*. Brookline Books, 1995.

- [27] Yaw-Huei Chen, Yi-Han Tsai, and Yu-Ta Chen. Chinese readability assessment using tf-idf and svm. In *Machine Learning and Cybernetics (ICMLC), 2011 International Conference on*, volume 2, pages 705–710. IEEE, 2011.
- [28] Kevyn Collins-Thompson and James P Callan. A language modeling approach to predicting reading difficulty. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 193–200, 2004.
- [29] W Bruce Croft, Donald Metzler, and Trevor Strohman. *Search engines: Information Retrieval in Practice*, volume 283. Addison-Wesley Reading, 2010.
- [30] Alice Davison and Robert N Kantor. On the failure of readability formulas to define readable texts: A case study from adaptations. *Reading Research Quarterly*, pages 187–209, 1982.
- [31] Felice Dell’Orletta, Simonetta Montemagni, and Giulia Venturi. Read-it: Assessing readability of italian texts with a view to text simplification. In *Proceedings of the Second Workshop on Speech and Language Processing for Assistive Technologies*, pages 73–83. Association for Computational Linguistics, 2011.
- [32] Biljana Drndarević, Sanja Štajner, Stefan Bott, Susana Bautista, and Horacio Saggion. Automatic text simplification in spanish: a comparative evaluation of complementing modules. In *Computational Linguistics and Intelligent Text Processing*, pages 488–500. Springer, 2013.
- [33] Lijun Feng. Automatic readability assessment for people with intellectual disabilities. *ACM Special Interest Group on Accessible Computing*, (93):84–91, 2009.
- [34] Lijun Feng, Martin Jansche, Matt Huenerfauth, and Noémie Elhadad. A comparison of features for automatic readability assessment. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 276–284. Association for Computational Linguistics, 2010.
- [35] Rudolph Flesch. A new readability yardstick. *Journal of Applied Psychology*, 32(3):221, 1948.
- [36] Thomas François and Cédric Fairon. An ai readability formula for french as a foreign language. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 466–477. Association for Computational Linguistics, 2012.
- [37] Itziar Gonzalez-Dios, Maria Jesús Aranzabe, Arantza Díaz de Ilarraza, and Haritz Salaberri. Simple or complex? assessing the readability of basque texts. In



- Proceedings of International Conference on Computational Linguistics*, volume 2014, 2014.
- [38] Ahmed MS Ibrahim, Christina R Vargas, Pieter GL Koolen, Danielle J Chuang, Samuel J Lin, and Bernard T Lee. Readability of online patient resources for melanoma. *Melanoma Research*, 26(1):58–65, 2016.
  - [39] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, and Chris Bizer. DBpedia - a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web Journal*, 2014.
  - [40] Colleen Lennon and Hal Burdick. The lexile framework as an approach for reading measurement and success. *Electronic publication on [https://cdn.lexile.com/m/resources/materials/Lennon\\_Burdick\\_2004.pdf](https://cdn.lexile.com/m/resources/materials/Lennon_Burdick_2004.pdf)*, 2004.
  - [41] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *ACL (System Demonstrations)*, pages 55–60, 2014.
  - [42] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
  - [43] James RP Ogloff and Randy K Otto. Are research participants truly informed? readability of informed consent forms used in research. *Ethics & Behavior*, 1(4):239–252, 1991.
  - [44] Llus Padr, Miquel Collado, Samuel Reese, Marina Lloberes, and Irene Castelln. Freeling 2.1: Five years of open-source language processing tools. In *Proceedings of 7th Language Resources and Evaluation Conference (LREC’10)*, La Valletta, Malta, May 2010.
  - [45] Llus Padr and Evgeny Stanilovsky. Freeling 3.0: Towards wider multilinguality. In *Proceedings of the Language Resources and Evaluation Conference (LREC 2012)*, Istanbul, Turkey, May 2012. ELRA.
  - [46] Chirag R Patel, Saurin Sanghvi, Deepa V Cherla, Soly Baredes, and Jean Anderson Eloy. Readability assessment of internet-based patient education materials related to parathyroid surgery. *Annals of Otology, Rhinology & Laryngology*, pages 523–527, 2015.
  - [47] Maria Soledad Pera and Yiu-Kai Ng. Automating readers’ advisory to make book recommendations for k-12 readers. In *Proceedings of the 8th ACM Conference on Recommender Systems*, pages 9–16. ACM, 2014.

- [48] Jennifer Petkovic, Jonathan Epstein, Rachelle Buchbinder, Vivian Welch, Tamara Rader, Anne Lyddiatt, Rosemary Clerehan, Robin Christensen, Annelies Boonen, Niti Goel, et al. Toward ensuring health equity: Readability and cultural equivalence of omeract patient-reported outcome measures. *The Journal of Rheumatology*, 42(12):2448–2459, 2015.
- [49] Richard D Robinson, Michael C McKenna, and Judy M Wedman. Issues and trends in literacy education. 2000.
- [50] Horacio Saggion, Sanja Štajner, Stefan Bott, Simon Mille, Luz Rello, and Biljana Drndarevic. Making it simplext: Implementation and evaluation of a text simplification system for spanish. *ACM Transactions on Accessible Computing (TACCESS)*, 6(4):14, 2015.
- [51] Seth Spaulding. A spanish readability formula. *The Modern Language Journal*, 40(8):433–441, 1956.
- [52] Sanja Štajner, Ruslan Mitkov, and Gloria Corpas Pastor. Simple or not simple? a readability question. In *Language Production, Cognition, and the Lexicon*, pages 379–398. Springer, 2015.
- [53] Sanja Štajner and Horacio Saggion. Readability indices for automatic evaluation of text simplification systems: A feasibility study for spanish. In *Proceedings of the 6th International Joint Conference on Natural Language Processing (IJCNLP 2013)*, Nagoya, Japan, pages 374–382, 2013.
- [54] Barry D Weiss, Mary Z Mays, William Martz, Kelley Merriam Castro, Darren A DeWalt, Michael P Pignone, Joy Mockbee, and Frank A Hale. Quick assessment of literacy in primary care: the newest vital sign. *The Annals of Family Medicine*, 3(6):514–522, 2005.

