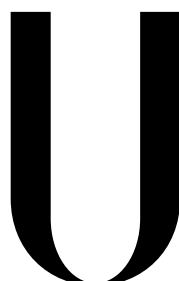


UNIVERSIDADE DE LISBOA
FACULDADE DE CIÊNCIAS
DEPARTAMENTO DE INFORMÁTICA



LISBOA

UNIVERSIDADE
DE LISBOA

Semantic annotation of electronic health records in a multilingual environment

Luís Campos

DISSERTAÇÃO

MESTRADO EM BIOINFORMÁTICA E BIOLOGIA COMPUTACIONAL
ESPECIALIDADE EM BIOINFORMÁTICA

2017

UNIVERSIDADE DE LISBOA
FACULDADE DE CIÊNCIAS
DEPARTAMENTO DE INFORMÁTICA



LISBOA

UNIVERSIDADE
DE LISBOA

**Semantic annotation of electronic health records
in a multilingual environment**

Luís Campos

DISSERTAÇÃO
MESTRADO EM BIOINFORMÁTICA E BIOLOGIA COMPUTACIONAL
ESPECIALIDADE EM BIOINFORMÁTICA

Tese orientada pelo Prof. Doutor Francisco Couto e pelo Dr. Vasco Pedro

2017

Resumo

Palavras Chave: palavras chave

Abstract

Keywords: keywords

Resumo Alargado

Acknowledgements

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objectives	3
1.3	Methodology	4
1.4	Contributions	4
2	Related Work	7
2.1	Text Mining	7
2.1.1	Named-entity Recognition	7
2.1.2	Natural Language Processing	8
2.1.2.1	Tokenization	8
2.1.2.2	Stemming and Lemmatization	9
2.1.3	Application of Text Mining on Radiology Reports	9
2.1.4	Ontologies	10
2.2	Translation	11
2.2.1	Terminology	11
2.2.2	Machine Translation	11
2.2.2.1	Word-Based Models	11
2.2.2.2	Phrase-Based Models	13
2.2.3	Post-editing	14
2.2.4	Machine Translation Services	14
2.2.4.1	Yandex	14
2.2.4.2	Google	15
2.2.4.3	Unbabel	15
2.2.5	Translation of Medical Text	15

CONTENTS

2.2.5.1	Multilingual Text Mining	15
2.2.5.2	Machine Translation of Doctor-Patient Commu- nication	16
2.2.5.3	Machine Translation of Public-Health Information	17
2.2.5.4	Machine Translation for Information Retrieval . .	17
2.2.5.5	Machine Translation of Other Types of Medical Text	18
2.2.6	Translation of Biomedical Lexicons	18
2.3	External Tools and Terminologies	19
2.3.1	RadLex	19
2.3.2	Open Biomedical Annotator	21
2.3.3	NOBLE Coder	21
2.4	Evaluation Metrics	24
2.4.1	Micro- and Macro- Evaluation Metrics	25
3	Framework	27
3.1	MRRAD (Multilingual Radiology Research Articles Dataset) Corpus	27
3.1.1	Web Crawl of the articles (1,2)	28
3.1.2	Note On Human Translations	28
3.1.3	Yandex Translation (3)	29
3.1.4	Google and Unbabel Translation (4,5)	29
3.2	Annotation	30
3.2.1	Direct Match - Annotation with Open Biomedical Annotator	30
3.2.2	All Match and Best Match - Annotation with NOBLE Coder	30
3.3	Evaluation	31
4	Experimental Results	33
4.1	NER Lexicon-based approach	33
4.1.1	Clinical Finding and Anatomical Entity Subtrees	36
5	Conclusions	41
	References	43
	Commands used to run NOBLE Coder	51

CONTENTS

Uploading RadLex ontology to NOBLE	53
Macro Evaluations	55

List of Figures

4.1	Micro Evaluation of Translations being tested (All Match)	34
4.2	Micro Evaluation of Translations being tested (Best Match) . . .	34
4.3	Micro Evaluation of Translations being tested (Direct Match) . .	35
4.4	Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (All Match)	37
4.5	Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (Best Match)	38
4.6	Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (Direct Match)	39
1	Macro Evaluation of Translations being tested (All Match)	55
2	Macro Evaluation of Translations being tested (Best Match) . . .	56
3	Macro Evaluation of Translations being tested (Direct Match) . .	56
4	Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (All Match)	57
5	Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (Best Match)	58
6	Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity sub- trees" (Direct Match)	59

List of Tables

2.1	Lexical translation probability table for the word broken	12
2.2	NOBLE matching strategies present in the GUI interface. Adapted from (Tseytlin <i>et al.</i> , 2016). This correspond to the options used in the GUI tool.	23
3.1	Number of articles by journal in parallel corpus	29
3.2	OBA parameters used	30
4.1	Number of RadLex terms found by document	33

Chapter 1

Introduction

1.1 Motivation

Radiology reports describe the results of radiography procedures and have the potential of being an useful source of information, which can bring benefits to health care systems around the world. However, these reports are usually written in free-text and thus it is hard to automatically extract information from them. Nonetheless, the fact that most reports are now digitally available make them amenable for using Text Mining tools. Another advantage of Radiology reports is that even if written in free-text, they are usually well structured.

A lot of work has been done on Text Mining of Biomedical texts, including health records (Pons *et al.*, 2016), but although Radiology reports are usually written in the native language of the radiologist, Text Mining tools are mostly developed for English. For example, (Hassanpour & Langlotz, 2016) created an information extraction system for English reports that depends on RadLex, a lexicon for radiography terminology, which is freely available in English. Given this dependence, the system cannot be easily applied to reports written in other languages. And even if the system was not dependable on an English lexicon, it is not certain that the results would be the same if text in another language was used, because of, for example, differences in syntax. This creates an obstacle to the sharing of Radiology information between countries, which is important to understand and effectively address health problems.

1. INTRODUCTION

There are mainly two possible solutions to this problem. One is to translate the lexicon itself (Bretschneider *et al.*, 2014) and the other is to translate the reports. Translating the lexicon has the advantage of not requiring continuous translation, i.e., after translating a lexicon to, for example, Spanish, we can then use it to process as many untranslated Spanish reports as needed. However, when a new version of the lexicon is released the changes need also to be translated, otherwise the translated lexicon would become outdated. In this work I explore the alternative option of translating the reports and check its feasibility. This approach has the advantage that the translated reports would be accessible to any doctor who understands English and the state-of-the-art Text Mining tools focused on English text can be applied without any need for adaptation.

If the translation is done by professionals trained in the translation of medical texts, we probably can assume that not much information is lost in translation. We call this type of translation Human Translation (HT). But expert translators are expensive and there are a lot of reports to translate, so this would have a really high monetary cost. Another option is to use Machine Translation (MT). Notwithstanding the lower translation quality, it is way cheaper and more feasible in a large scale. Finally, an option that tries to get the best of both worlds is using Machine Translation with Post-Editing (MT-PE) by humans. In this approach the text is automatically translated by a machine and then the translation is corrected by a human. Cheaper than the HT option and with better quality than the MT one.

The choice of translation approach is important because it will affect the quality of the output of the Text Mining tools. From the best of my knowledge, currently there is no publicly available study that provided a quantitative evidence that would help make this choice. This could be explained by the lack of a parallel corpus that could be used to study this. To the best of my knowledge, the most similar work to this one is (Castilla, 2007). He founds that a rule-based MT system has a good performance in translating Portuguese text to English for the purposes of applying a text mining tool (better described in 2.2.5.1). The author does not compare translation systems, something that is done on the present work.

In this work I focused on the Text Mining task of finding RadLex terms in text, a Named-entity recognition (NER) task. The identification of RadLex terms can be useful, for example, in Image Retrieval ([Gerstmair *et al.*, 2012](#)) systems, so the results of this thesis can have a direct impact on developer’s decisions on how to integrate translation in their applications.

1.2 Objectives

Thus, I aimed at addressing the following research question: lacking the resources to pay for human translation services, what kind of automatic (MT) or semi-automatic translation (MT+PE) approach should be used in the task of translating Portuguese Radiology-related text to English, for the purposes of finding RadLex terms in the translated text? I propose the following hypothesis:

Hypothesis: MT+PE is a good trade-off between quality and cost, compared with MT and HT, for translating Portuguese Radiology reports to English, for the purpose of identifying RadLex terms in the translated text.

For this to be true, these conditions have to hold:

1. MT+PE has to be cheaper than HT
2. The terms extracted from MT+PE translations have to be similar enough to the ones extracted from the HT translation
3. The terms extracted from MT+PE translations have to be more similar to the ones extracted from the HT translation than the ones extracted from MT translations

The first condition is known to be true. The last condition its important because if MT+PE quality is similar to MT quality, as MT cost is lower, then it is worth to just use MT. In this thesis I only try to answer to the quality issues, not doing a thorough economic analysis of the problem.

1. INTRODUCTION

1.3 Methodology

To test this hypothesis I have compared the RadLex terms extracted from MT and MT+PE translations to the terms extracted from HT translations, which I assumed to be a gold standard.

For this purposes I've created MRRAD, a parallel corpus containing 51 Portuguese scientific articles related to Radiology and corresponding HT, MT (Google and Yandex) and MT+PE (Unbabel) English translations. These translations were annotated with RadLex terms using the Open Biomedical Annotator and NOBLE Coder. More than one annotation approach was used to experiment with different kinds annotation approaches. For each translation and annotation approach I created the set of the RadLex terms that were found in that translations with that annotation approach. The terms found in the MT and MT+PE translations were then compared with the ones found in the HT translations.

1.4 Contributions

This thesis lead to the following contributions:

- **MRRAD Corpus**

- A Portuguese-English parallel corpus of research articles related to Radiology, called MRRAD (Multilingual Radiology Research Articles Dataset), containing for each article the original Portuguese document, the HT translation, two alternative MT translations and a MT+PE translation;

- **Main Scientific Results**

- Measurement of the performance of multiple automatic or semi-automatic translation approaches in the task of translating Portuguese Radiology-related text to English, for the purposes of finding RadLex terms in the translated text;

- **Bioinformatics Open Days 2017¹**

- Abstract submission and oral presentation describing this work;
- Organization and presentation of workshop on Biomedical Text Mining with other members of the LaSIGE team²;

- **Scientific Publications**

- ...

- **Other Open-Source Contributions**

- A Python binding of the BioPortal REST API³;
- Converter of NOBLE Coder annotation file to Webanno TSV 2 annotations files⁴;

- **Multilingual Report Annotator**

- Development of a proof of concept web application for translation and annotation of Radiology text;⁵

¹<http://bioinformaticsopendays.com/>

²<https://sites.google.com/view/biomedicaltextminingworkshop>

³<https://github.com/LLCampos/pybiportal>

⁴<https://gist.github.com/LLCampos/5f1680941984c4b63f986965e7384e6c>

⁵<http://www.lasige.di.fc.ul.pt/webtools/mra/>

Chapter 2

Related Work

2.1 Text Mining

Text Mining consists in the machine supported analysis of text (Hotho *et al.*, 2005). It can be used, for example, to help researchers cope with information overload (Cohen & Hersh, 2005) due to the big volume of scientific data in the form of unstructured literature. More related to this thesis, it can also be used to extract information from free-text Radiology reports (Pons *et al.*, 2016).

2.1.1 Named-entity Recognition

Named-entity recognition (NER) is a task of Text Mining that has the goal of locate and classify all the named-entities in a certain document. Named-entities are elements of the text that belong to one of certain predefined classes. For example, in the phrase *Atrial fibrillation has strong associations with other cardiovascular diseases* the term *Atrial fibrillation* is a named-entity that belongs to the class *Disease*. The task of NER can be further divided in two subtasks: identify where in the text are the named-entities and classify the named-entities.

The approaches of NER can be divided into three categories (Mohit, 2014): Rule-based approaches, Machine Learning based approaches and hybrid approaches.

- In rule-based approaches the identification and classification subtasks are based on rules crafted by humans. Usually domain specific.

2. RELATED WORK

- In Machine Learning based approaches the subtasks are turned into classification problems and Machine Learning algorithms are used to identify and classify named-entities. These approaches are easily ported to different domains other than the ones they were originally developed to be applied on.
- Hybrids approaches combines the two last approaches.

Lexicon based-approaches are a subset of the rule-based approaches. In this approach we already have a list of the named-entities (a lexicon) that we want to identify in the text. For example, if we want to identify chemical entities in text, we use a chemical-entities lexicon. The goal of the lexicon based-approach is then to identify, in text, mentions of terms presented in the lexicon. This could be done by direct matching, as implemented by the Open Biomedical Annotator¹ (Jonquet *et al.*, 2009). In this strategy, the system only tries to find in text terms that are also in the lexicon, not considering, for example, lexical variations. The recall can be lower than expected because lexical variants (like plurals), abbreviations and partial matchings of lexicon terms are not recognized in the text. For this purpose, more complex tools like NOBLE Coder² (Tseytlin *et al.*, 2016) or Concept Mapper (Stewart *et al.*, 2012) can be used.

2.1.2 Natural Language Processing

Because Text Mining has to manipulate text, it is not too surprising that it borrows tools from Natural Language Processing (NLP), a research fields that seeks to improve computational understanding of natural language. In the next sections I will briefly explain some of the NLP tasks relevant to this thesis.

2.1.2.1 Tokenization

Tokenization is one of the main tasks of NLP and consists in dividing a certain text in pieces called *tokens*. A token can be defined as "an instance of a sequence of characters in some particular document that are grouped together as a useful

¹<http://bioportal.bioontology.org/annotator>

²<http://noble-tools.dbmi.pitt.edu/>

semantic unit for processing" (Manning *et al.*, 2009b). So, for the sentence *the mother had a surgery*, it is possible to divide it in five tokens, one for each word, using the heuristic that each token is separated by a whitespace. For more complicated text, one could intuitively think that a good strategy would be to split on all non-alphanumeric characters, but this sometimes raises problems. This strategy would tokenize *isn't* in *isn* and *t* which is intuitively wrong. More complicated strategies are needed. Relevant to this work, the tokenization strategies used are language specific. One should not use an English tokenizer to tokenize Portuguese text (Branco & Silva, 2003).

2.1.2.2 Stemming and Lemmatization

Sometimes it is necessary to normalize lexical variations of a word to a base form, e.g., normalize the words *car*, *cars*, *car's* and *cars'* to just *car*. This can be useful, for example, in lexicon-based NER applications. If the word *car* is in the lexicon, it can make sense to consider lexical variations of the word *car* to be mentions of *car*. This can be accomplished by normalizing the words in the lexicon and in the text. This is done by using one of two techniques, Stemming or Lemmatization (Manning *et al.*, 2009a). In Stemming, crude rules are applied to cut off the suffixes of a word, the most popular stemmer being Porter's algorithm (Porter, 1980). On the other hand, Lemmatization does something similar but considers the context of the word.

2.1.3 Application of Text Mining on Radiology Reports

Text Mining tools can be used for automatic detection of important findings in Radiology Reports. For example, (Dreyer *et al.*, 2005) used an algorithm based on information theory to classify reports as having or not having important clinical findings and as having or not having recommendations for subsequent action. (Cotik *et al.*, 2015) did something similar for Spanish reports, using a translation of RadLex terms. These tools can also be used to detect the presence of more specific findings, as the presence of invasive mold diseases (Ananda-Rajah *et al.*, 2014) or invasive fungal diseases (Martinez *et al.*, 2015), both using a classifier based on a Support Vector Machine. Also possible is to extract general

2. RELATED WORK

information about reports (Hassanpour & Langlotz, 2016) and the data obtained can be used as input to other tools.

In literature it's possible to find some examples of Radiology reports/images search applications, that use NLP tools. The goals of these search tools include search for educational, research and clinical decision support purposes. One example of such a system is Render (Dang *et al.*, 2009), which even applies one of the information extraction system mentioned above (Dreyer *et al.*, 2005) to improve relevance of information retrieved.

Other applications include studying the appropriateness of existing Radiology reports templates, as done by (Hong & Kahn, 2013).

2.1.4 Ontologies

To answer the questions presented in Chapter 1, the RadLex ontology is used. An ontology is a "common, controlled knowledge representation designed to help knowledge sharing and computer reasoning" (Robinson & Bauer, 2011). It is a way to represent a subset of the real word which can be used as basis for communication between parties wanting to change information about that subset of the real word.

RadLex, for example, is a representation of the subset of the world related to Radiology which can be used as standard on how to talk about Radiology. Ontologies usually have a tree structure in which a class, representing some abstract entity in the real world, can have subclasses. For example, in RadLex, there is the class *clinical finding* which has subclasses *benign finding* and *pathophysiologic finding* (among others). This subclasses have a *is a* relationship with their parent classes: *benign finding* is a *clinical finding*. Other common relationship used in ontologies is the *part of* relationship.

Other popular examples of ontologies include the Gene Ontology¹, focused on genomics, SNOMED CT², a healthcare related ontology and ChEBI³, an ontology of small molecular entities.

¹<http://www.geneontology.org/>

²<http://www.snomed.org/snomed-ct>

³<https://www.ebi.ac.uk/chebi/>

2.2 Translation

2.2.1 Terminology

During this dissertation some terminology related to translation practices is used. In this section I briefly explain this terms (Koehn & Philipp, 2010).

Parallel Corpora - A *corpus* is just a set of texts (*corpora* is used if you want to refer to more than one of these sets). The term *parallel corpus* is used to refer to a set of texts paired with corresponding translations into other languages.

Language Pair - This term refers to the languages involved in a translation. For example, in a translation from Portuguese to English, we can say that the language pair is Portuguese-English, Portuguese being the *source language* and English the *target language*.

2.2.2 Machine Translation

Machine Translation (MT) is the use of computers to automatically translate natural language text. Currently, Statistical Machine Translation (SMT) is the most popular approach to MT. Other approaches included Rule-Based Machine Translation (RBMT) and Neural Machine Translation (NMT). RBMT involves the use of hand-crafted rules on how to do the automatic translation and NMT uses neural-networks and its use has recently been growing (Bentivogli *et al.*, 2016). I will now briefly review word-based and phrase-based which are both covered by the SMT approach. This is mostly based on (Koehn & Philipp, 2010).

2.2.2.1 Word-Based Models

These kind of models are not the state of the art anymore, but many of the principles and techniques of this approach are still in use today. The idea here is to translate the sentences word by word. Here is an example, translating English to Portuguese:

English - The bone is broken
Portuguese - O osso está partido

2. RELATED WORK

This is easy for a human to translate, but how would a computer know that *partido* is the translation of *broken* when *broken* has other potential translations? For example, the word *broken* could be interpreted as being financially ruined, as in "I've spent all the money in the casino, I'm completely *broken*". In that case, *broken* would be translated to *falido*. Of course, this doesn't make sense but the computer does not know that.

One way to teach the computer which translation to use would be to pick a large collection of English texts paired with the corresponding Portuguese translation and check how many times *broken* is translated to *partido* and how many times it is translated to *falido*. Lets assume that in our collection of texts the word *broken* is translated to *partido* 80% of the times and to *falido* 20% of the time. With this we could create a lexical translation probability table for the word *broken*. We could have a table like this for every word in the source text.

broken	
t	p(t s)
<i>partido</i>	0.8
<i>falido</i>	0.2

Table 2.1: Lexical translation probability table for the word broken

Here t stands for target, s stands for source and $p(t/s)$ is the probability that the target word is the translation of the source word. So, when the computer is translating the sentence above and arrives to the word *broken*, it checks the table and chooses *partido* as the translation because it has the higher probability of being the real translation. This type of estimation is called maximum likelihood estimation. What we are doing here is estimating what is called lexical translation probability distributions.

The example above was easy because the sentences were aligned word by word. This is not always the case. For example, the English expression *red swelling* should be translated to *inchaço vermelho*, not *vermelho inchaço*¹. Meaning, sometimes we must do some word reordering so that the translation is correct. This is accommodated by using an alignment model. But how can we generate an alignment model from a pair of collection of texts if we don't know which word

¹*red* can be translated to *vermelho* and *swelling* to *inchaço*

is aligned with which word? This is done by using the expectation maximization algorithm, which, in this case, iteratively applies the alignment model to the texts (expectation step) and learns the alignment model from the texts (maximization step) until convergence of the parameters in the algorithm.

With the lexical translation probability distributions and an alignment model we have a translation model. But this is not enough. A translation could be syntactically and semantically right but still not sounding right. For example, two possible translations of *chá forte* are *strong tea* and *powerful tea*. However, the second option doesn't sound right, it is not fluent. This problem is solved by using a language model. With an English language model, for example, we could calculate the probability that a certain sentence is correct English, considering all the data that was used to train the model. A language model would probably give a low probability to the phrase *powerful tea* because normally the word *powerful* is not used with the word *tea*.

We combine the language model and the translation model this way:

$$\arg \max_t \Pr(t|s) = \arg \max_t \Pr(s|t) \Pr(t) \quad (2.1)$$

We want to find the target word (t) with the higher probability of being the translation of the source word (s). $\Pr(t|s)$ represents the translation model and the $\Pr(t)$ represents the language model. This way of combining the translation and the language models is called noisy-channel model.

2.2.2.2 Phrase-Based Models

In this approach, instead of translating a sentence word by word we translate small words sequences at a time, sequences that we call phrases. These models have a better performance than the word-based models and this is not too surprising. Sometimes words are not the best unit of translation: there are cases when two words in the source sentence are translated into one word in the target sentence, for example. Another advantage is that translating phrases instead of words can help to solve ambiguities, as in the problem of deciding how to translate the text *chá forte* (see last section). We would check a parallel collection of texts and realize that most of the times *chá forte* is translated to *strong tea*. So, the

2. RELATED WORK

idea here is to divide the sentence in phrases, translate the phrases and do some reordering if necessary.

2.2.3 Post-editing

Post-editing (PE) is the task of editing, modifying and/or correcting a text that was pre-translated by use of MT, in order to improve the translation. (Somers, 2003) refers to the lower cost of MT+PE compared with HT to explain the growth of PE: companies want to become global but can't afford the cost of HT to translate from native language to the many languages they want to operate on.

(Koponen, 2016) tried to understand if MT+PE is really worth, compared with just HT, concluding that yes, most of the times it is worth it. But this depends on the quality of the MT, which in turn depends on, for example, the quality of the MT system and on the language pair.

Most of the research regarding PE refers to work done by professional translators. One approach that has been gaining traction is the use of the crowd to do the PE (Tatsumi & Aikawa, 2012). The advantages of this strategy include lower per-word cost and sometimes an higher speed, compared with HT. One big disadvantage is less assurance of quality.

2.2.4 Machine Translation Services

2.2.4.1 Yandex

Yandex¹ is a Russian search-engine company. Currently, Yandex.Translate (the name given to Yandex's MT system) uses a statistical approach. From their website², the system is composed by three components, a translation model, a language model and a decoder which is the part that actually does the translation.

I couldn't find any research paper evaluating the translation's quality of Yandex.Translate in the language pair Portuguese-English.

¹<https://yandex.com/>

²<https://tech.yandex.com/translate/doc/intro/concepts/how-works-machine-translation-docpage/>

2.2.4.2 Google

Google¹ is a company from the United States that offers a lot of technological services, including machine translation. For the language pair Portuguese-English, their translation services now use Neural Machine Translation² (see section 2.3.1).

2.2.4.3 Unbabel

Unbabel³ is a Portuguese start-up which sells translation services focused on conversational content like customer service or website copywriting, using an MT+PE approach. Although it is not mentioned in the Unbabel's API documentation, for the language pair Portuguese-English, currently Unbabel uses Google Translate's services in MT step of the MT+PE approach (personal communication). Next is an overview of Unbabel's translation pipeline:

1. Text is translated by MT (in this case, using Google Translate)
2. MT translated text is post-edited by users of the Unbabel platform. Users translate the text using Unbabel's web-interface or mobile app.
3. Translation resulting from last step is reviewed by an Unbabel's senior user, an user that was promoted for having good ratings

From now on I'm going to call this type of translation *Unbabel Translation*.

2.2.5 Translation of Medical Text

2.2.5.1 Multilingual Text Mining

There is not much research studying the effect of translation on Text Mining tools. (Castilla, 2007) is the most similar work to the one developed on this thesis and curiously, also studies translation of Portuguese medical text. In the main part of the study, Portuguese-written Radiology reports were translated to English using the SYSTRAN MT system, which uses a rule-based approach complemented with

¹<https://www.google.com>

²<https://blog.google/products/translate/found%2Dtranslation%2Dmore%2Daccurate%2Dfluent%2Dsentences%2Dgoogle%2Dtranslate>

³<https://unbabel.com/>

2. RELATED WORK

a specialized medical translation dictionary. Then the translation was processed by the Medical Language Extraction and Encoding System (MEDLEE) to extract information on the presence of mentions of certain medical conditions. The results were compared to reference results created by three radiologists on the original reports. The results are really positive, with values of sensitivity, specificity, positive and negative predictive values all above 88%. These results suggest that for this specific task of information extraction a MT translation retains a lot of information from the original text.

2.2.5.2 Machine Translation of Doctor-Patient Communication

Most of the work I found on medical translation focuses on translation of doctor-patient communication. This has the objective of breaking language barriers that sometimes exist between a doctor and a patient who do not speak the same language, with health-related consequences to the patient (Schyve, 2007). This could be done with trained medical interpreters but that option is costly compared with using MT and raises problems regarding patient confidentiality.

Several MT speech-to-speech translation systems for doctor-patient communication exist, but for most of them, evaluations are not found in the literature. One exception is (Bouillon *et al.*, 2005) which studies MedSLT, a multilingual spoken language translation system tailored for headache, chest pain and abdominal pain domains. However, (Bouillon *et al.*, 2005) only studies the appropriateness of the design choices within the system, not comparing its performance with any other system. Others example of systems of this type are Jibbig¹, Universal Doctor² and Transonics (Nagata & Pedersen, 2005).

(Kaliyadan & Pillai, 2010) did a small study on the use of Google Translate to translate between English and French during doctor-patient interaction in India medical offices, with promising results regarding patient satisfaction. Also using Google Translate, (Patil & Davies, 2014) studied the quality of the translation of 10 commonly used medical statements to 26 languages. Of all the 260 translations, 57.7% were right. The results were better for Western European languages than for others. Portuguese had the highest score, with 9 of the 10 sentence

¹<http://jibbiggo-translator-2-0.soft112.com/>

²<http://www.universaldocor.com/>

translated being right. Other work was also done on non-European languages, which have less resources (Kathol *et al.*, 2005; Musleh *et al.*, 2016).

Some researchers (G *et al.*, 2013; Kaliyadan & Pillai, 2010; Marta R. Costajussà, Mireia Farrús, 2012) suggest that MT should be used very cautiously in this situations, because of imperfect performance in a domain where accuracy is really important. One way to improve the systems could involve the use of existing public medical terms database (Eck *et al.*, 2004).

2.2.5.3 Machine Translation of Public-Health Information

In the USA, most of the public health information is written in English, although a substantial percentage of the population have limited English proficiency. One of the barriers for more widespread translation is the cost of translation services and a way of streamlining the process would be using MT+PE. (Kirchhoff *et al.*, 2011; Turner *et al.*, 2015) studied the feasibility of this system for translation from English to Spanish, with some promising results, and to Chinese, which was more problematic.

2.2.5.4 Machine Translation for Information Retrieval

The ACL 2014 Ninth Workshop on Statistical Machine Translation had a Medical Translation Task (Bojar *et al.*, 2014), which consisted in two subtasks: translation of sentences from summaries of medical articles and translation of queries entered by users of medical information search engines. This task was supported by the Khresmoi¹ project which develops a multilingual search and access system for biomedical information and documents, allowing the user to make search queries and read summaries of the results in their own language. The task had 8 participants, the winner being the UEDIN team (Durrani *et al.*, 2014) which used the Moses phrase-based system².

¹<http://khresmoi.eu/>

²<http://www.statmt.org/moses/>

2. RELATED WORK

2.2.5.5 Machine Translation of Other Types of Medical Text

Studies of the translation of other types of documents are also present in the literature. For example, (Wolk & Marasek, 2015) compares neural based with statistical machine translation of descriptions of medical products in the language pair Polish-English, obtaining mixed results.

More related to the work done on this thesis, (Castilla, 2007) studied the use of the MT application SYSTRAN to translate sentences from Radiology reports. The MT system uses a ruled-based approach and was complemented with a specialized medical translation dictionary. The translations were evaluated by an expert in the field, finding good scores for understandability, fidelity with original text and translation coverage of the original text.

(Zeng-Treitler *et al.*, 2010) tested if a general-purpose machine translation tool like the Babel Fish¹ is adequate to translate sentences of discharge summaries, surgical notes, admission notes and Radiology reports from English to Spanish, Chinese, Russian and Korean. They found that most of the times the translation is incomprehensible and inaccurate.

More recently, there was a Biomedical Translation Task during the ACL 2016 First Conference on Machine Translation (WMT16) in which the participants were asked to submit systems to translate titles and abstracts from scientific publications (Bojar *et al.*, 2016). The evaluators note that the quality of the machine translation is still poor in comparison to the reference translations. The only submissions to the English-Portuguese and Portuguese-English translation tasks (Aires *et al.*, 2016) were the ones with the worse results relative to the baseline system.

2.2.6 Translation of Biomedical Lexicons

One alternative solution to the one I am exploring in this thesis, translating the medical text to English, is to translate the lexicon, on which the task at hand depends on, to the language of the medical documents we want to study. For example, if a researcher have a Spanish corpus and wants to annotate it with terms of some lexicon, it will be a problematic task since most of the available

¹<https://www.babelfish.com/>

ontologies are not multilingual. To solve this the researcher could translate the ontologies she wants to use to the language of the corpus. This example is similar to (Cotik *et al.*, 2015), in which all RadLex terms were translated to Spanish using Google Translate and medical reports were annotated with this translated terms.

For the German language there is (Bretschneider *et al.*, 2014). Having in mind that translating all the entries of an ontology one wants to use would be expensive, the authors propose translating only a subset of the ontology, a subset relevant to the task at hand. They do this semi-automatically with the help of the corpus one wants to annotate. With this, the authors improved the annotation of German text with RadLex terms.

2.3 External Tools and Terminologies

Some of the work done during the thesis used and was inspired by some external tools and terminologies that I now briefly review.

2.3.1 RadLex

RadLex¹ is an ontology which focuses on Radiology-related terms. It was developed to standardize annotation, indexation, and retrieval of Radiology information resources in the digital world (Langlotz, 2006) and it helped to fill a gap in Radiology terminology (Langlotz & Caldwell, 2002; Woods & Eng, 2013). The RadLex terms were originally gathered from existing ontologies at the time, including the American College of Radiology (ACR) Index, SNOMED-CT, and the Foundational Model Anatomy and it is a highly dynamic ontology: its number of terms grew from around 8000 to around 75000 in just ten years. Being an ontology, RadLex can be visualized as a tree, which contains other subtrees. This characteristic can be used to extract subsets of the RadLex ontology. For example, if someone just wants to use the RadLex classes related to clinical findings she could just use the RadLex subtree containing just the children of the RadLex class "clinical finding".

¹<http://www.rsna.org/RadLex.aspx>

2. RELATED WORK

There are a few studies on the completeness of RadLex. (Marwede *et al.*, 2008) found that an old version of RadLex covered 84% of terms extracted manually from 250 thoracic CT reports, with higher coverage for terms in the *Findings* (90%) category and lower coverage for the *Modifier* category (78%). Curiously, in a study using more recent versions of RadLex (versions 3.1–3.5) (Woods & Eng, 2013) found a lower coverage of 62% using the same type of reports (they used less reports in this study, just 100). They find higher coverage for the categories of *anatomic objects* and *physiological conditions* and lower coverage for the categories of *imaging observations* and *procedures* (the categories used in both studies are not the same). The authors justify the lower coverage with the inclusion in the study of categories such as *procedures*, which did not had any match with RadLex terms. They also used a different methodology to find matches between manual extracted terms and Radlex terms. These studies analyzed the coverage of RadLex of terms mentioned in the contents of radiology reports. (Hong *et al.*, 2012), on the other hand, studied how well RadLex covers the terms of templates of structured Radiology reports developed by the Radiological Society of North America, finding that 41% of the terms found in the templates matched exactly to RadLex and that 26% matched partially. Since these analysis, new versions of RadLex were launched so the results and critics present in the studies are not necessarily relevant anymore.

One could use RadLex to assist in the matching of research articles manuscripts to reviewers profiles, like done by the RadioGraphics journal (Klein, 2013). Or to help in the visual analysis of neurography images (Wang *et al.*, 2015). Having said this, most of the examples described in the literature are of applications related to Information Retrieval (IR), the task of extracting some information resource from a collection of information resources. These resources can be images or websites, for example. One such example of a IR system using Radlex, is (Spanier *et al.*, 2016), who takes advantage of the tree structure of this ontology to create a new method of case-based image retrieval (M-CBIR). Most existing M-CBIR systems use low-level characteristics of medical images (like color, shape and texture) to induce similarity between them. But this is problematic since medical images which show the same type of content can have different low-level characteristics. One solution is to induce this similarity from the information

contained in the textual radiological reports that accompany the images and the authors take advantage of RadLex to do just that. This can help radiologists to find related medical cases in a certain database which then can help them in their decision-making process. Other approaches to IR systems using Radlex include the ones described in (Do *et al.*, 2010), (Kurtz *et al.*, 2014) and (Gerstmair *et al.*, 2012).

2.3.2 Open Biomedical Annotator

The Open Biomedical Annotator (OBA)¹ is an open-source tool for NER using a lexicon-based approach, made available by the North-American National Center for Biomedical Ontology (Jonquet *et al.*, 2009), which can be used to annotate text with concepts from ontologies. For example, if you go to the website, input a Radiology report and choose the ontology RadLex, the tool will return all the mentions in the text of terms belonging to the RadLex terminology. OBA uses MGrep, which implements a radix-tree based data structure that allows for a fast match between terms in a lexicon and terms in text. OBA can easily be used as a web-service and it is relatively fast. It uses a case-insensitive direct match approach, not considering lexical variations of words (see 2.1.1).

2.3.3 NOBLE Coder

NOBLE Coder² (Tseytlin *et al.*, 2016) is a software for NER using a lexicon-based approach. The lexicon is set by the user (it has to be in UMLS (RRF)³, OWL⁴ or OBO⁵ formats or be present in BioPortal⁶). The lexicon is processed into two hash-tables which are then used during by NOBLE to find, in an arbitrary text, mentions of terms found in the lexicon.

Unlike the system used by OBA, NOBLE can find mentions of lexical variations of the terms present in the Lexicon because it applies word Stemming. For

¹<http://biportal.bioontology.org/annotator>

²<http://noble-tools.dbmi.pitt.edu/>

³<https://www.ncbi.nlm.nih.gov/books/NBK9685/>

⁴<https://www.w3.org/OWL/>

⁵<http://www.geneontology.org/faq/what-obo-file-format>

⁶<http://biportal.bioontology.org/ontologies>

2. RELATED WORK

example, *lobe* is a term present in the RadLex terminology, but its plural, *lobes*, is not. However, NOBLE considers that *lobes* is a mention of the term *lobe*, which is right. But this can sometimes go wrong; for example, NOBLE considers that *headings* is a mention of the RadLex term *head*, which is wrong. So although this strategy can improve recall it does so at the cost of precision.

The NOBLE tool is flexible in what is considered a mention of a lexicon term, giving the user the power to adapt the tool for her specific purposes. This can be done by choosing to use or not a certain *matching option*. These include:

- **Subsumption** - Only match the longest mention. For example, *toe*, *toe skin* and *skin* are all RadLex terms. If the "Subsumption" option is set, in the text *toe skin*, only the term *toe skin* will be recognized. Otherwise, the terms *toe* and *skin* are also recognized.
- **Overlap** - If this option is used, matched terms can overlap each other. For example, if this option is not set, NOBLE will only recognize the terms *deep* and *lateral margin* in the text *deep lateral margin*. If it is set, it will also recognize the term *deep margin* which overlaps with the two other terms.
- **Contiguity** - Terms must be contiguous to be matched. For example, if set, in the text *multiple ducts lesions* both *multiple ducts* and *multiple lesions* are considered matches, although *multiple* and *lesions* are not adjacent to each other. It's possible to set how many irrelevant words can be between words belonging to a term (in Table 2.2, this is called *gap*).
- **Order** - Terms must be in the same order as in the Lexicon to be considered mentions. If not set, *lesions multiple* is considered a mention of the Radlex term *multiple lesions*.
- **Partial** - Partial match with terms in Lexicon are considered a Lexicon term mention. If set, *multiple* is considered a mention of *multiple lesions*.

The user can also choose to, for example:

- Skip single letter words

2.3 External Tools and Terminologies

- Skip stop words
- Use heuristics to filter out potential false positives
- When a term can be considered a mention of more than one concept in the Lexicon, select only the highest scoring one

Different combinations of these options are useful for different purposes. NOBLE already offers some built-in matching strategies, listed in Table 2.2.

Combination of matching options					
Task	Subsumption	Overlap	Contiguity	Order	Partial
<i>Best match</i>	Yes	Yes	Yes (gap=1)	No	No
<i>All match</i>	No	Yes	No	No	No
<i>Precise match</i>	Yes	Yes	Yes (gap=0)	Yes	No
<i>Sloppy/Partial match</i>	No	Yes	No	No	Yes

Table 2.2: NOBLE matching strategies present in the GUI interface. Adapted from (Tseytlin *et al.*, 2016). This correspond to the options used in the GUI tool.

The authors of the tool provide suggestion for what kind of task each strategy is more appropriate. For example, they suggest that the *best match* strategy is best for concept coding and information extraction and that the *all match* strategy is more suitable for information retrieval and text mining.

(Tseytlin *et al.*, 2016) compares the NOBLE tool with other Lexicon-based NER tools, finding that its performance in recognizing terms from lexicons its comparable with other similar software like Concept Mapper(Stewart *et al.*, 2012) or cTAKES¹², although it probably depends a lot on the corpus used.

One big advantage of NOBLE is its ease of use compared with other similar systems. Little or no programming skills are needed to use the software since it includes a GUI (Graphical User Interface) which allows an user to upload lexicons in a number of formats and easily annotate texts.

¹<https://cwiki.apache.org/confluence/display/CTAKES/cTAKES+3.0+-+Dictionary+Lookup>

²<https://cwiki.apache.org/confluence/display/CTAKES/cTAKES+3.2+-+Fast+Dictionary+Lookup>

2. RELATED WORK

2.4 Evaluation Metrics

For a certain task (for example, annotation of a corpus with terms related to diseases) it is useful to have standard evaluation metrics so that we can compare many systems and know which one is the best. In information retrieval and information extraction systems precision (P), recall (R) and F-score (F) are the measures that are mostly used. For example they were the measures used in a competition which involved a task similar to the example I gave above (Elhadad *et al.*, 2015).

To use these measures we need to have a reference, a gold-standard, which we assume represents the perfect performance in a certain task, the ground truth. In the example of extraction of disorder mentions, it could be an annotation done by a human expert. To calculate these measures we also need the number of true positives, false positives and false negatives. I will illustrate each one of these with the example of the annotation of disease mentions.

- True positive (TP) – The system being tested annotated a term also annotated in the reference.
- False positive (FP) – The system annotated a term that is not annotated in the reference.
- False negative (FN) – The system did not annotate a term that is annotated in the reference.

Precision corresponds to the fraction of the terms annotated by the system that are also annotated in the reference.

$$P = \frac{TP}{TP + FP} \quad (2.2)$$

If of the 10 terms annotated by the system, only 6 are annotated by the reference, then the system has a precision of 0.6. If every term extracted by the system is also extracted by the reference, then the system has a precision of 1, the best score possible. But the system can have a score of 1 if only annotates one right term, even though there are a lot of other terms annotated in the reference.

This system, although having a score of 1, would not be very useful. Recall is a measure that helps to solve this issue.

Recall calculates what fraction of all terms annotated in the reference are annotated by the system.

$$R = \frac{TP}{TP + FN} \quad (2.3)$$

If the system annotates 8 terms of the 10 that are annotated in the reference, then it has a recall of 0.8. If it annotates all of them, it has a recall of 1, the perfect score. But, as is the case with precision, this measure also has problems. If the system annotates all the terms in a corpus, it will have a perfect score in the recall measure, because it is sure to have annotated all the terms annotated in the reference, although it also annotated a lot of wrong terms.

As you can see, both measures have problems when used in isolation. One way to combine them is by using the F-score measure, that corresponds to the harmonic mean of precision and recall.

$$F - score = 2 * \frac{P * R}{P + R} \quad (2.4)$$

2.4.1 Micro- and Macro- Evaluation Metrics

Now imagine that you want evaluate your system on more than one document. How do you aggregate the metrics explained above? You can sum the TP, FP and FN values of each document and then use the Precision, Recall and F-Score formulas exposed above. With this approach, you would calculate the Micro Precision, Micro Recall and Micro F-score.

Another approach is to calculate Precision, Recall and F-Score for each document and then average for all documents. This would give you the Macro Precision, Macro Recall and Macro F-score values.

Chapter 3

Framework

3.1 MRRAD (Multilingual Radiology Research Articles Dataset) Corpus

From the best of my knowledge there is no parallel corpus of Radiology reports. So I created a Portuguese-English parallel corpus of research articles related to Radiology, assuming that the writing style and content of these research articles are similar to Radiology reports. For each research article the MRRAD corpus contains:

1. Original Portuguese text
2. Human Translated English text
3. Machine Translated English text (Yandex)
4. Machine Translated English text (Google)
5. Machine Translation + Post-Editing English text (Google + Unbabel)

In the next few lines I will explain how I've constructed the corpus.

3. FRAMEWORK

3.1.1 Web Crawl of the articles (1,2)

To obtain a list of research articles related to Radiology, that were available both in English and in Portuguese, I used the NCBO Entrez Programming Utilities (E-utilities)¹ to query the PubMed database with the search query “portuguese[Language] AND english[Language] AND radiography[MeSH Major Topic] AND hasabstract[text]” (search done on 11/12/2016). The last filter is used to avoid getting texts for which only the title is available.

Then I programmatically crawled each article PubMed page to get the URL where the full article could be found. Most of the articles were hosted in SciELO² and only articles hosted in there were included in the corpus. More, only articles for which the original language is Portuguese are included in the corpus.

Finally, I programmatically crawled the SciELO pages for each article to get both the original Portuguese texts and corresponding the English translations. From the HTML of each page I extracted everything from the abstract until, but not including, the references/bibliography.

Three of the articles were about surveys, not containing much vocabulary about Radiology (PMIDS: 19936506, 22002140, 23515770). They were excluded from the corpus. Other two contained encoding problems and were also excluded (PMIDS: 21793046 and 24263777).

The final result is a parallel corpus of 51 articles, distributed by journal as shown in Table 3.1.

To give a sense of the corpus size, the human English translations have a total of 163,423 words³ the longer article having 12,451 and the smaller 848. The articles have an average of 3,204 words each.

3.1.2 Note On Human Translations

It is not known for sure how exactly the original human translations were performed, since some of the articles are not recent and some of the journals did not answer our request for more information about the translation, but all the

¹<https://www.ncbi.nlm.nih.gov/books/NBK25501/>

²<http://www.scielo.br/>

³Tokenization done by NLTK’s word_tokenize function (<http://www.nltk.org/>)

3.1 MRRAD (Multilingual Radiology Research Articles Dataset) Corpus

Table 3.1: Number of articles by journal in parallel corpus

Journal	Number Of Articles
Arquivos Brasileiros de Cardiologia	24
Jornal Brasileiro de Pneumologia	14
Revista do Colégio Brasileiro de Cirurgiões	4
Brazilian Journal of Otorhinolaryngology	2
Arquivos Brasileiros de Cirurgia Digestiva	2
Revista Brasileira de Cirurgia Cardiovascular	2
Jornal da Sociedade Brasileira de Fonoaudiologia	1
Einstein (São Paulo)	1
Revista Brasileira de Reumatologia	1

answers received mentioned the use of specialized translation services. Having said this, it is being assumed that the translations are of high quality since they are published by scientific magazines.

3.1.3 Yandex Translation (3)

I used Yandex’s free Translate API¹ to machine translate the Portuguese version of the articles. Yandex is a Russian company which, among other things, sells automatic translation services, but it has a limited free service. It currently uses a Statistical approach to Machine Translation. Each translation request had a limit of 10,000 characters so we developed software to break the text to various pieces, without breaking the text in the middle of sentences, send the translation request for each piece and then join everything back.

3.1.4 Google and Unbabel Translation (4,5)

Both MT with Google and MT+PE with Unbabel were obtained using Unbabel’s API². I obtained Google’s Statistical Machine translation using the *mt_translation* endpoint of the API and Unbabel’s Machine Translation + Post-Editing using the *translation* API’s endpoint. The requests for Unbabel Translations have a limit of words, so I used a software similar to the one utilized for the Yandex Translations.

¹<https://tech.yandex.com/translate/>

²<http://developers.unbabel.com/>

3. FRAMEWORK

3.2 Annotation

All the English versions of the articles in the corpus were annotated three with RadLex terms, one time using a direct matching approach and two using two of the built-in matching strategies provided by NOBLE Coder. I'm calling the three approaches Direct Match¹, All Match and Best Match². Three different kinds of approaches were used to check what effect the annotation strategy have on the results.

3.2.1 Direct Match - Annotation with Open Biomedical Annotator

The articles were annotated with OBA using the REST API³. The default parameters were used, namely the ones shown in Table 3.2.

Table 3.2: OBA parameters used

Parameter	Value
expand_class_hierarchy	false
expand_mappings	false
minimum_match_length	3
exclude_numbers	false
whole_word_only	true
exclude_synonyms	false
longest_only	false

3.2.2 All Match and Best Match - Annotation with NOBLE Coder

NOBLE Coder was chosen against others similar tools because of its comparable quality and higher ease of use. Each of the articles was annotated twice with this tool, using two different matching strategies, Best match and All match.

¹See 2.1.1

²See 2.3.3

³http://data.bioontology.org/documentation#nav_annotator

More information on how NOBLE Coder was used can be found at the appendices.

3.3 Evaluation

For each document and annotation approach I created the set of the RadLex terms (identified by their RIDs) that were found in that document with that annotation approach. This is the data used in the assessment of translation solutions that follows.

The RadLex terms extracted from each MT or MT+PE translated article were compared against the ones extracted from the corresponding HT translated article, which was considered the gold standard. Both Micro- and Macro- Precision, Recall and F1-scores were calculated. This was done for each matching approach.

To facilitate the understanding of the results, we will now walk through a short example. Consider that we have one Portuguese document and corresponding HT English translation and MT English translation. Four terms of interest were identified in the HT translation, bone, cell, finger, colon¹. This is going to be our gold standard. In the MT translation, 2 terms of interest were found, brain, bone. One of these terms is also in the gold standard, which means $TP = 1$, but the other term is not, $FP = 1$. In the gold standard there are 3 terms that were not found in the MT translation, which means $FN = 3$. After calculations (see 2.4), this gives us a Precision score of 0.5, a Recall score of 0.25 and F-Score of 0.33. These methods measure how similar are the terms annotated on the MT or MT+PE texts to the terms annotated on the HT texts.

¹We use here human understandable names instead of RIDs so that the example is easier to follow

Chapter 4

Experimental Results

4.1 NER Lexicon-based approach

Translation	Direct Match	All Match	Best Match
Human	119.55	177.92	145.0
Yandex	116.06	173.92	145.16
Google	120.8	179.49	147.61
Unbabel	120.92	178.86	148.16

Table 4.1: Number of RadLex terms found by document

Table 4.1 presents the number of RadLex extracted by document using the different annotation approaches. One of the highlights here is that the All Match approach consistently found more terms than the Best Match approach, which itself found more terms than the Direct Match approach. This was expected since the All Match approach is the most flexible one in what it considers to be a mention of a RadLex term. The Best Match approach is more strict than the All Match approach but less than the Direct Match approach, considering lexical variations and word reordering, for example. But in all cases we can see that many terms are being extracted from each document.

As seen in Figures 4.1, 4.2 and 4.3 (Supplementary Table 1 presents the data in table format), the terms extracted from Google translations are more similar to the ones extracted from HT translations than the ones from Yandex transla-

4. EXPERIMENTAL RESULTS



Figure 4.1: Micro Evaluation of Translations being tested (All Match)



Figure 4.2: Micro Evaluation of Translations being tested (Best Match)



Figure 4.3: Micro Evaluation of Translations being tested (Direct Match)

tions. This could be just because the human translators used Google Translator to help them in their translation process. This argument loses strength if we assume Google Translate translation outputs changed since the articles were human translated (publication years of the articles used range from 2003 to 2013), but data could not be found to corroborate this assumption.

The terms extracted from Unbabel and Google translations are really similar, the F-Scores being almost equal. That the translations are similar is not too surprising since the Post-Editing phase at Unbabel is done after MT translation using Google. What could be surprising is that Unbabel does not have a higher score. One conclusion to take from this is that Post-Editing step on the MT+PE is not useful for this task. The results are similar when a Macro Evaluation is done (see Supplementary Table 2).

In the Introduction to the thesis I've proposed the following hypothesis:

Hypothesis: MT+PE is a good trade-off between quality and cost, compared with MT and HT, for translating Portuguese Radiology reports to English, for the purpose of identifying RadLex terms in the translated text.

4. EXPERIMENTAL RESULTS

I've written that for this to be true, "The terms extracted from MT+PE translations have to be more similar to the ones extracted from the HT translation than the ones extracted from MT translations". This does not hold. So, for this task, if someone had to choose between Google and Unbabel, this someone would be better off using Google since it is cheaper.

I've also written that for the hypothesis to be true, "The terms extracted from MT+PE translations have to be similar enough to the ones extracted from the HT translation". The terms extracted from any of these translations are not extremely different but they are also not equivalent to the ones extracted from the human translation. It could be the case that for some applications only translations close to human quality are acceptable, while for other applications a mediocre translation would be good enough. Therefore, the suitability of the MT and MT+PE translations probably depends on the practical usage for these translations and annotations.

To better understand the results I will now provide a detailed analysis on the annotations for the "clinical finding" and "anatomical entity" subtrees of RadLex. These are two of the subtrees that probably would be more important when applying RadLex to a Information Retrieval system, a type of application for which the results of this study can be useful.

4.1.1 Clinical Finding and Anatomical Entity Subtrees

Depending on the type of annotation approach and translation it was found between 35.25 and 55.55 "clinical finding" or "anatomical entity" terms per document (See Supplementary Table 3). As seen in Figures 4.4, 4.5 and 4.6 (see Supplementary Table 4 to see the data in table format), the scores obtained are similar to the ones obtained for all terms, with Yandex translation extracted terms being the less similar to the HT translation extracted terms and Google and Unbabel having similar scores. Similar results were found when Macro evaluation was performed (see Supplementary Table 5). But why these scores?

In an attempt to better understand the results, I did an analysis of the False Positives and False Negatives errors committed by the MT and MT+PE translations, focusing on the terms belonging to the "clinical finding" or "anatomical

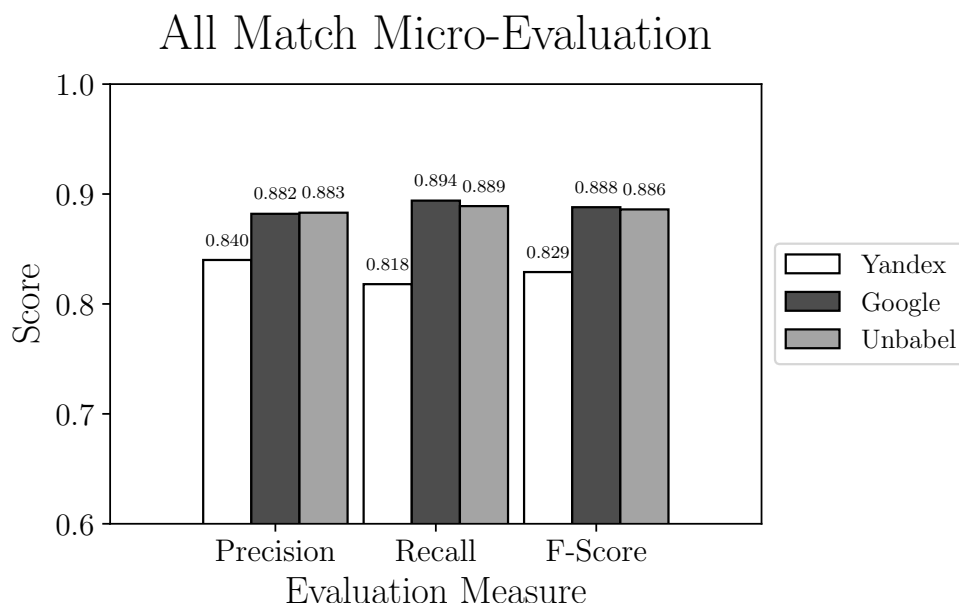


Figure 4.4: Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (All Match)

entity" RadLex subtrees. From preliminary analysis I knew that some of the FPs and FNs are not caused by a bad translation but for other causes, for example, a different translation which is correct but causes a different annotation, e.g., translating *parênquima pulmonar* to *pulmonary parenchyma* instead of to *lung parenchyma*. Both translations are correct but the second one leads to the extraction of the term lung while the first does not. Still, I expected a higher number of real errors in the Yandex translations compared with the Unbabel or Google translations, since both of these types of translation had better scores.

I did an analysis on the FPs and FNs errors committed by Yandex, Google and Unbabel translations in 9 random documents and each error was classified by type (See Supplementary Table 6). The results from the Best Match Approach were used. As predicted, the percentage of errors of Yandex due to a wrong translation (25% of 100 FPs or FNs) was higher than the percentage of errors of Google and Unbabel (22.09% of 86 and 21.18% of 85 FPs or FNs, correspondingly), but only slightly (See Supplementary Tables 7 and 8). The reasons for the others FPs

4. EXPERIMENTAL RESULTS

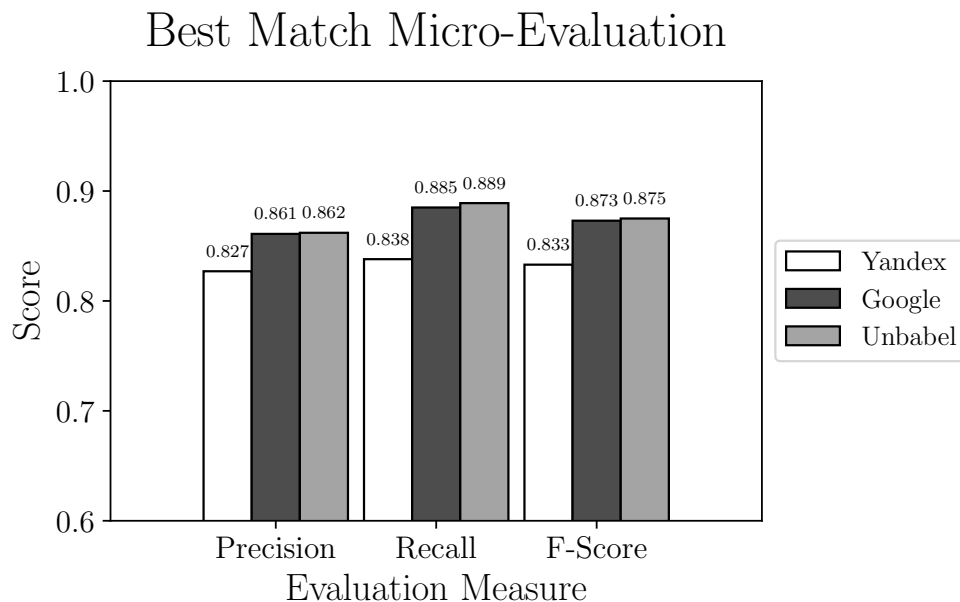


Figure 4.5: Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (Best Match)

and FNs included, among others, cases i) of different translations which are both correct but lead to different annotations, as described above and ii) in which the word extracted does not have the same meaning in the text as it has in RadLex. For example, the case of extracting the anatomical term *hand* from "(...) on the other hand, it has to be considered that (...)" , in which the word *hand* is used metaphorically. This happens because a rule-based approach is being used, which does not consider the context of the term.

There were a lot of these ii) cases, which maybe would happen less if a Machine Learning approach was used. I thought about this but the problem is that, to the best of our knowledge, there is no training and test data large enough, i.e., English Radiology text resulting from human translation of Portuguese text and annotated with RadLex terms by experts.

Next I analyzed what kind of translation errors were done when the FPs and FNs were caused by real translation errors (See Supplementary Table 9). These subcategories included cases in which:

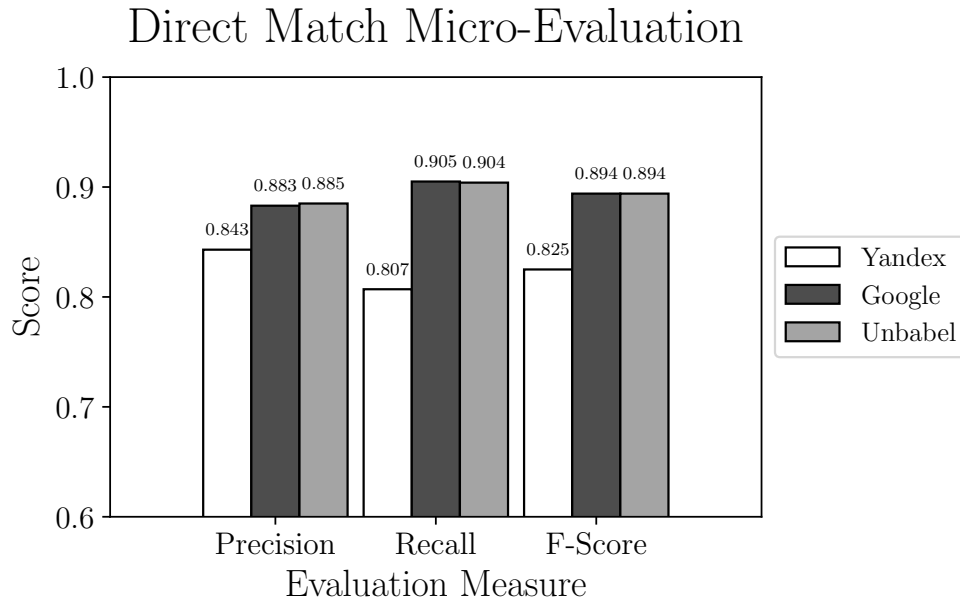


Figure 4.6: Micro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (Direct Match)

- There was an extra word in the translation;
- There was a missing word in the translation;
- A wrong hyphenation was used;
- An acronym was not translated;
- The test translation used a term that was too general;
- A wrong lexical variation was used;
- The most correct medical term was not used;

Each of these cases had a low number of occurrences and so it is not worth a deeper analysis. One interesting thing to note is that in the Yandex translations there were some cases (six) in which the original Portuguese word was not even translated. This never happened in the Google and Unbabel translations that were analyzed. This could be explained by the fact that probably Yandex focuses

4. EXPERIMENTAL RESULTS

on different languages than Google and so their Portuguese-English translation and/or language models are not so well trained. But most of the errors correspond to just to a general wrong choice of terms to use as a translation. For example, translating *média* to *middle* instead of *mean* or *lesões de via biliar* to *lesions via bile* instead of *lesions to the biliary tract*. This type of problems could probably be solved by training the models used by Google and Yandex with more data, specifically data related to medicine.

One could expect that Unbabel translations would have a lot less mistakes than Google's but this is not always the case. There are situations where errors are even added during the Post-Editing step. A review of the errors makes us propose that this could be due to the lack of medical knowledge of Unbabel current editors. For example, a stroke is something that occurs in the brain but in one case it was used as something that happens in the heart - someone with some knowledge on medicine would not make this error. But the truth is that Unbabel currently do not have a focus on medical content. I predict that if they did and invested in growing a crowd of users with a better knowledge of medical language, this would lead to better results.

Chapter 5

Conclusions

References

- AIRES, J., LOPES, G.P. & GOMES, L. (2016). English-Portuguese Biomedical Translation Task Using a Genuine Phrase-Based Statistical Machine Translation Approach. *Proceedings of the First Conference on Machine Translation, Volume 2: Shared Task Papers*, **2**, 456–462. [18](#)
- ANANDA-RAJAH, M.R., MARTINEZ, D., SLAVIN, M.A., CAVEDON, L., DOOLEY, M., CHENG, A. & THURSKY, K.A. (2014). Facilitating surveillance of pulmonary invasive mold diseases in patients with haematological malignancies by screening computed tomography reports using natural language processing. *PloS one*, **9**, e107797. [9](#)
- BENTIVOGLI, L., BISAZZA, A., CETTOLO, M. & FEDERICO, M. (2016). Neural versus Phrase-Based Machine Translation Quality: a Case Study. In *To appear: EMNLP-2016*. [11](#)
- BOJAR, O., BUCK, C., FEDERMANN, C., HADDOW, B., KOEHN, P., LEVELING, J., MONZ, C., PECINA, P., POST, M., HERVE, S.A., SORICUT, R., SPECIA, L. & TAMCHYNA, A. (2014). Findings of the 2014 Workshop on Statistical Machine Translation. *2014 Workshop on Statistical Machine Translation*, 12–58. [17](#)
- BOJAR, O., CHATTERJEE, R., FEDERMANN, C., GRAHAM, Y., HADDOW, B., HUCK, M., JIMENO YEPES, A., KOEHN, P., LOGACHEVA, V., MONZ, C., NEGRI, M., NEVEOL, A., NEVES, M., POPEL, M., POST, M., RUBINO, R., SCARTON, C., SPECIA, L., TURCHI, M., VERSPOOR, K. & ZAMPIERI, M. (2016). Findings of the 2016 Conference on Machine Translation. In *Proceedings of the First Conference on Machine Translation*, vol. 2, 131–198. [18](#)

REFERENCES

- BOUILLON, P., RAYNER, M., CHATZICHRISAFIS, N., HOCKEY, B.A., SANTA-HOLMA, M.E., STARLANDER, M., NAKAO, Y., KANZAKI, K. & ISAHARA, H. (2005). A Generic Multi-Lingual Open Source Platform for Limited-Domain Medical Speech Translation. In *Proceedings of the 10th Conference of the European Association for Machine Translation (EAMT)*, 50–58. 16
- BRANCO, A. & SILVA, J. (2003). Tokenization of Portuguese: resolving the hard cases. Tech. rep., University Of Lisbon. 9
- BRETSCHNEIDER, C., OBERKAMPF, H., ZILLNER, S., BAUER, B. & HAMMON, M. (2014). Corpus-based Translation of Ontologies for Improved Multilingual Semantic Annotation. *Coling*, 1–8. 2, 19
- CASTILLA, C. (2007). *Instrumento de Investigação Clínico-Epidemiológica em Cardiologia Fundamentado no Processamento de Linguagem Natural*. Ph.D. thesis, Universidade de São Paulo, São Paulo. 2, 15, 18
- A survey of current work in biomedical text mining. 7
- COTIK, V., FILIPPO, D. & CASTAÑO, J. (2015). An Approach for Automatic Classification of Radiology Reports in Spanish. *Studies in Health Technology and Informatics*, **216**, 634–638. 9, 19
- DANG, P.A., KALRA, M.K., SCHULTZ, T.J., GRAHAM, S.A. & DREYER, K.J. (2009). Informatics in radiology: Render: an online searchable radiology study repository. *Radiographics : a review publication of the Radiological Society of North America, Inc*, **29**, 1233–46. 10
- DO, B.H., WU, A., BISWAL, S., KAMAYA, A. & RUBIN, D.L. (2010). Informatics in Radiology: RADTF: A Semantic Search-enabled, Natural Language Processor-generated Radiology Teaching File. *RadioGraphics*, **30**, 2039–2048. 21
- DREYER, K.J., KALRA, M.K., MAHER, M.M., HURIER, A.M., ASFAW, B.A., SCHULTZ, T., HALPERN, E.F. & THRALL, J.H. (2005). Application of Recently Developed Computer Algorithm for Automatic Classification of Unstructured Radiology Reports: Validation Study. *Radiology*, **234**, 323–329. 9, 10

REFERENCES

- DURRANI, N., HADDOW, B., KOEHN, P. & HEAFIELD, K. (2014). Edinburgh’s Phrase-based Machine Translation Systems for WMT-14. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, 97–104. 17
- ECK, M., VOGEL, S. & WAIBEL, A. (2004). Improving statistical machine translation in the medical domain using the unified medical language system. *Proceedings of the 20th international conference on Computational Linguistics - COLING ’04*, 792–es. 17
- ELHADAD, N., PRADHAN, S., GORMAN, S.L., MANANDHAR, S., CHAPMAN, W.W. & SAVOVA, G. (2015). SemEval-2015 Task 14 : Analysis of Clinical Text. *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, 303–310. 24
- G, R., M, F., R, A., O, E. & K, P. (2013). Using machine translation in clinical practice. *Canadian family physician Medecin de famille canadien*, **59**, 382–383. 17
- GERSTMAYER, A., DAUMKE, P., SIMON, K., LANGER, M. & KOTTER, E. (2012). Intelligent image retrieval based on radiology reports. *European Radiology*, **22**, 2750–2758. 3, 21
- HASSANPOUR, S. & LANGLITZ, C.P. (2016). Information extraction from multi-institutional radiology reports. *Artificial Intelligence in Medicine*, **66**, 29–39. 1, 10
- HONG, Y. & KAHN, C.E. (2013). Content analysis of reporting templates and free-text radiology reports. *Journal of Digital Imaging*, **26**, 843–849. 10
- HONG, Y., ZHANG, J., HEILBRUN, M.E. & KAHN, C.E. (2012). Analysis of RadLex coverage and term co-occurrence in radiology reporting templates. *Journal of Digital Imaging*, **25**, 56–62. 20
- HOTH, A., NÜRNBERGER, A. & PAASS, G. (2005). A Brief Survey of Text Mining. *LDV Forum - GLDV Journal for Computational Linguistics and Language Technology*, **20**, 19–62. 7

REFERENCES

- JONQUET, C., SHAH, N.H. & MUSEN, M.A. (2009). The open biomedical annotator. *Summit on translational bioinformatics*, **2009**, 56–60. 8, 21
- KALIYADAN, F. & PILLAI, S.G. (2010). The use of Google language tools as an interpretation aid in cross-cultural doctor-patient interaction: A pilot study. *Informatics in Primary Care*, **18**, 141–143. 16, 17
- KATHOL, A., PRECODA, K., VERGYRI, D., WANG, W., RIEHEMANN, S., INTERNATIONAL, S.R.I. & PARK, M. (2005). Speech Translation for Low-Resource Languages : The Case of Pashto. *Syntax*, 2273–2276. 17
- KIRCHHOFF, K., TURNER, A.M., AXELROD, A. & SAAVEDRA, F. (2011). Application of statistical machine translation to public health information: a feasibility study. *Journal of the American Medical Informatics Association : JAMIA*, **18**, 473–478. 17
- KLEIN, J.S. (2013). A Look Back at 2012 and Plans for 2013. *RadioGraphics*, **33**, 1–2. 20
- KOEHN & PHILIPP (2010). *Statistical Machine Translation*. Cambridge University Press. 11
- KOPONEN, M. (2016). Is machine translation post-editing worth the effort? A survey of research into post-editing and effort. *Journal of Specialised Translation*, 131–148. 14
- KURTZ, C., DEPEURSINGE, A., NAPEL, S., BEAULIEU, C.F. & RUBIN, D.L. (2014). On combining image-based and ontological semantic dissimilarities for medical image retrieval applications. *Medical Image Analysis*, **18**, 1082–1100. 21
- LANGLOTZ, C.P. (2006). RadLex: a new method for indexing online educational materials. *Radiographics : a review publication of the Radiological Society of North America, Inc*, **26**, 1595–1597. 19

REFERENCES

- LANGLOTZ, C.P. & CALDWELL, S.A. (2002). The completeness of existing lexicons for representing radiology report information. *Journal of digital imaging : the official journal of the Society for Computer Applications in Radiology*, **15** Suppl 1, 201–205. [19](#)
- MANNING, C.D., RAGHAVAN, P. & SCHÜTZE, H. (2009a). Stemming and lemmatization. In *Introduction to Information Retrieval*, chap. 2.2.4. [9](#)
- MANNING, C.D., RAGHAVAN, P. & SCHÜTZE, H. (2009b). Tokenization. In *Introduction to Information Retrieval*, chap. 2.2.1. [9](#)
- MARTA R. COSTA-JUSSÀ, MIREIA FARRÚS, J.S.P. (2012). Machine Translation in Medicine. In *ARSA - PROCEEDINGS IN ARSA - ADVANCED RESEARCH IN SCIENTIFIC AREAS*, 1995–1998, EDIS - Publishing Institution of the University of Zilina. [17](#)
- MARTINEZ, D., ANANDA-RAJAH, M.R., SUOMINEN, H., SLAVIN, M.A., THURSKY, K.A. & CAVEDON, L. (2015). Automatic detection of patients with invasive fungal disease from free-text computed tomography (CT) scans. *Journal of Biomedical Informatics*, **53**, 251–260. [9](#)
- MARWEDE, D., SCHULZ, T. & KAHN, T. (2008). Indexing thoracic CT reports using a preliminary version of a standardized radiological lexicon (RadLex). *Journal of Digital Imaging*, **21**, 363–370. [20](#)
- MOHIT, B. (2014). Named Entity Recognition. In I. Zitouni, ed., *Natural Language Processing of Semitic Languages*, chap. 7, 221–246, Springer. [7](#)
- MUSLEH, A., DURRANI, N., TEMNIKOVA, I., NAKOV, P., VOGEL, S. & AL-SAAD, O. (2016). Enabling Medical Translation for Low-Resource Languages. *Proceedings of the 16th Conference on Intelligent Text Processing and Computational Linguistics (CICLING)*. [17](#)
- NAGATA, M. & PEDERSEN, T. (2005). Proceedings of the ACL Interactive Poster and Demonstration Sessions. In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*. [16](#)

REFERENCES

- PATIL, S. & DAVIES, P. (2014). Use of Google Translate in medical communication: evaluation of accuracy. *BMJ (Clinical research ed.)*, **349**, g7392. 16
- PONS, E., BRAUN, L.M.M., HUNINK, M.G.M. & KORS, J.A. (2016). Natural Language Processing in Radiology: A Systematic Review. *Radiology*, **279**, 329–343. 1, 7
- PORTER, M. (1980). An algorithm for suffix stripping. *Program: electronic library and information systems*, **14**, 130–137. 9
- ROBINSON, P.N. & BAUER, S. (2011). *Introduction to bio-ontologies*. CRC Press. 10
- SCHYVE, P.M. (2007). Language differences as a barrier to quality and safety in health care: The joint commission perspective. *Journal of General Internal Medicine*, **22**, 360–361. 16
- SOMERS, H. (2003). *Computers and translation: a translator's guide*. John Benjamins Publishing Company. 14
- SPANIER, A.B., COHEN, D. & JOSKOWICZ, L. (2016). A new method for the automatic retrieval of medical cases based on the RadLex ontology. *International Journal of Computer Assisted Radiology and Surgery*. 20
- STEWART, S.A., VON MALTZAHN, M.E. & ABIDI, S.S.R. (2012). Comparing metamap to mgrep as a tool for mapping free text to formal medical lexicons. In *CEUR Workshop Proceedings*, vol. 895, 63–77. 8, 23
- TATSUMI, M. & AIKAWA, T. (2012). How Good Is Crowd Post-Editing? Its Potential and Limitations. . . . *2012 Workshop on* 14
- TSEYTLIN, E., MITCHELL, K., LEGOWSKI, E., CORRIGAN, J., CHAVAN, G. & JACOBSON, R.S. (2016). NOBLE - Flexible concept recognition for large-scale biomedical natural language processing. *BMC bioinformatics*, **17**, 32. xxi, 8, 21, 23

REFERENCES

- TURNER, A.M., DEW, K.N., DESAI, L., MARTIN, N. & KIRCHHOFF, K. (2015). Machine Translation of Public Health Materials From English to Chinese: A Feasibility Study. *JMIR public health and surveillance*, **1**, e17. [17](#)
- WANG, K.B., SALUNKHE, A., MORRISO, J., LEE, P., MEJINO, J., DETWILER, L.F., BRINKLEY, J.F., SIEGEL, E.D., RUBIN, D. & CARRINO, J. (2015). Ontology-based image navigation: Exploring 3.0-T MR neurography of the brachial plexus using AIM and radlex. *Radiographics*, **35**, 142–151. [20](#)
- WOLK, K. & MARASEK, K. (2015). Neural-based Machine Translation for Medical Text Domain. Based on European Medicines Agency Leaflet Texts. In *Procedia Computer Science*, vol. 64, 2–9. [18](#)
- WOODS, R.W. & ENG, J. (2013). Evaluating the completeness of radlex in the chest radiography domain. *Academic Radiology*, **20**, 1329–1333. [19](#), [20](#)
- ZENG-TREITLER, Q., KIM, H., ROSEMBLAT, G. & KESELMAN, A. (2010). Can multilingual machine translation help make medical record content more comprehensible to patients? In *Studies in Health Technology and Informatics*, vol. 160, 73–77. [18](#)

Commands used to run NOBLE Coder

```
$ java -jar NobleCoder-1.0.jar -terminology radlex \  
-input [portuguese reports path] -output [output path] \  
-search all-match\textit{
```

```
$ java -jar NobleCoder-1.0.jar -terminology radlex \  
-input [portuguese reports path] -output [output path] \  
-search best-match
```


Uploading RadLex ontology to NOBLE

I had to edit the RadLex ontology .owl file before it could be correctly processed and uploaded to NOBLE Coder. In the original .owl file the properties "Preferred_name" and "Synonym" are considered to be *DatatypeProperty* but I had to change both to *AnnotationProperty*. That is, where in the file was

```
<owl:DatatypeProperty rdf:ID="Preferred_name">  
</owl:DatatypeProperty>
```

I've had to change it to:

```
<owl:AnnotationProperty rdf:ID="Preferred_name">  
</owl:AnnotationProperty>
```

And the analogous thing for the "Synonym" property.

Macro Evaluations



Figure 1: Macro Evaluation of Translations being tested (All Match)

. MACRO EVALUATIONS

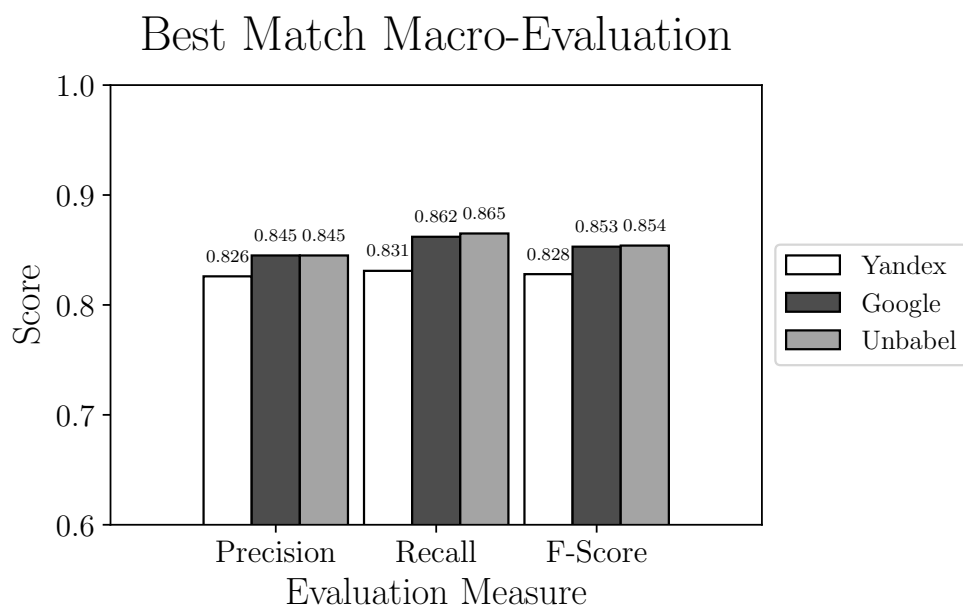


Figure 2: Macro Evaluation of Translations being tested (Best Match)

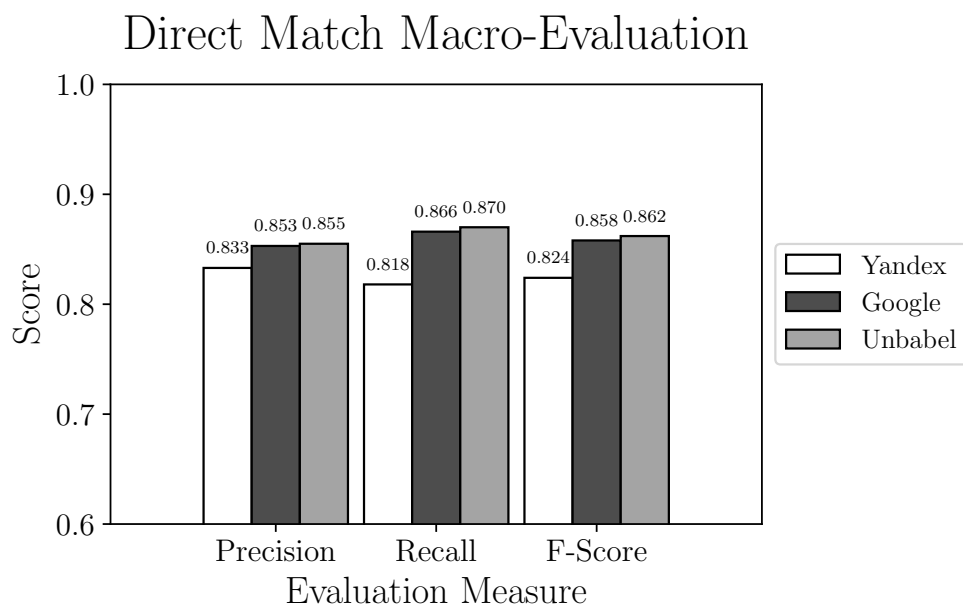


Figure 3: Macro Evaluation of Translations being tested (Direct Match)

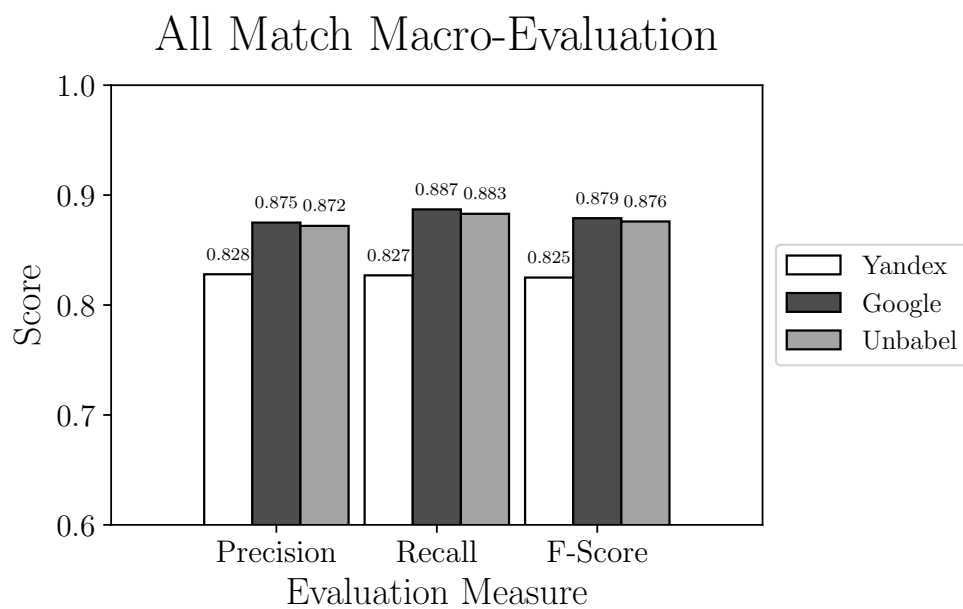


Figure 4: Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (All Match)

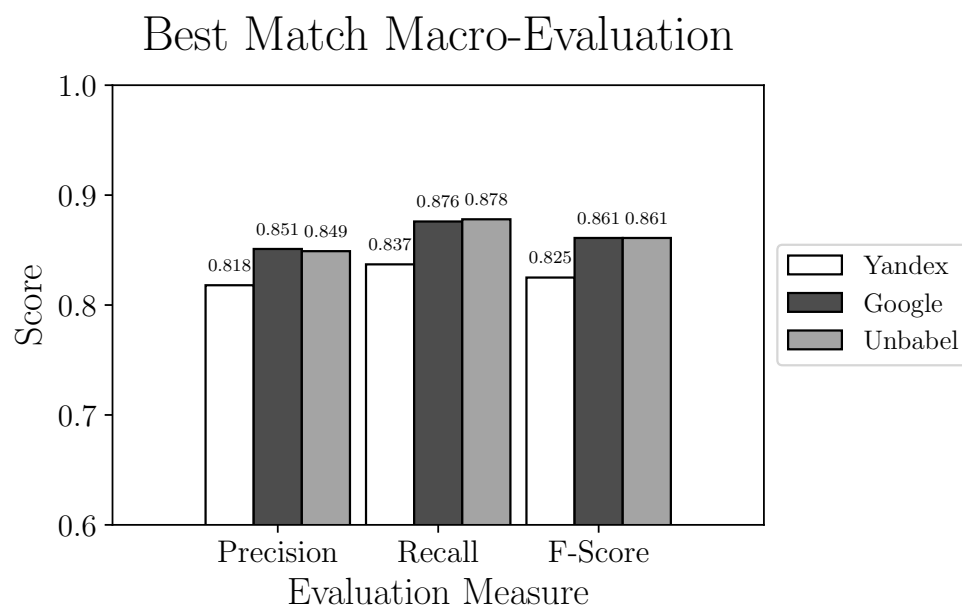


Figure 5: Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (Best Match)

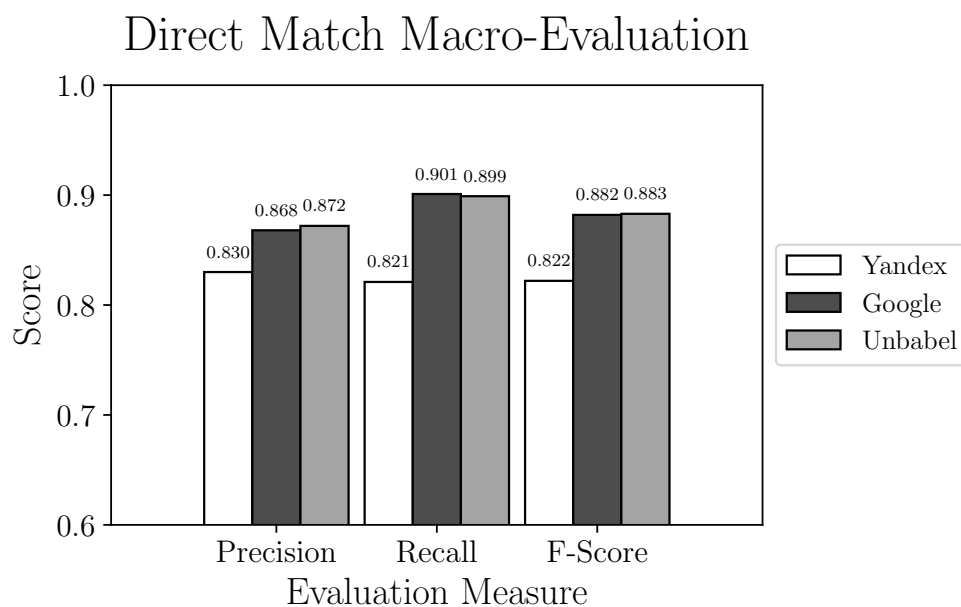


Figure 6: Macro Evaluation of translations being tested, considering just the terms from RadLex "clinical finding" and "anatomical entity subtrees" (Direct Match)