

Universidad de Valladolid

FACULTAD de FILOSOFÍA Y LETRAS DEPARTAMENTO de FILOLOGÍA INGLESA Grado en Estudios Ingleses

TRABAJO DE FIN DE GRADO

The Usefulness of Specialized Monolingual Corpora in Translation activity

Sandra Ivonne Silva Aedo

Tutor: Beatriz Méndez Cendón

Curso 2013/14

ABSTRACT

The aim of the project is to point out the great usefulness that specialized monolingual corpora can provide translators in their challenging task of medical translation. To carry out our purpose, a specialized corpus is built. We take into account the criteria of corpus compilation and corpus analysis. Once it is built, we compare the usage of our corpus with lexicographical resources, such as: specialized dictionaries and glossaries, in order to indicate the benefits that translators can obtain from corpora. Finally, our conclusion is reachedbe on the basis of the importance of using corpora in translation activity.

Keywords: medical translation, corpora, phraseology, KWIC concordance, collocations.

RESUMEN: El propósito del proyecto es señalar la gran utilidad que el corpus especializado monolingüe puede ofrecer a los traductores en la difícil tarea de la traducción médica. Para ello, compilamos un corpus especializado y, a continuación, analizamos los criterios que hemos tenido en su compilación y análisis. Más adelante, llevamos a cabo una comparación de su uso con recursos lexicográficos frecuentes, como los diccionarios o glosarios especializados, para indicar los beneficios que los traducotres pueden obtener de los corpus. Por último, ofrecemos una conclusión basada en la importancia de usar un corpus en la actividad de la traducción.

Palabras clave: traducción medica, corpus, fraseología, concordancia KWIC, colocaciones.

Index

1.	Introduction
2.	English Specialized language12
3.	Corpora in translation15
	3.1. Definition of corpora15
	3.2. Importance of corpora15
	3.3. Types of corpora17
4.	Corpus compilation and analysis
	4.1. Corpus design and compilation
	4.1.1. Selection of suitable material
	4.1.2. Selection of computer software
	4.2. Corpus analysis
	4.2.1. Term extraction
	4.2.2. Collocations
5.	Specialized monolingual corpora in translation activity
6.	Conclusion71
7.	References

List of tables

Table 1: Collocates of cell	
Table 2: Collocates of tumor	
Table 3: Collocates of cancer	
Table 4: Collocates of gene	
Table 5: Collocates of mutation	
Table 6: Collocates of <i>imaging</i>	
Table 7: Collocates of chemotherapy	
Table 8: Collocates of DNA	
Table 9: Collocates of apoptosis	
Table 10: Collocates of carcinoma	
Table 11: Collocates of antibody	45
Table 12: Collocates of metastasis	45
Table 13: Collocates of stage	
Table 14: Collocates of PET	
Table 15: Collocates of toxicity	46
Table 16: Collocates of MRI	
Table 17: Collocates of surgery	46
Table 18: Collocates of recurrence	46
Table 19: Collocates of biopsy	
Table 20: Collocates of lymphoma	47
Table 21: Collocates of screenings	
Table 22: Collocates of phenotype	47
Table 23: Collocates of prognosis	
Table 24: Collocates of genome	
Table 25: Collocates of angiogenesis	
Table 26: Collocates of adenocarcinoma	
Table 27: Collocates of enzyme	
Table 28: Collocates of antigen	

Table 29: Collocates of oncogene 49
Table 30: Collocates of staging
Table 31: Collocates of <i>leukemia</i> 49
Table 32: Collocates of <i>differentiation</i>
Table 33: Collocates of <i>biomarker</i> 50
Table 34: Collocates of translocation
Table 35: Collocates of mortality
Table 36: Collocates of immunohistochemistry 50
Table 37: Collocates of <i>nodule</i>
Table 38: Collocates of mastectomy
Table 39: Collocates of chromosome
Table 40: Collocates of <i>infusion</i>
Table 41: Collocates of necrosis 51
Table 42: Collocates of carcinogenesis
Table 43: Collocates of brachytherapy 52
Table 44: Collocates of OS
Table 45: Collocates of SD 52
Table 46: Collocates of carcinogen
Table 47: Collocates of <i>prostatectomy</i>
Table 48: Collocates of gene 53
Table 49: Collocates of <i>immunotherapy</i>
Table 50: Collocates of dissection
Table 51: Collocates of nucleus
Table 52: Collocates of macrophage
Table 53: Collocates of <i>lumpectomy</i> 54
Table 54: Collocates of chemoprevention 54
Table 55: Collocates of <i>colonoscopy</i>
Table 56: Collocates of X-ray
Table 57: Collocates of resection 55

Table 58: Collocates of osteosarcoma	. 55
Table 59: Collocates of breast cancer	. 55
Table 60: Collocates of lymph node	. 55
Table 61: Collocates of radiation therapy	. 56
Table 62: Collocates of stem cell	. 56
Table 63: Collocates of primary tumor	. 56
Table 64: Collocates of clinical trials	. 57
Table 65: Collocates of cell proliferation	. 57
Table 66: Collocates of <i>bone marrow</i>	. 57
Table 67: Collocates of cell cycle	. 57
Table 68: Collocates of tumor volume	. 57
Table 69: Collocates of drug resistance	. 58
Table 70: Collocates of stem cell transplantation	. 58
Table 71: Collocates of cancer stem cell	. 58
Table 72: Collocates of DNA repair	. 58
Table 73: Collocates of squamous cell carcinoma	. 59
Table 74: Collocates of DNA methylation	. 59
Table 75: Collocates of partial response	. 59
Table 76: Collocates of risk factor	. 59
Table 77: Collocates of colon cancer	. 60
Table 78: Collocates of progression-free survival	. 60
Table 79: Collocates of <i>local recurrence</i>	. 60
Table 80: Collocates of survival rate	. 60
Table 81: Collocates of colorectal cancer	. 60
Table 82: Collocates of multiple myeloma	. 61
Table 83: Collocates of non-small cell lung cancer	. 61
Table 84: Collocates of side effects	. 61
Table 85: Collocates of in situ	. 61
Table 86: Collocates of CT scan	. 61

Table 87: Collocates of tumor suppressor gene	. 62
Table 88: Collocates of carcinoma in situ	. 62

1. Introduction

The production of a text into another language, most specifically in a specialized text, covers a set of aspects regarding terminology and translation. Therefore, it does not only involve the fact of transferring words from one language into another, but the meaning of the source text must be expressed according to the rules of the target language in order to be accepted by the audience. In this regard, the translators play the role of intermediary between two languages and cultures.

Every translator follows a translation process to carry out this assignment. Following the classification of Sven Tarp (2007), the translation process can be classified into three phases: "preparation, translation and revision".

The preparation phase deals with the collection of data in order to acquire knowledge of the topic involved.

The second phase (the translation) refers to the understanding of the source language (SL). This phase translators may find difficulties since they find new terms. In order to solve their problems, they look for the equivalences of the terms in dictionaries or glossaries. The problem lies in the fact that these sources provide them with several options and they may select the wrong ones or the less appropriate ones. As a consequence their translation may not be accepted by the LSP target readers. It is due to the fact that they lack of information and therefore they do not satisfy the needs of translators.

The last phase is the production of the target text (TT), in which translators need, apart from the equivalence of the terms, grammatical and syntactic information about the terms. It is also difficult that translators find this information in the lexicographical resources mentioned above. In most cases, these resources are not helpful because they are not outdated. To summarize, a good translation into L2 does not only involve using the appropriate terminology but also use the right grammatical structures in order to make the target text sounds naturally. In other words, translators have to achieve a translation that must be as close as possible to original texts produced within a particular field of LSP.

As it is well known, the Internet is a useful source that has changed translation activity together with the increasing availability of corpora. It is true that a corpus cannot respond to all answers of translators, but provides them with information that the previous resources do not offer. A corpus, therefore, provides several advantages for their users. For example, since it contains thousands of words, translators can identify and extract collocational, conceptual and terminological information using some software programs. The fact that the texts are gathered means that they can be accessed easily by translators.

For the reasons mentioned previously, the aim of our project is to compile a corpus that consists of original texts in a particular subject field in order to point out that it can become a considerable help for the translators. We also demonstrate how software programs can be used by the translators in order to solve their translation problems, and in this way, they can produce a high-quality translation.

Our project is divided into four parts: the first part deals with English specialized language and its main characteristics. We believe that the understanding of specialized language is essential for translation. In the second section, we introduce a definition of corpora and explain the importance of corpora in translation studies. We also describe different types of corpora. In the third chapter, we explain the characteristic of our specialized corpus and its compilation process, which is followed by the corpus analysis that consists of the retrieval of terms and collocations. Finally, the last part deals with, an analysis in which examine the role that a specialized monolingual corpus plays in translation activity by using a Spanish article. In this regard, our aim is to point out its usefulness, but also to highlight the differences that exist between using this tool and common resources (e.g. dictionaries and glossaries). The project ends, with a conclusion in which we summarize the main points of our work.

2. English Specialized language

The translation of a specialized text is an activity that is totally different from a text belonging to general language. One of the reasons is that specialized language has its own terminology and its own phraseology and it makes that specialized language different from the general language, although they also share similar characteristics.

Before focusing on the main topic of our project, we would like to indicate what English specialized language is and to demonstrate the differences and similitudes between general and specialized languages since we consider that this information should be learnt by translators before they begin to translate a specialized text.

First, English specialized language is part of LSP (Language for Specific Purposes), which may be defined as a language used by a specific group in order to communicate and discuss within a particular field by using a specialized knowledge. This term comprises English specialized languages, French specialized languages, German specialized languages, Spanish specialized languages, etc. Therefore, it is said that LSP covers all global specialized languages (Méndez Cendón 2012-13).

Second, although LSP differs from LGP (Language for general purposes) there is an interrelation between them. It is due to the fact that they share the same grammatical structures. In other words, there aren't two different grammars and, therefore, we find that grammatical structures of LSP exist in LGP. They also share many words thanks to **de-terminologization**. It is a process by which specialized words start being used in general language, and it may be for several reasons such as the emergence of scientific knowledge in the society (Meyer and Mackintosh, 2000). For example, "virus" is a term from the medical field that is nowadays used in the computer field and by the people in general to refer to a program that affects their computers.

LSP has its own characteristics and these are (Méndez Cendón 2012-13):

- Different degrees of abstraction within LSP: topic, communicative situation and users.
- Monofunctional character because it used by a restricted number of users who learn specialized language voluntarily.
- Knowledge of specialized language is mandatory for members who use it in an expert community.
- Each subject field has its own specific terminology.

As mentioned above LSP and LGP share grammatical characteristics, but in LSP we find a high frequency of the following structures (Méndez Cendón 2012-13):

- Specific terminology
- Prefixes and suffixes
- Premodifiers (complex noun phrases)
- Compound nouns
- Passive constructions
- Reduce relative clauses
- Nominalizations

Third, LSP is classified according to specialized fields, communicative context and users or participants. Pearson (2002:27) classifies LSP according to users: expert, semi–experts and non – experts. As there are three types of users, also there are different levels of LSP communication and these are the following:

Expert- to expert communication: The writer and the reader share a high knowledge of specialized language, and therefore they understand the terminology and complex structures that are used in texts.

Expert- to semi-experts communication: The writer and the reader do not share the same level of knowledge. However, the reader may be familiar with some of the terms used within the field. In this type of communication, the writer provides the users with some explanations in order to make them understand the terms.

Expert- to non-expert communication: The situation in which the reader is not familiar with terminology. For this reason, the author uses words belonging to general language in order to describe the concepts in an easier way.

Finally, it is important to say that if the user is interested in learning LSP he needs to acquire two types of knowledge: linguistic and conceptual as Bowker and Pearson (2002: 30) point out. Linguistic knowledge refers to terms, collocations, grammatical structures and stylistic features, whereas conceptual knowledge deals with the information about the specialized concepts.

3. Corpora in translation

In this section, we provide an overview about the concept of corpora, followed by their importance in translation and finally their types.

3.1. Definition of corpora

There are several definitions for a corpus and maybe one of the easier to understand is the one that Bowker and Pearson formulate in their book *Working with Specialized Language: A Practical Guide to Using Corpora* (2002:9). They state that a corpus can be described as "a large collection of authentic texts that have been gathered in electronic form according to a specific set of criteria".

Therefore, the texts that a corpus comprises have to be selected according to a specific purpose to be used as a representative sample of language, or a particular aspect of language.

3.2. Importance of corpora

The usage of corpora has been introduced in translation studies recently. Electronic corpora have changed translation training and, scholars have published several studies about this. Mona Baker (1993) was one of the first to discuss the importance of corpora in translation. She says that they can help translators to investigate translations. In other words, translators would analyze the texts and extract relevant information for their translations.

Baker (1996) also indicates that comparable corpora are very useful for translators because they allow them to extract valuable information when translating, due to the fact that this type of corpora contains original texts and not translations. On the other hand, Maeve Olohan points out:

Pearson (2003) and Bowker (1999) state that comparable corpora can be a useful resource for students and teachers involved in translator training, used for checking terminology and collocates,

identifying text-type-specific formulations, validating intuitions and providing explanations for appropriateness of certain solutions to problems. (2004: 172).

Therefore, this type of corpora provides translators with different types of information which help in their translation assignment.

Silvia Bernardini is another of the scholars who gives importance to the usage of corpora in translation. She indicates that the corpus helps translation students to obtain a better translation. She says "The result is a better – documented, more accurate, as well as more fluent translation" (Bernardini 2004:20).

Moreover, the development of corpus analysis tools such as Wordsmith and AntConc have also led to several translation scholars publish studies based on corpora in translation. Some of them are Belinda Maia (2005) and Natalie Kübler (2003).

Maia (2005) argues that the use of general corpora allows translation students to go beyond the dictionary because they can observe words in context. Additionally, Kübler (2003) focuses especially on technical translations and indicates: "Learning to use corpora and corpus analysis tools can give future translators the technical skills that were usually not associated with translation".

Therefore, there is no doubt that the usage of corpora has played an important role in the translation field and continues to do so. Thanks to them, the translators can investigate style and usage, as well as search for information about contexts in which a word occurs. As Bowker and Pearson (2002) already mentionated, a corpus can be used as a translation resource.

However there is no just one type of corpus, but there are different types of corpora and translators should determine the type of corpus they will compile.

3.3. Types of corpora

Corpora are usually classified according to their purpose. The types of corpora are "general corpora", "specialized corpora", "parallel corpora", "comparable corpora" and "monitor corpora".

As for "general and specialized corpora", according to Baker (1995), a general corpus is used for a wide range of purposes, giving all the relevant varieties of the language to the researchers. An example of this type is British National Corpus (BNC) for English Language and Real Academia Española (RAE) for Spanish language. In contrast, a specialized corpus is smaller and it is due to the fact that it focuses on a specific aspect of language.

Other types of corpora are "parallel" and "comparable". They are different because a parallel corpus is "composed of a set of the source texts and their translations" (Bravo Gozalo, 1998: 225), whereas a comparable corpus only contains original texts in two or more languages. In order to be a comparable corpus, both subcorpora have to share common features such as size, date of publication, genre, subject – field, etc. If not, they are just a collection of texts.

A "monitor corpus" is one that increases in size regularly in order to provide the users with the changes in language. In the words of Sinclair (1987:21), it is "a dynamic rather than a static phenomenon, consisting of very large amounts of electronically-held texts which will pass through the computer". An example of this type is COCA (Corpus of Contemporary American English).

Apart from the types of corpora listed above, there are other types in terms of mode, language and time.

Corpora can be classified according to mode, which deals with the way in which we want to recollect the information and there are two types: "written corpora" and "spoken corpora". Atkins et al. (1992) point out that whether mode is not specified, it will be written by default.

According to time, corpora may be "synchronic" and "diachronic". It is a synchronic corpus when the purpose of the study focuses on one particular period of time whereas it is diachronic when its aim is related with the development of time.

Other factor involved in the classification of corpora is language: monolingual, bilingual or multilingual. A "monolingual corpus" refers to the type of corpus that includes only one language whereas "bilingual corpus" deals with two languages. A "multilingual corpus" comprises at least three languages.

In translation studies, "comparable" and "parallel corpora" are often the most useful tool for translators. "Parallel corpora" are used as a bilingual dictionary because the translator can identify terms and phrases with their equivalences into another language. And although, comparable corpora contain only original texts, they can also be used by translators to look for equivalent terms and phraseology and also by those who are interested in the topic that they deal with.

It is not an easy assignment to design "parallel corpora" since it is difficult to find the translations of source texts to compile a good parallel corpus.

As for the corpus designed and compiled for this project, it is monolingual. The aim is showing that the usage of this corpus can help translators in their translation assignment.

4. Corpus compilation and analysis

4.1. Corpus design and compilation

Once the translators decide the type of corpora that they will create, they must consider different criteria for their corpus. These criteria are important because they may influence the results of their analysis. According to Bowker and Pearson (2002: 45-52), the criteria are size, medium, subject, text type, authorship, language and date of publication.

Size is one of factors that the translators must take into account. Sinclair (1991:18) argues that it must be "as large as possible, and should keep on growing" while Bowker and Pearson (2002:45) point out that there are no hard and fast rules that can be followed to determine the ideal size of a corpus. And Biber (1993) says that size depends on the purpose for which the corpus is intended. Therefore, the translators are the one who have to determine how big their corpus is going to be, and their decision will be based mainly on the aim of their study.

The size of the corpus of this project is 231, 989 words. It is not big if we compare it to a general corpus. Obviously, there are reasons that justify this difference. First, we found some difficulties to collect the samples since most the texts are not available on the Web. And second, this type of corpora is usually smaller than general corpora since it analyzes a particular aspect of language. Although the size of the corpus is small, it includes a reasonable number of texts in order to achieve the expected results. For example, if the aim of the translators is to find synonyms of specified terms and the size of their corpus is too small, this would be a problem in their analysis.

However, the size of a corpus is not the only relevant aspect but also the **size of the texts**. As well as the size of the corpus is based on the purpose of the study, the size of the texts also depends on this. If the translators only select chunk texts, they then delete part of the information which may be important for their study. For this reason, full texts are the best option rather than chunk texts. We would also like to mention the question of the **copyright**

in building a corpus. If the compiler is going to include whole texts he/she should ask permission to the editorial to use those texts for research or translation purposes. In this project we are just using the four typical sections of a medical research paper to compile the corpus: Introduction, Materials and Methods, Results and Discussion (IMRAD superstructure). This means that we are omitting the References.

The total **number of texts** is another criterion related to the size of the corpus. Whether the translators are interested in identifying the way in which the authors use specialized language, it would be advisable that he will include a great number of texts from a variety of authors in their corpus. This way, they will get information about what terms and concepts are typically used in a specific subject field. The texts included in our corpus of the project are written by different authors since the purpose is to identify terminology and grammatical structures commonly used by experts.

It is also important to determine whether the corpus will contain only **written material** or **spoken material**, or both. It also depends on the study that the translators will do. The medium chosen for the corpus that we built is written material and it is due to the main two reasons: first, the purpose is to identify terms and grammatical structures that experts use in scientific texts and second, it is much easier to collect this type of material than the other option.

It is said that all the texts included in a corpus have to deal with the specialized **subject** that the user will analyze. However, Bowker and Pearson (2002:50) point out that if your project sets out to study particular features of specialized language, it may not be necessary for all your texts to be about the same subject. Our purpose is to analyze the specialized subject of Oncology and thus all our texts deal with this subject.

It is also necessary to include texts in the corpus that belong to the same **text type**. The problem lies in the fact that the texts are written in the different style. For example, if we read a text from a newspaper we notice that the language style is totally different from a

text written in a magazine. As this study is based on a specific aspect of language, it was necessary to compile a corpus that included texts that share the same characteristics in their structure. The samples of language for the corpus are research articles, written by experts and addressed to experts.

Moreover, **language** is an issue that the translators take into account when they are compiling their corpus. The texts that they are going to include in their corpus must be original texts and not translations. The reason why original texts are the best option is that they contain expressions written by the experts in the source language. This information is essential for the translators, who aim to produce a translation that be read as an original text written in the target language.

And finally, **date of publication** of the texts also depends on the aim of the study. The translators are thus who determine whether they want to include up-to-date texts in their corpus or not. As stated above, the purpose of creating a corpus is to help the translators to be faced with problems that arise when they are translating, and we consider that up- to-date texts are essential since terms become obsolete in scientific fields.

4.1.1. Selection of suitable material

The translator is required to include appropriate information in his corpus to obtain a helpful analysis. No doubt that one of the easiest ways of searching for material is on the Internet since it contains an amount of information on almost any subject. However, this information given is not always suitable.

Bowker and Pearson (2002:61) state two main tools that the user can use for searching his texts and are: **search engines** and **subject directories**. The former are tools that allow them to enter the search term on the Web, for example, Google, Bing, and Yahoo search; whereas the latter are web-sites organized by people, who classify them according to subject field.

We consider that subject directories are better than search engines because the former contain a selective list according to the field that you are interested. For the selection of the texts of our corpus, we used the database of the University of Valladolid's library because thanks to the information retrieved we could make sure the quality of the texts.

This database has a list of several sections such as Science and Technology, Legal Sciences, ProQuest Database, Journal Citations Report, etc. The section chosen was Journal Citations Report, which gives the users two types of selection: by categories and by journal or category title. Due to the fact that we did not know the titles of journals, we chose the first option by selecting the field of oncology. After our selection, we were given a detailed list indicating the main journals frequently read in the field of oncology. We selected the journals paying attention to their impact factor and to their total cites. We also took into account those that came from USA because our concern was to design a corpus that only included American English texts since the main oncology journals are North American.

The journals were: *Neoplasia, Cancer Journal, Oncotarget, Lancet Oncology*, and *Radiation Oncology*. Once the journals were selected, the next step was to download their articles from the Internet. The number of articles collected was limited by copyright restrictions in the case of *Oncotarget, Cancer Journal, Lancet Oncology* and *Radiation Oncology*. However, *Neoplasia* was the only journal in which their articles were immediately downloaded thanks to its accessibility.

In the image below there is a table that illustrates the medical journals we used for our corpus and the number of articles collected from each journal:

Medical journals	Number of texts
Neoplasia	17
Oncotarget	9
Cancer Journal	7
Lancet Oncology	6
Radiation Oncology	4

Image 1: Medical Journals

The corpus thus involved 43 texts collected from 5 medical journals.

Moreover, most the texts downloaded from the online journals had PDF format and it was necessary to copy them into a word processor. Subsequently, they were saved as a plain text file due to the fact it is the only file format type that AntConc program allows us to use for the analysis of the corpus.

Type of corpus	Specialized
Size	231.989 words
Number of texts	43 texts written by different authors
Medium	Written
Subject	Oncology
Authorship	Texts written by experts
Language	Texts written in American English by experts
Publication date	Texts from 2010 to 2014

The following image shows the main characteristic of the **monolingual corpus**:

Image 2: Characteristics of our corpus

4.1.2. Selection of computer software

As we already know, a corpus can be used as a valuable aid for translators; the question is how to obtain the information from corpora. Once the translators have determined the criteria for the design of their corpus, they now have to decide the computer software that they will use.

There are different software program packages available on Internet, for example: **WordSmith Tools** and **AntConc program**. The former is developed by Mike Scott and the latter by Laurence Anthony. They both provide their users with similar tools to work as researchers of language. However, they differ from each other because the former is not totally free; you just have an online demo version. AntConc is a program that does not require any payment and is downloaded easily. This program offers seven tools, such as: *Word List, Keyword list, Collocates, Clusters, Concordance, Concordance Plot* and *File*

View. Each of them has menu preferences available. The user can use them in a simple way thanks to the tutorials and information that the program provides¹. These tools help the user with the analysis of his corpus.

For all these reasons we thought that AntConc was the appropriate software package chosen for the project.

¹ More information is available at <u>http://www.antlab.sci.waseda.ac.jp/software/README_AntConc3.2.4.pdf</u> (accessed date 09/03/2014).

4.2. Corpus analysis

In this section, we discuss the procedure involved in the analysis of our specialized monolingual corpora. For the analysis, we have used the following tools:

- Word List
- Word Clusters
- Concordances

The first step was to identify the terms that frequently occur in the specialized subject of Oncology. Our goal was to use the terms as "the search node" (Pearson 1998: 191) in order to observe their usage in context. The question that immediately arose was how to extract them. On the Web, there are online tools called **term extractors** that may help the translator with this task. Some of them are free and do not require installation. However, the results that these term extractors provide are not totally suitable. Their disadvantage is the high percentage of noise and silence. We can confirm this drawback because we used online extractors, such as: "Termine" and "Maui-indexer" in our Applied Linguistics classes (Pizarro Sánchez 2013-14). We learned that some of these term extractors presented these previous problems. Bowker and Pearson (2002: 169) state that "noise refers to unwanted items that are erroneously retrieved (i.e. patterns that are *not* terms), while silence refers to cases where patterns that *are* terms do not get retrieved".

This question may be solved easily if the user knows the usages of corpora, since nowadays they can be used as a **term extraction tool**.

4.2.1. Term extraction

We are going to explain in this section how the translator can extract terms from the corpus. Terms are not extracted automatically from the corpus, this process is **semi-automatic**. After the extraction, some candidate terms may not be terms in the end so it is necessary that these term candidates are verified by the user. We speak about term candidates, when we refer to those "words or phrases that appear to be terms" (Bowker and Pearson 2002:145).

As we already know, a corpus contains thousands of words; this means that it contains an extensive amount of information. However, it may hinder our process of identifying term candidates since it may be a more complex and time-consuming process. A good way of reducing the information of the corpus is to use a stop list. It is "a list of items that you would like the computer to ignore" (Pearson 2002:172). In other words, it is a list that includes those words that the user does not want to appear in the frequency lists, such as prepositions, articles, or conjunctions. So we decided to use a stop list² to delete those words that we considered useless to extract.

In specialized language, a term may be a single- word (e.g. 'cell') or multiple words (e.g. 'cancer stem cell'). Pearson (1998: 125) states that "single word terms are the most common, followed by two word noun + noun and adj. + noun combinations". We decided that our analysis is based on her approach. Therefore, we selected the nouns as single candidate terms and noun + noun and adjective + noun combination as multiword term candidates. They were extracted through the tools that were mentioned before.

Pearson (1998:123) points out that "frequency is an important criterion for assessing the eligibility of a term candidate". It is thus advisable that the user establishes the minimum number of times that a lexical unit may occur in the corpus in order for it to be considered as a candidate term. However, Pearson (1998:123) also argues that "low frequency should not prelude a term candidate from being considered". The problem lies in the fact that some terms can be used by their abbreviated form. On the other hand, acronyms are commonly used in medical texts. For that reason, the user is required to pay attention to both aspects.

² The stop list was obtained from <u>http://jmlr.org/papers/volume5/lewis04a/a11-smart-stop-list/english.stop</u> (accessed date 13/03/2014)

The easiest way to identify single - word term candidates is by using a word list, which shows the terms with their frequency of occurrence. It counts how many times each lexical unit appears in the corpus.

According to what is said before, we established that the minimum number of occurrences to be 8 and we used a concordance tool to find those terms that could be omitted. These terms could not be extracted because they are used by their abbreviated form rather than by their full name. To find those terms, we entered as keywords "tomography", "imaging", "survival" and "disease". The following acronyms were found:

- PET (positron emission tomography)
- MRI (magnetic resonance imaging)
- OS (overall survival)
- SD (stable disease)

4.2.1.1. Simple term candidates

The following table shows term candidates that meet the minimum frequency stated above and acronyms found in the corpus that we have compiled. We extracted this list of terms using a word list software program:

SIMPLE TERM CANDIDATES	FREQUENCY
Cell	3332
Tumor	1552
Cancer	1375
Gene	635
Therapy	464
Drugs	459
Tissue	380
Line	369

D - d'-d'	226
Radiation	326
Protein	309
Pathway	299
Lung	282
Factor	255
Mutation	242
Melanoma	239
Imaging	235
Human	231
Chemotherapy	217
DNA	195
Apoptosis	185
Carcinoma	173
Antibody	163
Proliferation	148
Metastasis	145
Stage	145
PET	131
Toxicity	129
MRI	126
Stem	126
Surgery	123
Signals	98
Biopsy	91
Recurrence	87
L	I]

	· · · · · · · · · · · · · · · · · · ·
Screenings	86
Docetaxel	85
Osteosarcoma	79
Blood	77
Lymphoma	74
Phenotype	69
Genome	60
Prognosis	59
Angiogenesis	57
Bone	56
Methylation	52
Adenocarcinoma	50
Marrow	49
Transplantation	49
Enzyme	48
Skin	46
Cisplatin	45
Mortality	43
Antigen	40
Oncogene	40
Staging	40
Leukemia	38
Differentiation	36
Layer	35
Locoregional	33
L	I

Translocation31Transfer28Immunohistochemistry27Myeloma25Nodule25Mastectomy25Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy13Epithelium13Nucleus12Macrophage12	Biomarker	32
Immunohistochemistry27Myeloma25Nodule25Mastectomy25Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Brachytherapy18Immunotherapy17OS15SD14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Translocation	31
Myeloma25Nodule25Mastectomy25Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy13Epithelium13Nucleus12	Transfer	28
Nodule25Mastectomy25Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Immunohistochemistry	27
Mastectomy25Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy13Epithelium13Nucleus12	Myeloma	25
Resection24Brain24Infusion23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Nodule	25
Brain24Infusion23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy13Epithelium13Nucleus12	Mastectomy	25
Infusion23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Resection	24
Dissection23Dissection23Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Brain	24
Chromosome18Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Infusion	23
Necrosis18Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Dissection	23
Carcinogenesis18Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Chromosome	18
Brachytherapy18Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Necrosis	18
Immunotherapy17OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Carcinogenesis	18
OS15SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Brachytherapy	18
SD14Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	Immunotherapy	17
Carcinogen14Lumpectomy14Prostatectomy13Epithelium13Nucleus12	OS	15
Lumpectomy14Prostatectomy13Epithelium13Nucleus12	SD	14
Prostatectomy 13 Epithelium 13 Nucleus 12	Carcinogen	14
Epithelium 13 Nucleus 12	Lumpectomy	14
Nucleus 12	Prostatectomy	13
	Epithelium	13
Macrophage 12	Nucleus	12
	Macrophage	12

Chemoprevention	11
Colonoscopy	11
X-ray	10
Implantation	10

Image 3: Single term candidates extracted from a word list

4.2.1.2. Multiword term candidates

On the other hand, word cluster software programs enable the user to detect multiword term candidates. The selection of multiword term candidates may be based on a statistical approach (Bowker and Pearson 2002:170) which consist of establishing the minimum number of times that a series of lexical units must be repeated in the corpus. This minimum value should be indicated by the user, as well as the minimum and maximum length of lexical units in order to indicate multiword term candidates.

Our selection was based on this approach and we established that the minimum number will be 8 for the times that a series of lexical units are repeated and the minimum and maximum will be 2, 3, and 4 for the length of words. We indicated several values for the length of words in order to identify more term candidates.

MULTIWORD TERM CANDIDATES	FREQUENCY
Cell line	292
Breast cancer	205
Lung cancer	171
Radiation therapy	159
Tumor cells	156
In vivo	140
In vitro	123
Prostate cancer	111
Lymph node	108
Melanoma cell	99
Stem cell	79

The image below shows the cluster list:

Primary tumor	65
Signaling pathway	65
Clinical trials	62
Cell proliferation	53
Pancreatic cancer	50
Bone marrow	47
Cell cycle	45
Tumor tissue	42
Tumor volume	36
Drug resistance	35
Stem cell transplantation	34
Cell type	32
Risk factor	32
Cancer stem cell	31
DNA repair	31
Squamous cell carcinoma	30
DNA methylation	27
Partial response	26
Cell migration	24
Colon cancer	23
progression-free survival	22
Cell culture	21
Blood vessels	20
Survival rate	20
Local recurrence	20

Locoregional recurrence	19
Colorectal cancer	17
Multiple myeloma	16
Side effects	15
In situ	14
Non – small cell lung cancer	12
Tumor suppressor gene	10
Carcinoma in situ	9
CT scan	8

Image 4: Multiword term candidates extracted from a cluster list

There is another slower way to identify more term candidates. The user just takes the single – word candidates or multiword candidates and then uses a concordance tool. This tool allows him to identify lexical units associated to other lexical units, which may constitute a term.

4.2.1.3. Identification of key terms

As the process is semi- automatic, the user has the task to indicate which term candidates are the key terms and eliminate those lexical units that are not. **Specialized glossaries** are essential for him/her to achieve this assignment.

We searched specialized glossaries on the Web and selected those that are related with the field of oncology. They helped us to verify which term candidates were the key terms. The glossaries were found on these websites: Genomic Health (www.genomichealth.com), Cancer Monthly (www.cancermonthly.com), American Cancer Society (www.cancer.org) and American Society for Cancer Research (www.aacr.org)³. They are official medical centers and institutions from the US which give information about causes, treatment, prevention of cancer, etc. to people who suffer this disease or those who are interested in the topic. The information included in these sources comes from experts and it is therefore reliable. It was the main reason why they were selected as reference for indicating the key terms.

The following image shows our term candidates that occur in these previous glossaries

Term candidates	Genomic Health	Cancer Monthly	American Cancer Society	American Society for Cancer Research
Cell				\checkmark
Tumor				N
Cancer		N		N
Gene				N
Drugs				
Stem				
Human				
Tissue				
Therapy				
Line				

³ Accessed date: 27/03/2014

Protein	1			
Pathway	+			
Lung				
Radiation				
Mutation				1
Melanoma	+			× ·
Factor				
Imaging				1
Chemotherapy				1
DNA		V		N
		N		
Apoptosis				N N
Carcinoma		N		
Antibody				1
Proliferation				
Metastasis			1	
Stage				1
PET			V	
Toxicity				V
MRI			\checkmark	
Surgery			N	
CT scan			N	
Signals				
Recurrence	1			1
Biopsy	1	V	N	
Docetaxel				
Lymphoma			N	
Blood				
Screenings	1			
Phenotype				1
Clone				
Bone				
Prognosis			N	
Genome			V	
Angiogenesis		1	V	1
Marrow			-	-
Transplantation				
Methylation				
Skin				
Adenocarcinoma	1	1		
Enzyme			1	
Cisplatin				
Antigen	1	1	V	1
Oncogene				1
Staging	V	1	1	-
		1	-	

Differentiation		V		
Layer		-		
Locoregional		1		
Biomarker		1		
Translocation		-	1	V
Mortality			v.	
Transfer				
Brain				
Immunohistochemistry			1	
Myeloma				
Nodule			1	
Mastectomy		1	v v	
Chromosome		1	•	1
Infusion		1		,
Necrosis		1	1	
Carcinogenesis		N.	Y	
Brachytherapy		1		
OS		V		1
SD				V
		N		
Carcinogen			N	N
Prostatectomy			V	
Epithelium			,	N
Immunotherapy			V	N
Dissection			V	
Nucleus			N	V
Macrophage			N	
Lumpectomy				
Chemoprevention			\checkmark	\checkmark
Colonoscopy	\checkmark	\checkmark	\checkmark	
X-ray			\checkmark	
Resection	V		V	
Implantation				
Cell line				
Tumor cells				
Lung cancer				
In vivo				
Breast cancer				\checkmark
In vitro				
Prostate cancer				
Lymph node			\checkmark	
Melanoma cell				
Radiation therapy		\checkmark		
Stem cell				
Primary tumor		1		

Signaling pathway				
Clinical trials	V	1	1	\checkmark
Cell proliferation		N		
Bone marrow		N		
Pancreatic cancer				
Cell cycle			N	\checkmark
Tumor tissue				
Multiple myeloma			1	
Tumor volume			1	
Drug resistance			N	\checkmark
Cell type				
Stem cell			1	
transplantation				
Tumor suppressor		\checkmark	V	
gene				,
Cancer stem cell				\checkmark
DNA repair			\checkmark	
Squamous cell			\checkmark	\checkmark
carcinoma				,
DNA methylation				N
Partial response		\checkmark		\checkmark
Risk factor			\checkmark	
Colon cancer		\checkmark		
Cell migration				
Progression-free				\checkmark
survival				
Cell culture				
Blood vessels				
Survival rate		N	\checkmark	
Locoregional				
recurrence	,			,
Local recurrence	V		,	V
Colorectal cancer		- , · · · · · · · · · · · · · · · · · ·	N	V
Non – small cell lung		N	\checkmark	V
cancer		-		ļ
Osteosarcoma		N	,	L ,
Side effects	,		1	V
In situ	V		N	
Carcinoma in situ	\checkmark	\checkmark	\checkmark	V

Image 5: Term candidates extracted from online medical glossaries

Therefore, the term candidates that we considered key terms are "cell", "tumor", "cancer", "gene", "mutation", "imaging", "chemotherapy", "DNA", "apoptosis", "carcinoma", "antibody", "metastasis", "stage", "PET", "toxicity", "MRI", "surgery", "recurrence", "screenings", "phenotype", "biopsy", "lymphoma", "prognosis", "genome", "adenocarcinoma", "enzyme", "antigen", "oncogene", "angiogenesis", "staging", "leukemia", "differentiation", "biomarker", "translocation", "mortality", "immunohistochemistry", "nodule", "mastectomy", "chromosome", "infusion", "necrosis", "carcinogenesis", "brachytherapy", "OS", "SD", "carcinogen", "prostatectomy", "epithelium", "immunotherapy", "dissection", "nucleus", "macrophage", "lumpectomy", "chemoprevention", "colonoscopy", "X-ray", "resection", "osteosarcoma", "breast cancer", "lymph node", "radiation therapy", "stem cell", "primary tumor", "clinical trials", "cell proliferation", "bone marrow", "cell cycle", "multiple myeloma", "tumor volume", "drug resistance", "stem cell transplantation", "tumor suppressor gene", "cancer stem cell", "DNA repair", "squamous cell carcinoma", "DNA methylation", "partial response", "risk factor", "colon cancer", "progression-free survival", "survival rate", "local recurrence", "colorectal cancer", "non-small cell lung cancer", "side effects", "in situ", "carcinoma in situ" and "CT scan".

4.2.2. Collocations

Olohan (2004:63) states that a concordance software program is "the most common tool for data extraction". She also points that "each instance is displayed with its immediate co-text" (2004:63). As we indicated previously, the terms comprised in image 5 will be used as the keywords. Therefore, the next step in our analysis was to use the terms as the search pattern to extract the information about their usage in context. We used a concordance software program for our analysis.

We believe that this tool has two main advantages: The first one is that it is a helpful tool for the translators because they may verify whether the chosen translation equivalent that they have found in dictionaries or glossaries is the right one or not. And second, this tool allows the users to acquire knowledge of target language (TL) such as: new expressions. Our analysis focused on the combinations formed by **verb** + **noun**, **noun**+ **verb**, **adjective**+ **noun** and **noun** + **noun**, in which the noun is referred to as the base and the verb, the adjective and the noun as the collocates. Similarly, we established a minimum of frequency to select collocations. This means that the candidate collocates should have a have a minimum of frequency of 3 occurrences to be considered as such (Méndez Cendón, 2002).

Collocates can be defined as those words that go along with the search term, for instance:

"Mammographic screening"

This example was extracted entering the search term "screening".

A table shows our analysis, which is composed of three sections: the first section is called "term" in which the single-word terms and multiword terms are included, the second one comprises their collocations, and the last section is called notes. This section was created to

introduce all information that we considered necessary for the translation task, such as: synonyms, abbreviated form and compound terms.

Tourn		Natas			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
					→ Abbreviated form
					√ Synonym
					* Compound terms

Image 6: Table of information regarding the terms

As shown in Image 6, we used three different symbols to distinguish the information that will be included in Notes. The first symbol is \dashv to indicate abbreviated form, the second one is $\sqrt{}$ to point synonyms and the last one is * to show compound terms.

To begin with our analysis, we first use single - word terms as the search pattern and then we obtain their KWIC concordances lines.

		Coll	locations		
Term	Verb+	Term+	Adj.+ term	Term+ noun	Notes
	term	verb			
Cell	Induce, affect, treat, control, promote, inhibit, decrease, cause, use, enhance, initiate, produce	Treat, use consist, produce, undergo, increase, contain, activate	Leukemic endothelial, epithelial, basal, blast, primary, apoptotic, stromal*, myeloid, cancerous, dentritic*, malignant, clonogenic, mononuclear, mesenchymal,	line, type, division, death, in vitro, invasion, culture, carcinoma, differentiation, growth, adhesion, interaction, migration, apoptosis, viability, membrane, resistance, phenotype, generation,	*Bone marrow stromal cells (BMSCs) *Tumor dendritic cells (TUDCs) *Melanoma initiating cells (MIC) *Senescent melanoma cells (SSMC) *Human umbilical cord
				lymphoma, sensitivity	endothelial cells

			(HUVECs)
Table 1: Collocates of the	search term cel	11	

			Collocations		
Term	Verb+	Term+	Adj.+ term	Term+ noun	Notes
	term	verb			
Tumor	Generate, form, treat reduce, cause, inhibit, promote, initiate	Arise, use produce, reach,	gastric, recurrent, circulating*, solid, colonic, colorectal, epithelial, human, aggressive, immunosuppressive, metastatic, oral, subcutaneous, pancreatic, original, orthotopic, intestinal,	Cell, vessels, recurrence, DNA, clone, tissue, mass, malignancy, angiogenesis, invasiveness, stroma, stage, grown, nodule, environment, formation, development, necrosis, resection, expansion, killing,	*Circulating tumor cells (CTCs)

 Table 2: Collocates of the search term *tumor*

	Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Cancer	Develop, prevent, identify, detect, increase, inhibit	Cause, occur, include, increase	Esophageal, epithelial, gastric, hematologic, oesophagogastric, pancreatic, rectal, ovarian, invasive, human, malignant,	Cell, tissue, resistance, treatment, death, tumor, therapy, development, screening, growth, detection, xenograft, progression,		

				mortality						
Table 3: Coll	locates of the	search term ca	incer	Table 3: Collocates of the search term <i>cancer</i>						

Collocations Term Notes Term+ verb Term+ noun Verb+ term Adj.+ term Gene Use, Contain Circadian, Mutation, *Gene ontology identify, housekeeping fusion, ontology*, (GO) expression*, *Gene regulation expression profiling (GEP)

Table 4: Collocates of the search term gene

Term		Notes			
1 et m	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TAULES
Mutation	Identify,	Identify,	Somatic,	Detection	
	active,	occur	oncogenic,		
			intronic,		
			genetic		

Table 5: Collocates of the search term *mutation*

Term	Verb+	Term+	Adj.+ term	Term+	Notes
	term	verb		noun	
Imaging	Use	Play, detect, perform	Spectroscopic, photoacoustic,	Test, agent, modality	
		periorin	multiparametric,	modulity	
			intravital,		
			optical,		
			molecular,		
			morphologic		

Table 6: Collocates of the search term *imaging*

Term	Verb+	Term+	Adj.+ term	Term+ noun	Notes
	term	verb			
Chemotherapy	Reduce,	Provide	Neoadjuvant,	Treatment,	
	resist,		neurotoxic	drug, agent,	
				resistance	

Table 7: Collocates of the search term *chemotherapy*

	Collocations				
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
DNA	use, increase,		Genomic	Hypermethylation, sample, damage, fragmentation microarray	

Table 8: Collocates of the search term *DNA*

Term		Collocations					
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes		
Apoptosis	Induce,		Induced,	Assay			
	trigger,		endothetial	_			
	mediate						

 mediate

 Table 9: Collocates of the search term *apoptosis*

Term	Verb+	Term+	Adj.+ term	Term+	Notes		
	term	verb		noun			
Carcinoma			lobular,	Cell	*Invasive		
			gastric		ductal		
			ovarian,		carcinoma		
			invasive,		(IDC)		
			mammary,		*Basal-cell		
			hepatocellular,		carcinoma		
Table 10: Co	Table 10: Collocates of the search term <i>carcinoma</i>						

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INULES
Antibody	Control, use	Use	Monoclonal,		
-			polyclonal,		
			polyclonal, secondary		
			primary		

Table 11: Collocates of the search term *antibody*

Term		Colloc	cations		Notes
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Metastasis					

Table 12: Collocates of the search term *metastasis*

Term		Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes		
Stage			Clinical,	Melanoma,			
			early,	disease, III,			
			unresectable,	IV, I			
			advanced				

Table 13: Collocates of the search term *stage*

Term		Notes			
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INDIES
PET			Interim,	Scan,	√Positron
			early	imaging,	emission
				finding	tomography

Table 14: Collocates of the search term *PET*

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Toxicity	Reduce, minimize		Cardiac, systemic, haematological		

 haer

 Table 15: Collocates of the search term *toxicity*

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INOLES
MRI	Undergo,	Use	Enhanced	Screening	√Magnetic
	use				resonance
					imaging

Table 16: Collocates of the search term MRI

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Surgery	Undergo		Excisional		*breast-
					conserving
					surgery

Table 17: Collocates of the search term *surgery*

Term		Notes			
ICIII	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THOLES
Recurrence			Local,		
			biochemical,		
			locoregional,		

Table 18: Collocates of the search term *recurrence*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TIOLES
Biopsy	Use	Yield,	Human	specimens	
		control			

Table 19: Collocates of the search term *biopsy*

Term		Colloc	cations		Notes
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Lymphoma			Anaplastic, lymphocytic	Cell	

Table 20: Collocates of the search term lymphoma

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Screenings			Mammographic	Mammography	

 Table 21: Collocates of the search term screenings

Term	Verb+	Term+ verb	Adj.+ term	Term+ noun	Notes
	term				
Phenotype	Rescue		Tumorigenic,		
			invasive,		
			antiapoptotic,		
			stemness,		
			mesenchymal,		
			epithelial,		

Table 22: Collocates of the search term *phenotype*

Torm	Term					
1 erm –	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Prognosis	Confer		Dismal			

Table 23: Collocates of the search term *prognosis*

Term		Notes			
Verb	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INOLES
Genome	Perform, use		Human,		*Whole
			whole*		genome
					sequencing
					(WES)

 Table 24: Collocates of the search term genome

Term		cations		Notes	
ICIIII	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Angiogenesis	Induce,				
	inhibit				

Table 25: Collocates of the search term *angiogenesis*

		(
Term	Verb+	Term+	Adj.+ term	Term+	Notes
	term	verb		noun	
Adenocarcinoma			Pancreatic, ductal,		*Pancreatic
			gastric,		ductal
			oesophagogastric,		adenocarcinoma
			gastrooesophageal,		(PDA)
					*Adenocarcinoma
					in situ (AIS)

 Table 26: Collocates of the search term adenocarcinoma

Term		Notes			
I CI III	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TULES
Enzyme					

Table 27: Collocates of the search term enzyme

Term		Notes			
I CI III	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Antigen	Make		specific*	Retrieval	*Prostate –
					specific
					antigen (PSA)

Table 28: Collocates of the search term antigen

Term		Notes			
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INDIES
Oncogene					

Table 29: Collocates of the search term *oncogene*

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INULES
Staging			Accurate		

Table 30: Collocates of the search term *staging*

Term	Verb+	Term+	Adj.+ term	Term+	Notes
	term	verb		noun	
Leukemia			Lymphoblastic*,		*Chronic
			myeloid,		lymphoblastic
					leukemia
					(CLL)
					*Acute
					lymphoblastic
					leukemia
					(ALL)

Table 31: Collocates of the search term leukemia

Term	Verb+	Term+	Adj.+ term	Term+	Notes
	term	verb		noun	
Differentiation			Adipogenic		

 Table 32: Collocates of the search term differentiation

Term	Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Biomarker						

Table 33: Collocates of the search term *biomarker*

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Translocation			Nuclear, chromosomal		

 Table 34: Collocates of the search term *translocation*

Term		Notes			
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THORES
Mortality					

Table 35: Collocates of the search term *mortality*

Term	Verb+	Term+	Adj.+	Term+	Notes
	term	verb	term	noun	
Immunohistochemistry					JHC

Table 36: Collocates of the search term immunohistochemistry

Term		Notes			
ICIIII	Verb+ term Term+ verb Adj.+ term Term+ noun				
Nodule					

Table 37: Collocates of the search term *nodule*

Term		Notes			
ICIIII	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Mastectomy	Undergo				

Table 38: Collocates of the search term *mastectomy*

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TOLES
Chromosome					

 Table 39: Collocates of the search term chromosome

Term		Notes			
ICIIII	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THOLES
Infusion			intravenous		

Table 40: Collocates of the search term infusion

Term		Notes			
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Necrosis					

Table 41: Collocates of the search term necrosis

Term	Verb+	Term+	Adj.+	Term+	Notes
	term	verb	term	noun	
Carcinogenesis					

 Table 42: Collocates of the search term carcinogenesis

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Brachytherapy					

 Table 43: Collocates of the search term brachytherapy

Term		Notes			
ICIIII	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
OS					√Overall
					survival

Table 44: Collocates of the search term *OS*

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INDIES
SD					√Stable
					disease

Table 45: Collocates of the search term SD

Term		Notes			
Verb+ term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Carcinogen					

 Table 46: Collocates of the search term carcinogen

Term	Verb+	Term+ verb	Adj.+ term	Term+ noun	Notes
	term				
Prostatectomy			Radical		

Table 47: Collocates of the search term *prostatectomy*

Torm	Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Epithelium						

Table 48: Collocates of the search term gene

		Colloc	cations		
Term	Verb+	Term+	Adj.+	Term+	Notes
	term	verb	term	noun	
Immunotherapy					

 Table 49: Collocates of the search term immunotherapy

Term		Notes			
1 erm	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Dissection	Undergo				

Table 50: Collocates of the search term dissection

Term	Collocations					
Verb+ term		Term+ verb	Adj.+ term	Term+ noun	Notes	
Nucleus						

Table 51: Collocates of the search term *nucleus*

Term		Notes			
101111	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TABLES
Macrophage			associated*		*Tumor-
					associated
					macrophages
					(TAMs)

Table 52: Collocates of the search term *macrophage*

Term		Colloc	ations		Notes
lerm	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TULES
Lumpectomy					

Table 53: Collocates of the search term *lumpectomy*

Term	Verb+	Term+	Adj.+	Term+	Notes
	term	verb	term	noun	
Chemoprevention					

Table 54: Collocates of the search term *chemoprevention*

Term		Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes		
Colonoscopy							

Table 55: Collocates of the search term *colonoscopy*

Term	Term					
1 erm –	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
X-ray						

Table 56: Collocates of the search term X-ray

Term		Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes		
Resection	Undergo		Surgical				

Table 57: Collocates of the search term resection

Term	Verb+	Term+ verb	Adj.+ term	Term+ noun	Notes
	term				
Osteosarcoma	Isolate		Human,	Cell, stem	
			primary	cell	

Table 58: Collocates of the search term osteosarcoma

The following tables illustrate that complex noun groups⁴ are used as the search pattern.

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INUICS
Breast cancer	Increase		invasive,	Cell, risk	
			human,	mortality,	
				screening,	
				xenograft,	

Table 59: Collocates of the search term *breast cancer*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	NOLES
Lymph node			Axillary,	Metastasis,	
			sentinel,	biopsy,	
			regional,	dissection,	
			mammary,	invasion,	
			metastatic,		

Table 60: Collocates of the search term *lymph node*

⁴ We use the term 'complex noun group' to designate not only compound terms but also those multiword combinations which can designate more than one concept in the medical field.

		Collocations				
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Radiation therapy	Receive		Optimized, modulated, postoperative	regimen	, JRT √Radiotherapy, *Stereotactic ablative radiation therapy (SABR) *Intensity- modulated radiation therapy (IMRT)	

Table 61: Collocates of the search term *radiation therapy*

		Collocations					
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes		
Stem cell			Embryonic*, mesenchymal*	Transplant, markers	*Embryonic stem cells (ES cells) *Mesenchymal stem cells (MSCs)		

 Table 62: Collocates of the search term stem cell

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Primary				Cell, sample,	
tumor					

 Table 63: Collocates of the search term primary tumor

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Clinical trials			Randomized, early		

 Table 64: Collocates of the search term *clinical trials*

Term		Collocations						
ICIM	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes			
Cell	Inhibit							
proliferation								

Table 65: Collocates of the search term *cell proliferation*

Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Bone marrow			Normal, healthy	Microenvironment*, biopsy	*Bone marrow microenvironment (BMME)

Table 66: Collocates of the search term *bone marrow*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	NOLES
Cell cycle	Control			Regulation,	
				distribution,	
				progression,	

 Table 67: Collocates of the search term *cell cycle*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THOLES
Tumor		Calculate			
volume					

Table 68: Collocates of the search term *tumor volume*

Term	Verb+	Term+ verb	Adj.+ term	Term+ noun	Notes
	term				
Drug			Mediated*,		*Environment
resistance					- mediated
					drug
					resistance
					(EMDR)

 Table 69: Collocates of the search term *drug resistance*

Term	Verb+	Term+	Adj.+ term	Term+	Notes
	term	verb		noun	
Stem cell transplantation	Undergo		Haemotopoietic, autologous*		*Autologous stem cell transplantation (ASCT)

Table 70: Collocates of the search term stem cell transplantation

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Cancer stem					JCSC
cell					

Table 71: Collocates of the search term *cancer stem cell*

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Hotes
DNA repair	Increase		Enhanced	Efficiency	

 Table 72: Collocates of the search term DNA repair

		Collocations				
Term	Verb+	Term+	Adj.+ term	Term+	Notes	
	term	verb		noun		
Squamous			Oropharyngeal,		JSCC	
cell			oral		*Oropharyngeal	
carcinoma					squamous cell	
					carcinoma	
					(OSCC)	
					*Oral	
					squamous cell	
					carcinoma	
					(OSCC)	

 Table 73: Collocates of the search term squamous cell carcinoma

Term		Colloc	cations		Notes
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TULES
DNA					
methylation					

Table 74: Collocates of the search term DNA methylation

Term		Notes			
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Partial	Achieve				, ∠PR
response					

Table 75: Collocates of the search term *partial response*

Term		Colloc	cations		Notes
lerm	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	TULES
Risk factor			Cardiac		

Table 76: Collocates of the search term *risk factor*

Term			Notes		
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Colon cancer					

Table 77: Collocates of the search term *colon cancer*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes
Progression-			Median		JPFS
free survival					

Table 78: Collocates of the search term *progression-free survival*

Term		Notes			
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THOLES
Local					
recurrence					

Table 79: Collocates of the search term *local recurrence*

Term		Collocations				
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	Notes	
Survival rate		Improve	Overall			

 Table 80: Collocates of the search term survival rate

Term		Colloc	cations		Notes
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INOLES
Colorectal					, LCRC
cancer					

 Table 81: Collocates of the search term colorectal cancer

Term		Colloc	cations		Notes
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INOLES
Multiple					JMM
myeloma					

Table 82: Collocates of the search term *multiple myeloma*

Term		Colloc	cations		Notes
	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Non-small					JNSCLC
cell lung					
cancer					

Table 83: Collocates of the search term non-small cell lung cancer

Term		Colloc	cations		Notes
1 ci m	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
Side effects			Severe		

Table 84: Collocates of the search term side effects

Term		Colloc	ations		Notes
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	notes
In situ					

Table 85: Collocates of the search term in situ

Term		Colloc	Notes		
Term	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
CT scan		Allow			

Table 86: Collocates of the search term CT scan

Term		Colloc	cations		Notes
10111	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	THUES
Tumor					
suppressor					
gene					

 Table 87: Collocates of the search term tumor suppressor gene

Term		Colloc	cations		Notes
1 CI III	Verb+ term	Term+ verb	Adj.+ term	Term+ noun	INULES
Carcinoma in			Ductal*,		*Ductal
situ			lobular		carcinoma in
					situ (DCIS)

 Table 88: Collocates of the search term carcinoma in situ

5. Specialized monolingual corpora in translation activity

In this section we are focusing on demonstrating the role that played the information extracted from KWIC concordance lines in translation activity. Therefore, the importance of compiling a specialized monolingual corpus is described. In order to indicate the usefulness of the corpus as a tool for translators, we selected a sample of the article entitled *"Tratamiento perioperatorio del cáncer gástrico"* published in the journal of *Revisiones en*

Cáncer (Madrid 2013).

Las dificultades del desarrollo de esta estrategia terapéutica estriban por una parte en las limitaciones de las técnicas diagnósticas preoperatorias y, por otra parte, en la ausencia de fármacos con suficiente actividad antineoplásica. Los estudios en monoquimioterapia no ofrecen unas tasas de respuesta suficientes para generar evidencias (1) (Tabla I). Es por ello que el tratamiento del cáncer de estómago, tanto en la enfermedad avanzada como en el tratamiento neoadyuvante de la enfermedad locorregional se basan en la combinación de diferentes fármacos que hacen posible la resecabilidad de los tumores localmente avanzados o irresecables al diagnóstico, obteniendo una alta tasa de cirugías radicales (CR0).

TABLA I

MONOTERAPIA EN CANCER GASTRICO DISEMINADO (1	MONOTERAPIA	EN CÁNCER	GÁSTRICO	DISEMINADO (1
--	-------------	-----------	----------	---------------

Fármaco	п	Respuesta objetiva (%)
5 fluorouracilo	416	21
UFT	188	28
S-1	101	45
Mitomicina	211	30
Adriamicina	141	17
Epirubicina	80	19
Cisplatino	139	19
Etopósido	35	20
Docetaxel	157	22
Paclitaxel	82	15
Irinotecán	66	23

ESTUDIOS FASE II

Ajani y cols. (2) desarrollaron un estudio fase II en el que se incluyeron 48 pacientes con cáncer gástrico

Image 7: Spanish article

localmente avanzado potencialmente resecable, a los cuales se les administraban 3 ciclos de quimioterapia preoperatoria basada en etopósido, doxorubicina y cisplatino, y dos ciclos adicionales con el mismo esquema de quimioterapia tras la cirugía, consiguiendo que el 85% de los pacientes fueran intervenidos y en el 77% de los casos, lográndose una cirugía radical CR0.

Estos mismos investigadores realizaron un segundo estudio, utilizando el mismo régimen, pero esta vez administrando dos ciclos de quimioterapia previos a la cirugía y tres posoperatoriamente con un 100% de cirugías, de las cuales un 72% fueron CR0 (3).

Otro estudio fase II de quimioterapia neoadyuvante en pacientes con cáncer gástrico localmente avanzado, conducido por Rougier (4), se planteaba la administración de quimioterapia previa a la cirugía con 5-fluorouracilo (5-FU) en infusión continua (IC) y cisplatino. Un 56% de los pacientes presentó <u>lma respuesta</u> objetiva y el 60% de los casos pudieron ser <u>sometidos</u> a una cirugía radical con intención curativa. La mediana de supervivencia a los tres años fue del 38%, superior a los controles históricos, y en relación a la tolerancia, este estudio, al igual que en los estudios fase II anteriormente citados, objetivó que el uso de tratamiento quimioterápico neoadyuvante era bien tolerado y no comportaba un aumento de la morbimortalidad operatoria.

En el estudio de Wilke (5), 34 pacientes con cáncer gástrico localmente avanzado y no resecables, estadiados por laparoscopia, recibieron tratamiento quimioterápico neoadyuvante con etopósido, adriamicina y cisplatino. Posteriormente, y después de un «second look» se realizaba la cirugía radical en aquellos pacientes que presentaban una respuesta parcial o completa de la enfermedad. Se administraban 2 ciclos más de la misma quimioterapia después de la cirugía. La toxicidad fue aceptable, sin incremento en la morbimortalidad, consiguiendo una mediana de supervi-

As shown in Image 2, we selected those words and verbs that may present a problem for the translator. The words are highlighted in yellow while the verbs are underlined in red color.

Once we chose what we considered translation difficulties, our purpose was to reach a solution. Some words underlined are terms that we had already found in our previous analysis. It is the case of "respuesta parcial" (partial response), "toxicidad" (toxicity), "fármacos" (drugs), "infusión" (infusion), "quimioterapia" (chemotherapy) and "cáncer" (cancer).

However, our task was to translate "infusión continua," "quimioterapia preoperatoria" and "cáncer gástrico." What we did was to analyze their collocations on the right to find the appropriate adjective. The adjective "gastric" was found for cancer and therefore we considered that the translation equivalence was "gastric cancer".

But we did not find any adjective for "preoperatoria" and "continua". We trusted our intuition since we had studied these two adjectives in the subject "Specialized Translation". Therefore, our guess was that the equivalences might be "preoperative" and "continuous". We extracted concordance lines to confirm our hypothesis.

Concord	ance	Concordance Plot	File View	Clusters	Collocates	Word List	Keyword List	
Hit	KWIC							
1	NSA	BP B18 and B27 r	andomized	trials of	preoperat:	ive system	ic therapy,	in who
2	he e	radication of ly	mph node (disease by	preoperat:	ive system	ic therapy.	In sum
3	s 1)	whole-brain tum	or locali:	zation for	preoperat:	ive and in	traoperative	macro
4	micr	ocalcifications.	The role	of MRI in	preoperat:	ive stagin	g is controv	ersial
5	tice	s in CRC with re	spect to	screening,	preoperat:	ive evalua	tion, survei	llance
6	Alth	ough CT and MRI	are widel	y used for	preoperat:	ive whole-	body staging	, they

Image 8: Concordance for the search term *preoperative*

Concord	ance	Conc	ordance Plot	File View	Clusters	Collocates	Word List	Keyword List	
Hit	KWIC								
1	ir u	se as	treatment	and will	require a	a continuous	infusion	or frequent	administration.

Image 9: Concordance for the search term continuous infusion

The images above illustrate that these adjectives are used in research articles. Therefore, "preoperative chemotherapy" and "continuous infusion" were the right equivalents. Now,

our task was to translate "tratamiento neoadjuvante". Our intuition told us that its equivalence could be "neoadjuvant treatment". In order make sure that it was right, we extracted concordance lines with "neoadjuvant" as the keyword but there were no results.

As we can observe the image below, this compound noun in English is "neoadjuvant therapy".

Concord	Ince Concordance Plot File View Clusters Collocates Word List Keyword List
Hit	KWIC
1	ch could have been given as adjuvant or neoadjuvant therapy, or for advanced
2	1 has a 5-year survival rate below 20%. Neoadjuvant chemotherapy has been pro
3	o have received systemic therapy in the neoadjuvant setting. Because the rese
4	in patients undergoing mastectomy after neoadjuvant chemotherapy. These studi
5	of disease both at diagnosis and after neoadjuvant chemotherapy are relevant
6	ience a pathologic complete response to neoadjuvant chemotherapy appeared to
7	nced breast cancer (Fig. 2) who undergo neoadjuvant therapy before breast can
8	on of extent of disease before starting neoadjuvant chemotherapy. Sagittal fa
9	the intensity of [18F]FDG uptake during neoadjuvant therapy, as determined by
10	vival very low. For the last 4 decades, neoadjuvant chemotherapy and limb sal

Image 10: Concordance for the search term *neoadjuvant*

Moreover, "enfermedad" may be a problem for the translators since it can be translated into English as "illness, "sickness" or "disease." In order to avoid mistakes, the KWIC concordance showed that "disease" had a higher frequency of occurrence than the others. Thus, it is term used by experts in research articles. The following image shows that "locoregional disease is the appropriate equivalent term for "enfermedad locorregional".

Concord	Concordance Plot File View Clusters Collocates Word List Keyword List							
Hit	KWIC							
1	tastases from persistent reservoirs of locoregional disease. In recent years,							
2 3	tastases from persistent reservoirs of locoregional disease. In recent years,							
3	ve the same risk of harboring residual locoregional disease after mastectomy a							
4 5	ents might harbor a burden of residual locoregional disease that systemic the							
5	se-free survival and a 2.3% benefit in locoregional disease-free survival from							

Image 11: Concordance for the search term *locoregional disease*

As for, "intención" it has two English equivalents: "intent" or "intention". The task was to translate "intención curativa" but we did not know what equivalent to choose. We thought of two possibilities of translation and used them to generate concordances. "Curative intention" was the first option but we did not find results, whereas an occurrence was found for "curative intent as the image below shows.

Concordance		Conco	ordance Plot	File View	Cluste	rs (Collocates	Word List	Keywo	ord List	
Hit	KWIC										
1	erapy	and	radiothera	py adminis	stered	with	curative	intent.	A meta-	analysis	of

Image 12: Concordance for the search term *curative intent*

The need of knowing the translation equivalents into English leads the translators to consult a glossary or a dictionary. For that reason, we used Linguee⁵ that is a tool consisting of a dictionary and a search engine. As it gives examples of usage, we think that the translators may consider it essential for their translation task.

We entered the search term "estudio fase II" and as the image 2 illustrates, it offers two possible equivalents: "phase II trial" and "study phase II". Our doubt was whether the right equivalent was for a translation. Therefore, Linguee did not help us to solve our problems. So, we had to generate concordances with the possible translation equivalents until we finally found the right one. After this process, we observed that the right equivalent is "phase II trial".

⁵ <u>http://www.linguee.es</u> (accessed date 24/04/2014)

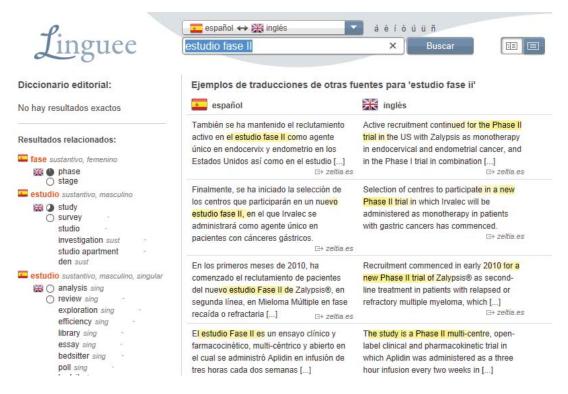


Image 13: Results for the search estudio fase II

Another online translation aid tool is the Proz glossary⁶. Our search term was "morbimortalidad". As we can see in the image 13 it offers two possibilities but we did not know which equivalent should be used. Then, we decided to extract KWIC concordance lines for both "morbidity and mortality" and "morbimortality". There were two results for the first possibility whereas there were no results for the second one. Thus, our conclusion was that the translator should use "morbidity and mortality" in his translation.

⁶ <u>http://www.proz.com/search/</u> (accessed date 25/04/2014)

4 total results	KudoZ:	1 KOG: 1	Glossaries: 1 Glosspost: 0 G	вк:0 \	Vikiwords: 1		
KudoZ (1	results)						
	Field		Term		Asked by		
(Medical)		morbimortalida	d (2)	valcobar			
1 results found	. Showing resu	ts 1 - 1					
		ary (1 resu	hown in parentheses next to the term				
(Media	:al)	morbimortalidad	* morbimortality		valcobar		
1 results found	-	s (1 results	.)				
Glossary		Term	* translation		Owner		
Nursing	morbimorta	lidad 🕈 morbidity	and mortality	🞗 Nigel Greenwood			
1 results found	1 results found. Showing results 1 - 1						
	ch so you may d properly.	see some repeate	t all entries are stored in Unicode. A few d results. Steps will soon be taken to co				
Language(s))	Fields	Term translation		Contributor		
Spanish to English			morbimortalidad → morbimortality	r	ı/a		
1 results found	Showing resul	ts 1 - 1					

Image 14: Results for the search morbimortalidad

The last compound noun to translate was "mediana de supervivencia". In this case, we tried to make a guess again and think about "median survival" as the equivalent term. The following concordance lines shows that our guess was right.

Concor	dance Concordance Plot File View Clusters Collocates Word List Keyword List
Hit	KWIC
1	oradiotherapy, with a 10.2 month improvement in median survival. By contrast, no benefit
1 2 3 4 5 6	he SWOG 002325 trials suggest that the 20 month median survival estimate for the control
3	1 benefits: It led to a substantial increase in median survival, from 67 days in the coh
4	pected, the U87 sh-control group showed a short median survival (16 days; Figure 3B). In
5	s (Figure 5C). As expected, ?EGFR decreased the median survival of nude mice when compar
6	death from disease ($P < .001$) and with a lower median survival ($P = .020$; Figure 1, C as

Image 15: Concordance for the search term *median survival*

After we have solved the problems regarding terms and compound nouns previously, our task was to translate the verbs that are highlighted. However, we thought that it was not only necessary to look for their equivalents, but also to look for words that co-occur with them since the translators take into account the whole sentence in order to make their translation sound naturally.

The verbs we were going to analyze are "desarrollar", "administrar", "realizar", "presentar" and "someter". The problem lied in the fact that there were several options of translation. Furthermore, their equivalents had to be selected according to the words that co-occurred with them, as mentioned above. It was the reason why we considered that it was a good idea to extract some KWIC concordance lines, using as the search term the word associated with the verb.

First, we began with "desarrollar" and the search term was "phase". The concordance lines showed that "perform" could be used as the translation equivalent. On the other hand, we observed that "realizar" is a synonym of this verb in the Spanish text and therefore the translator may use the same translation equivalent.

As for "administrar", we made a guess and extracted KWIC concordance lines with the key term "administer" to validate this. As the image below shows, it is the right verb.

Concord	Concordance Plot File View Clusters Collocates Word List Keyword List							
Hit	KWIC							
1	kely to benefit from radiotherapy and to administer treatment in ways that max							
2	kely to benefit from radiotherapy and to administer treatment in ways that max							
	become the most common means by which to administer breast boost treatments, a							
4	gle molecular biomarker target, and then administer all these nanoparticles a							
5	s been criticized because of the need to administer a radioactive tracer. Bes							

Image 16: Concordance for the search term *administer*

The next verb we analyzed was "presentar" and the search term was "partial response". We observed that "achieve" was the appropriate verb since occurred together with the term.

Conce	ordance Concordance Plot File View Clusters Collocates Word List Keyword List
Hit	KWIC
9	sponse with lymphocytosis, the median time to true partial response or con
10	uration by best response CR=complete response. PR=partial response. AE=adv
10 11 12 13 14	ated IgHV might have been predisposed to achieve a partial response or con
12	ore rapid resolution of lymphocytosis. However, if partial response with :
13	nt during follow-up. Of 13 patients who achieved a partial response with :
14	achieved a complete response and seven achieved a partial response with a

Image 17: Concordance for the search term partial response

Finally, "someter" was also analyzed in the corpus. We extracted KWIC concordance lines with "surgery" as the search word and found that the translation equivalent was"undergo". Then, more concordance lines were extracted with "undergo" in order to observe its collocations as the image below shows.

Conco	ordance	Concordance Plot	File View	Clusters	Collocates	Word List	Keyword List
Hit	KWIC		5				
1	beet	n proposed that	epithelial	tumor cel	ls undergo	EMT, rel	easing mesenchymal
1 2	CSCs	. The ability of	epithelia	1/non-CSC	to undergo	EMT and	acquire CSC proper
3	hope	s of diminishing	a tumor's	capacity	to undergo	metastas	is. TGF-? signalin
	of	cancer cells wit	hin the bu	ilk of tume	ors undergo	differen	t degree of EMT, a
4 5	eath	domain of Bim i	nduced mit	ochondria	to undergo	fission,	by increasing Bcl
	ide	ntify a subgroup	of patien	ts who cou	ld undergo	breast-c	onserving surgery
6 7 8 9	ustr	ian trial that r	andomized	869 women	to undergo	breast-c	onserving surgery
8	ient	s with 1 to 3 in	volved lym	mph nodes w	tho undergo	mastecto	my should strongly
9	also	should be consi	dered for	patients w	tho undergo	breast-c	onserving surgery.
10	gest	ed that rather t	han recove	ery, patier	nts undergo	adaptati	on to chronic symp
11	use :	nanoparticles me	asuring le	ss than 5	nm undergo	rapid re	nal clearance whil
12	tion	, cancer often h	as necroti	c sites th	at undergo	calcific	ation. Identificat

Image 18: Concordance for the search term undergo

6. Conclusion

Some translation students could think that second language acquisition makes them able to translate a technical or scientific text without any problem. But as we have seen at the beginning of our project, specialized language differs from global language. Therefore, they should be aware of that difference and acquire knowledge about the rules of LSP to write a specialized text.

As we have already pointed out the translators face with a set of problems when they translate a specialized text and in most cases they do not how to reach a solution despite the help of resources (e.g. dictionaries, glossaries, databases, and so on). So, this fact led us to build a specialized corpus.

It is true that many years ago, it was more difficult to compile a corpus due to the fact that it was a time-consuming process and the users were required to spend many hours collecting information in libraries. However, nowadays the translator can compile an electronic corpus rather fast. This is due to the great amount of information that is available online.

As we already know, comparable and parallel corpora are used in translation studies. But they are not so easy to compile. It is difficult to find translations of a source text to build a parallel corpus; whereas what sometimes make comparable corpus complicated to build is to meet the same characteristics for the two sub-corpora. In this project we compiled a monolingual corpus since our aim was to check its usefulness for medical translation.

The design criteria of a corpus always depend on the purpose for which you compile it. In our case, the corpus was compiled to extract linguistic information about the terms used in the IMRAD superstructure of medical research articles in order to indicate the benefits that the translators can obtain from it in case of translation problems. Our corpus consists of 45 full texts which comes to a total of 231.989 tokens. As was stated above, the Internet allows us to collect information easily. Nevertheless, not all information is reliable; consequently,

we have to be cautious choosing it. This fact was the main reason why we used the database the University of Valladolid's library to compile our corpus so that it will yield good results. We chose the journals that were listed at the top of the impact factor rank in the field of Oncology.

Regarding the texts included in our corpus, we chose to include the four typical parts of medical research paper: Introduction, Materials and Methods, Results and Discussion (IMRAD superstructure). Therefore, we did not choose chunk texts, because some important terminological and phraseological information could be excluded from our analysis. The texts deal with the same subject field (Oncology) and belong to the same text type (the research article) because we were interested in analyzing expert-to-expert texts belonging to a specific field.

It is a synchronous corpus, since the papers are from 2010 to 2014; therefore it can be considered as a representative sample of current specialized language. They are original texts, not translations. We believe that translations may be a problem for our analysis because these are not always expressed in a natural way.

In the selection of the program software, we chose AntConc mainly due to two reasons: it does not require any payment since it is free online, and it offers several tools that are necessary in corpus analysis.

The first step in our corpus analysis was the term extraction. Instead of using online term extractors, we made a decision to use our own corpus as a term extractor on the basis of our experiences by using these term extractors in our Applied Linguistics classes (Pizarro Sánchez 2013-14). This way we avoided the problems of *noise* and *silence*. Before extracting term candidates, the *stop list* was used to delete those words that we considered that were not necessary to extract. We extracted single term candidates and multiword term candidates by using the Wordlist and Cluster tools based on the approach by Pearson (1998:25). But we also used a concordance tool to extract those term candidates that are

sometimes more difficult to detect. Once term candidates were extracted, four specialized glossaries related with oncology were used to indicate the key terms. Although it was a little long process, we finally obtained the term list, which helped to solve our terminological problems.

The second step in our analysis was the collocate extraction. We have focused on lexical units (nouns) since we were interested in studying the use of noun terms in specific contexts. We have learnt in our Specialized Translation classes that there are not hardly any dictionaries and glossaries that include medical phraseological information in English (Mendéz Cendón 2012-13). This was the reason why we were interested in analyzing the collocations of nouns. As such, our analysis was focused on the patterns "verb + noun", "noun+ verb", "adjective+ noun" and "noun + noun". Apart from collocations, we also considered that it was important to extract another type of information, such as synonyms or acronyms since these have a high frequency of occurrence in specialized texts.

Once the KWIC concordance lines were generated, we noticed that the information retrieved from these lines revealed us what verbs, adjectives or nouns co-occur with the terms. Although the extraction and analysis of collocations may be a long process, the result is totally satisfactory. Therefore, the phraseological data extracted from our corpus can become "a type of documentation or source", which is, very helpful for us or for those who are dedicated to translate scientific research articles related to the topic of oncology.

A Spanish article was used in our concluding analysis because we needed to make a comparison between the usage of a specialized monolingual corpus with typical resources, such as dictionaries and glossaries in order to point out the advantages and disadvantages that can arise in the translation assignment. The same way that we look up a word in a dictionary, glossary or thesaurus, we can also look it up in a corpus. However, the results are different according to information provided. It is true that the previous resources provide you with translation equivalents and some examples of usage, in case of bilingual dictionaries, but their information is often limited. On the other hand, a corpus gives you

translation equivalents, synonyms, acronyms, phraseology, examples in several contexts and the definition of the term in some cases.

Moreover, we think that it was essential to explain possible ways of finding an equivalent in the corpus. For example, when our task was to find an equivalent for a particular verb we used lexical units that co-occur with the search term. In this regard, we consider that the translators should be creative and resourceful when they extract concordance lines, especially for those who are not familiar with corpora since translation equivalents are sometimes difficult to detect.

In addition to what has been mentioned previously, the following points illustrate the importance of building specialized monolingual corpora.

Firstly, translators specialized in scientific or technical fields usually have problems to find equivalents in dictionaries or glossaries because these do not contain information or they are outdated, whereas an up-to-date corpus can be used as term extractor, and the terminological problems will be solved easily.

Secondly, electronic corpora are very useful because you can detect and extract real and authentic examples of language, which allows you to write a translation that sounds natural, like an original text.

Thirdly, dictionaries and glossaries offer several possibilities of translation and sometimes the translator does not know what the appropriate equivalent is in the target language. On the other hand a corpus helps the specialized translator to select always the right translation equivalent. Then, there is the question of choosing the appropriate equivalent in your corpus in case that you decide to follow your intuition and you are not sure whether it is right or not.

Finally, the usage of a specialized monolingual corpus makes the translators to understand and acquire more knowledge about a particular subject field (in our case is Oncology). As we have said, the task of the translators involves several challenges since they play the role of mediators between two languages (the source language and the target language) and cultures. Through this task, they may be faced with terminological and phraseological problems, which can affect the quality of their translation. Undoubtedly, the usage of corpora can ease their translation task because they are allowed to have access to a great amount of data. The information extracted from concordance lines helps them to solve their doubts about the meanings of words and, what is more important, to learn how these words combine with others within a sentence.

To conclude, we would like to highlight the usefulness of specialized monolingual corpora as a translation tool, by saying that they are a requisite for the translation activity. They play a significant role in the field of translation and therefore we consider that they should be integrated in the translation classes as part of the teaching curriculum.

7. References

Atkins, S. et al. (1992) "Corpus Design Criteria," *Literary and Linguistic Computing*, 7: 1-16.

Baker, M. (1993) "Corpus Linguistics and Translation Studies: Implications and Applications". In Baker, Mona/ Francis, Gill/ Tognini-Bonelli, Elena (eds.) *Text and Technology: In honor of John Sinclair*. Amsterdam/Philadelphia: John Benjamins, pp. 233-250.

Baker, M. (1995) "Corpora in translation studies: An Overview and Some Suggestions for Future Research," *Target* 7, pp. 223-243.

Baker, M. (1996) "Corpus-based Translation Studies: The Challenge the Lie Ahead," In Somers, Harold (eds.) *Terminology, LSP and Translation: Studies in Language Engineering in Honour of Juan C. Sager.* Amsterdam/Philadelphia: John Benjamins Publishing Company, pp. 175 -186.

Bernardini, S. (2004) "Corpora in the Classroom: An Overview and Some Reflections on Future Developments." In John Sinclair (ed.) *How to use corpora in language teaching*. Amsterdam: John Benjamins, pp. 16-36.

Biber, D. (1993) "Representativeness in Corpus Design". *Literary and Linguistic Computing*, 8, pp. 243-257.

Bowker, L., and J. Pearson (2002). *Working with Specialized Language: A Practical Guide to Using Corpora*. London: Routledge.

Bravo Gozalo, J.M. and P. Fernández Nistal (1998). "La lingüística del corpus, las nuevas tecnologías de la información y los Estudios de Traducción en la década de 1990". In

Fernández Nistal, P. and J.M. Bravo Gozalo (eds.). *La traducción: Orientaciones Lingüísticas y Culturales*. Valladolid: Universidad de Valladolid, pp. 205-257.

De Santiago, Paula. "Introducción a la Traducción Automática." Class notes 2013-2014. Universidad de Valladolid.

Kübler, N. (2003) "Corpora and LSP translation". In Zanettin, Federico/ Bernardini, Silvia/ Stewart Dominic (eds.) *Corpora in Translator Education*. Manchester: St. Jerome Publishing, pp. 101-112.

Maia, B. (2005) "Terminology and Translation – Bringing Research and Professional Training together through Technology," *Meta: Translator's Journal*, 50. Èrudite. 22 Feb. 2014 <<u>http://nelson.cen.umontreal.ca/revue/meta/2005/v50/n4/019921ar.pdf</u>>

Mendéz Cendón, B. (2002) *Estrategias fraseológicas en el género discursivo de los artículos médicos en lengua inglesa*. Tesis Doctoral. Alicante: Biblioteca Virtual Miguel de Cervantes.

Mendéz Cendón, Beatriz. "Traducción de Lenguajes Especializados II (Inglés/Español: Ciencia y Técnica)." Class notes 2012-2013. Universidad de Valladolid.

Meyer, I. and K. Mackintosh (2002) "When terms move into our everyday lives: An overview of de-terminologization." *International Journal of Theoretical and Applied Issues in Specialized Communication*, 6. John Benjamins Publishing Company, pp. 111-138.

Olohan, M. (2004) Introducing Corpora in Translation Studies. London/New York: Routledge.

Pearson, J. (1998) *Terms in Context*. Amsterdam/Philadelphia: John Benjamins Publishing Company.

Pizarro, Isabel. "Lengua Inglesa: Lingüística Aplicada III." Class notes 2013-2014. Universidad de Valladolid.

Sinclair, J. (1987) Looking Up, London and Glasgow: Collins.

Sinclair, J. (1991) Corpus Concordance Collocation, Oxford: Oxford University Press.

Tarp, S. (2007) "¿Qué requisitos debe cumplir un diccionario de Traducción del siglo XXI?" In Fuertes-Olivera, P.A. (ed.) *Problemas Lingüísticos en la Traducción Especializada*. Valladolid. Universidad de Valladolid, pp. 227-256.