# Phraseological approach to automatic terminology extraction from a bilingual aligned scientific corpus

Frérot Cécile[1*], Rigou Géraldine[2], Lacombe Annik[3]

[1] University of Paris 7-Denis Diderot, 2 place Jussieu, F-75251 Paris Cedex 05, France, frerot@cicrp.jussieu.fr

[2] INRA-CRJ, Unité Centrale de Documentation, Translation and Terminology Unit, F-78352 Jouy-en-Josas Cedex, France, rigou@jouy.inra.fr, tel: +33 (0) 1.34.65.24.54, fax: +33 (0) 1.34.65.22.72

[3] INRA-CRJ, Unité Centrale de Documentation, Translation and Terminology Unit, F-78352 Jouy-en-Josas Cedex, France, alacombe@jouy.inra.fr, tel: +33 (0) 1.34.65.24.55

* This work was carried out during a training period at the INRA Translation and Terminology Unit as part of a post-graduate degree (DESS Industrie de la Langue et Traduction Spécialisée), UFR Etudes Interculturelles de Langues Appliquées

**Abstract**: The linguistic knowledge represented in specialised dictionaries should not be restricted to a collection of terms (single and complex terms). It should include phraseological units, i.e. more or less fixed multiword expressions –often called collocations– which cause serious problems for translators and technical writers, since translating them into the target language is rendered difficult by the syntactic and lexical characteristics of each Language for Specific Purposes (LSP). Corpus-based terminology acquisition tools make it easier to identify and collect these units and to generate terminology resources that better meet the needs of users.

We present the results of an automatic terminology extraction from a French-English aligned corpus with the aim of developing a scientific translation and writing tool for French researchers with emphasis on the phraseological dimension. This experiment was conducted at the Translation and Terminology Unit of the French National Institute for Agricultural Research (INRA), in collaboration with two researchers working in automatic terminology extraction and bilingual terminology alignment. LEXTER extracted term candidates from a French-language scientific corpus and TRINITY identified the translation candidates in the English corpus made up of the translations of the French texts. The translator-terminologists exploited the results of the extraction using the hypertext interface for validation of a database.

After describing the different stages of the experiment, from preparing the corpus to data processing, we explain how we validated and exploited the results. We focus on collocations and the problems linked to their identification and terminographic description from a translation perspective and consider the problems inherent to phraseology in LSP when efforts are made to improve natural language processing (NLP) tools.

**Keywords**: terminology extraction, aligned corpus, phraseology, LSP, scientific translation.

## 1   Introduction

Despite considerable terminographic work carried out to collect the vocabulary of specialised subject fields, deficiencies still exist in emerging and/or evolving subject fields. Available terminological resources do not fully meet the needs of technical writers and translators. Resources are often limited to terminological units, whether they are simple or complex terms, and exclude phraseological units. However, a large number of problems are raised by the linguistic environment of terms, i.e. the syntagmatic extension of the term to the sentence (Blais, 1993). The ability to deal with this linguistic environment is absolutely necessary for language specialists eager to respect the type of language used by the experts. Consequently, the study of lexical units other than nouns (e.g. verbs, adverbs and adjectives) is of utmost importance. Nevertheless, in spite of a growing interest in problems linked to identifying and collecting phraseology in LSP (Language for Special Purposes), there are still few electronic or print resources available for translators.

Corpus-based terminology acquisition tools such as terminology extractors make it easier to generate terminological resources. These tools help identify simple and complex terminological units as well as expressions –often known as collocations– specific to a linguistic community and including other lexical units such as verbs, adjectives and adverbs. Since we placed the emphasis on usage linked to collocations, we took into account the definition given by Pesant and Thibault (1993) who define collocations as being different parts of speech which appear together in a text and combine to form expressions fixed by usage. Two examples of collocations extracted from our corpus are *technologically valuable lactic acid bacteria* and *radio-tagged river trout.*

The present work describes the results of an automatic terminology extraction from a bilingual aligned scientific corpus. We aimed both at evaluating how well these tools performed and at analysing how they could be used to develop a scientific writing tool meeting the needs of French researchers at the National Institute for Agricultural Research (INRA).

## 2   Terminology extraction

### 2.1   Context of the experiment

The experiment was conducted by the Translation and Terminology Unit in collaboration with Didier Bourigault[1] who developed LEXTER, a tool for terminology extraction, and David Hull[2] who developed TRINITY, a tool for bilingual word alignment. The experiment aimed at increasing the terminological database currently used by the translators of the Unit. In the long term, we intend to develop a computer-based English writing tool for the French researchers of the Institute. Nowadays, the vast majority of publications are written in English. It is of vital importance to publish one's results, especially in emerging subject fields where competition is fierce, and the level of English is one of the criteria for acceptance of papers. However, whereas researchers often have a good command of their field's terminology, they often encounter great difficulty in expressing their ideas clearly and, as a result, in writing a syntactically coherent text. At INRA, they have the support of in-house translators to write their papers. In order to compensate for the lack of terminological data, the translators search on-line bibliographical databases constituting immense corpora, from which they extract the terminology and phraseology they cannot find in dictionaries.

The experiment aimed to collect terminology from a French-English translation corpus by validating the results of an automatic extraction. The analysis of noun phrases (NP) extracted by LEXTER helped us identify the terms and the extensions of these NPs. We then turned to the French list of nouns and verbs to collect verb phrases belonging mainly to general language and frequently used by researchers.

### 2.2   Preparation of the corpus

This extraction was carried out from an aligned bilingual corpus made up of French to English translations. The corpus represented a total of 340,000 words and included, in decreasing order: research papers, press releases and publications aimed at the general public, a software user's guide, presentation leaflets, a licence contract and summaries of monographs. The corpus comprised texts from different subject fields of which the most represented were: agricultural sciences, soil sciences, hydrobiology, environment, biometry and modelling, plant breeding and genetics, plant disease and

weed science. Paragraph marks in both the French and English texts were used to manually align the corpus. We also removed from the corpus any element likely to generate noise during the extraction process (bibliographical references, mathematic symbols, figures, etc.).

## 2.3 Data processing

Xerox tools were used to align the corpus sentence by sentence and tag it.

### 2.3.1 LEXTER[3]

LEXTER extracted term candidates (TC) from the tagged French corpus. TCs designate words and multiwords likely to be terminological units or collocations. The first stage in the process is a morphological analysis of the texts followed by the identification of noun phrases of maximal length, using linguistic rules of boundary detection. In the second stage, LEXTER uses a range of parsing rules to extract from each noun phrase of maximal length a set of sub-phrases which are likely to constitute terms by virtue of their position in the noun phrase (Bourigault, 1994). After statistical filtering, LEXTER supplies a list of TCs.

### 2.3.2 TRINITY[4]

TRINITY is an alignment system using statistical word alignment techniques to automatically construct a bilingual word and phrase lexicon from a collection of translated sentences. Statistical word alignment algorithms use common regular co-occurrence patterns between source and target language words to establish links between occurrences of these words in individual sentence pairs (Hull, to be published).

### 2.3.3 RESULTS: HYPERTEXT INTERFACE FOR VALIDATION

LEXTER automatically generates a list of term candidates corresponding to different grammatical categories (noun, adverb, adjective, verb) as well as adjective phrases and noun phrases. The LEXTER hypertext interface for validation gives access to a list of term candidates either in decreasing frequency or alphabetical order, depending on the lexical units mentioned above (see Appendix A). It is also possible to view the contexts in which the term candidates appear. Thanks to the output of the TRINITY system, the term candidates and the translation candidates are presented in their aligned contexts (see Appendix B). The number of occurrences is available for each item of data. For each noun phrase, the hypertext interface gives access to its extensions, i.e. the term candidates in which the noun phrase is the syntactic head or expansion. The interface also makes it possible to use a validation scale which we determined on the basis of French and English word segmentation as well as terminology and translation relevance. Lastly, French term candidates and their English equivalents can be entered into a lexicon if they are considered relevant.

## 3 Analysis of the results

### 3.1 Noun phrase analysis

#### 3.1.1 CREATION OF TERMINOLOGICAL ENTRIES

The structure of the lexicon supplied by the hypertext interface made it possible to enter only French term candidates and their English equivalents together with comments. Since we aimed to collect additional information useful to both translators and researchers, we created a terminological entry that would best meet their needs and based its structure on the existing entry of the translators' database. Therefore, we kept some of the fields (namely: *Term, Variant* –several, *Context, Subject field, Descriptor, Linguistic note)* and added two collocation fields (*Noun collocation* and *Verb collocation*). We removed the *Definition* field (see Appendix C). The aim was not to formulate terminological definitions but to collect defining statements in order to understand the terms. However, it should be noted that identifying defining statements proved quite unproductive. The reason for this unproductivity lies in the fact that the main type of publications contained in the corpus was the research paper: research papers relate to the expert-expert communicative setting and concepts are rarely defined since experts are assumed to have the same or very similar level of expertise.

---

[3] For further details, see BOURIGAULT D 1994
[4] For further details, see HULL D (to be published)

### 3.1.2 LEVELS OF ANALYSIS

We selected noun phrases according to three approaches: terminology, translation and scientific writing. The specific needs of translators and researchers were also considered.

Translators must use the appropriate terminology as well as the linguistic expressions specific to a subject field which reflect the language as used by researchers. This is rendered difficult by the diversity of INRA's subject fields. The selection of noun phrases is influenced by the linguistic skills and scientific expertise of the translator. This usually depends on the translator's experience (namely beginner, experienced, sub-contractor). Since translators need to understand the concept referred to by the author, we collected defining statements whenever possible.

As regards researchers, they are more concerned with phraseology than terminology. Two main reasons explain why they often have great difficulty in clearly expressing their ideas in English: their command of the target language can be quite limited and, secondly, they tend to apply French syntax, hence the likelihood of loan translations. Their knowledge of LSP phraseology is often very poor even when it comes to general language expressions such as *to put forward hypotheses, to conduct research, to take into account,* hence the focus on collecting verb phrases belonging to LGP (see 3.2). Generally speaking, researchers encounter difficulty when combining specialised terms with words from general language. From a theoretical point of view, this raises the issue of "the separation of terms on the one hand and the linking elements from LGP on the other [which] may not be maintained uncritically" (Picht, 1987).

After identifying these needs, we selected noun phrases on the basis of the issues raised by:

- translators: are these noun phrases or collocations specialised enough to be selected? Is this collocation an original feature linked to the author's style or a rather set expression worth being collected?

- terminologists: does this noun phrase designate a specific concept?

- researchers: which words are used to describe the steps of an experiment? Which is the appropriate verb used to describe results, a comparison or cause and effect? Which expressions would a native English speaker use? Is there a specific expression used to convey a given idea? What is the English equivalent of that specific French verb?

We analysed each noun phrase as follows: the term candidates and the translation candidates were analysed, as well as head and expansion productivity (see Appendix D), and the French and English contexts were read through to collect defining statements and collocations LEXTER would not have extracted.

### 3.1.3 IDENTIFICATION AND RECORDING OF COLLOCATIONS

The analysis of collocations has pointed out the lexical and syntactic differences between the two languages, as well as their unpredictability (Heid and Freibott, 1991) when translating from French into English. We bore these linguistic facts in mind when selecting collocations, which we define as usual associations of several words linked by prepositions and referring to one or several notions. We focused on prepositions since translating the phrases they are part of is a major problem for researchers. This definition excludes complex terms with patterns such as *noun + noun* or *adjective + noun.*

The main difficulty consists in differentiating a term from a collocation when analysing a sequence of lexical units. The difference involves a change from a conceptual organisation to a linguistic environment. Consequently, collocations are subjected to greater variation than terms, which is also related to the author's own style, to the creative power of phraseology (Blampain, 1993). It is difficult to determine whether experts in a given subject field often use a collocation or whether it is pure stylistic "extravagance". Frequency was not considered to be a criterion of selection, as most of the corpus belongs to a very specialised discourse, the terminology and phraseology of which are not widespread (Béjoint and Thoiron, 1992). Moreover, the corpus is of limited size and is therefore not representative of scientific discourse.

To record collocations, it is necessary to consider the intellectual processes involved in translating and writing. Translators and researchers probably try to answer the following question: What words or verbs usually combine to describe a given process, to show results? (Cohen, 1992) rather than What words or verbs usually combine? Consequently, concept is crucial and provides some answers as to whether users should have to access one collocation referring to several notions together or each of these notions separately.

In the following collocation: *addition de lipides à la ration de la truie en lactation*, the various notions referred to mean that three different entries need to be created (*addition de lipides / ration / truie en lactation*), but the unpredictability of the collocation led us to enter *addition de lipides à la ration de la truie en lactation* into the *Noun collocation* field of the three terminological entries. As it is more relevant and reliable to search a term rather than a collocation when querying a database, a collocation will be entered under the keyword (the base of the collocation) as well as the co-occurrent.

Example: Noun collocation: *dry matter content* under the entries *content* and *dry matter*.

To record all the collocations related to a term in a single field would make the content of the entry difficult for users to read and access to the information would be delayed. Therefore, we decided to create two collocation fields to improve the way in which noun collocations and verb collocations are identified.

### 3.2   Verb analysis

The extraction was performed from the French corpus and concerned simple verbs. TRINITY in this case did not align the translation candidates in English, which involved us:

- reading the French contexts to analyse the linguistic environment of the verb and identify complex verb phrases,

- reading the English contexts to identify the equivalents.

#### 3.2.1   CREATION OF VERB ENTRIES

The structure of the verb entry differs from that of the NP terminological entry. We integrated the following fields: *French verb*, *English equivalent*, *Linguistic note* but removed the *Definition*, *Subject field* and *Descriptor* fields. In fact, we focused on the information relating to the ways in which these verbs function and are used in context. We created a dozen fields to record linguistic combinations in which the verb appears (see Appendix E). These combinations are more or less fixed and constitute sentence segments that are interesting from a translation point of view.

In order to facilitate access to information, we created an entry for each verb phrase containing a support verb such as *mettre* which produces approximately 20 complex verb structures.

#### 3.2.2   SELECTION OF VERBS

From the list of "verb candidates", we selected those frequently used in original research papers due to their predominance among the documents translated or reviewed by the Translation Unit. An original research paper displays an IMRaD structure composed of the following sections:

**I**ntroduction
**M**aterials and methods (experimental protocol)
**R**esults
(**a**nd)
**D**iscussion

We first excluded specialised verbs referring to a specific notion in a subject field for several reasons: these verbs only create minor difficulties in translation, there is very little morphological or syntactic variation between both languages and researchers usually know the English equivalents of specialised verbs, as they do for specialised terms. This hypothesis was confirmed by analysing the list of specialised French verbs, which was relatively limited, and their English equivalents, as shown by the following examples: *cloner (to clone), coder (to code).*

We selected verbs belonging to general language which create difficulties for researchers when translating them, as well as support verbs such as *mettre* which produce approximately 20 complex verb structures (*mettre en comparaison, mettre en évidence…*). Here are some of the verbs we analysed: *montrer, constituer, réaliser, conduire, mener, estimer, permettre, varier, sembler, apparaître, entraîner, rechercher.*

#### 3.2.3   TERMINOGRAPHIC DESCRIPTION

The verb entry should not be limited to a lexicon, i.e. a verb and its equivalent(s) in the target language (L'Homme, 1993). Problems in translating complex verbs are related to the syntactic structures specific to French and English. Users therefore need more information than is usually

provided for a simple or complex term.  It is necessary to record how the verb functions from a linguistic point of view and to:

- describe how the verb functions syntactically (What kind of complement can be used alongside this verb?  What preposition is required with this verb?  Is it a transitive or intransitive verb?),

- provide semantic information by proposing a synonym to specify the meaning of the French verb (the French synonym is usually the most similar morphologically to the English equivalent),

- give an example of the verb in context, as it occurs in a sentence.

Interestingly, the English equivalents of the complex verb structures having a support verb as their nucleus, such as *mettre*, are often simple verbs (*mettre en évidence: to demonstrate*), and the French verb has several equivalents in English (*entraîner: to cause, to lead to, to result in*).

### 3.2.4 PROSPECTS FOR VERB ANALYSIS

The analysis of NPs had evidenced errors in the extraction performed by LEXTER.  For example, the NP list consists of a large number of sequences for which *compte de* is a head (*[tenir] compte de l'accumulation préférentielle d'assimilats, [rendre] compte de l'orientation des particules, [prendre] en compte des contraintes techniques*).  These sequences, had they been well identified, would have allowed verb phrases such as *rendre compte de, tenir compte de, prendre en compte* to be extracted. Therefore, from the list of nouns extracted by LEXTER, we selected those presumably having high co-occurrence productivity with verbs.

Examples:       *hypothèse:* admettre une ~, émettre une ~, faire l'~ que, rejeter une ~, etc.

*résultats:* commenter des ~, discuter des ~, obtenir des ~, diffuser des ~, etc.

*gènes*: identifier des ~, localiser des ~, introgresser des ~, porter des ~, etc.

Of the nouns with high co-occurrence productivity with verbs, it appears that some belong to general language and are combined with verbs from general language (*admettre une hypothèse*), while some are terms (*gène*) and are combined either with verbs from general language (*identifier des gènes*) or with specialised verbs (*introgresser des gènes*).  These examples illustrate the fact that general language intermingles with LSP in scientific documents and raises the issue of the supposed dichotomy between LGP and LSP, and the widely acknowledged point of view which considers as a boundary the limit between highly specialised jargon and words from general language (Bourigault and Jacquemin, 2000).

On the basis of the list of nouns extracted by LEXTER, we will further analyse complex verb structures starting with nouns from general language, since this type of structure raises major problems for INRA researchers when writing their papers.  We will continue with verb structures including terms.

Recording these different types of verb structures raises the problem of how they can be accessed. For the time being, the structures we have analysed do not allow us to define the best approach to recording them: it is still unclear whether access should be given to the whole verb structure or to the noun and verb separately depending on what information we will provide.  According to the information available in the corpus, we could provide two kinds of entries: a specific term with the various verbs it can be combined with, or a verb and the kinds of complements with which it can be used (L'Homme, 1997).  This will probably need to be discussed with future users as researchers and translators may not go about searching databases in the same way.


## 4  Conclusion and prospects

In the light of the corpus size and heterogeneity, the quality of the extraction, in French as well as in English, is irrefutable.  Terminology extractors are useful to translators considering the very little time they have to enter terms into a terminological database.  Moreover, this kind of extraction allows experienced translators to record NPs they would not have entered otherwise, judging them as being too basic. As for verbs, using LEXTER enables verb phrases to be identified rapidly within a context.  It would be interesting to perform a similar automatic extraction with English verbs.

We will continue exploiting the results of the extraction and, when the database has a sufficient number of entries, we will carry out a test in the field with researchers so as to make sure it meets their needs.

The work carried out on verbs could be improved to better describe verbs in scientific discourse which, in turn, would make it possible to determine appropriate conceptual classes of co-occurrents and classify verbs according to their usage. A more detailed analysis of verbs in LSP could contribute to improving corpus-based terminology acquisition tools and, more generally, to improving NLP.

Once the users' needs are taken into consideration and the analysis is no longer restricted to noun terms, other parts of speech (verb, adverb, adjective) emerge: it is then necessary to adopt a phraseological approach.

## 5 References

Bejoint H, Thoiron P 1992 Macrostructure et microstructure dans un dictionnaire de collocations en langue de spécialité. *Terminologie et traduction* 2/3: 513-522.

Blais E 1993 Le phraséologisme. Une hypothèse de travail. *Terminologies nouvelles* 10: 50-56.

Blampain D 1993 Notions et phraséologie. Une nouvelle alliance ? *Terminologies nouvelles* 10: 43-49.

Bourigault D 1994 *Lexter, un logiciel d'extraction de terminologie. Application à l'acquisition des connaissances à partir de textes*. Thèse en informatique linguistique, Ecole des Hautes Etudes en Sciences Sociales, Paris.

Bourigault D, Jacquemin C 2000 Construction de ressources terminologiques. In J-M. Pierrel, *Ingénierie des langues*. Paris, Hermès Sciences Publications, pp 215-233.

Cohen B 1992 Méthodes de repérage et de classement des cooccurrents lexicaux. *Terminologie et traduction* 2/3: 505-511.

Heid U, Freibott G 1991 Collocations dans une base de données terminologique et lexicale. *Meta* 36(1): 77-91.

Hull D (unpublished) Software tools to support the construction of bilingual terminology lexicons. In Bourigault D, Jacquemin C, L'Homme M-C (eds). *Recent advances in Computational Terminology*. London, John Benjamins.

L'homme M-C 1993 Le verbe en terminologie : du concept au contexte. *L'actualité terminologique* 26(2) 17-19.

L'homme M-C 1997 Méthode d'accès informatisé aux combinaisons lexicales en langue technique. *Meta* 42(1): 15-23.

Pesant G, Thibault E 1993 Terminologie et cooccurrence dans la langue du droit. *Terminologies nouvelles* 10: 23-35.

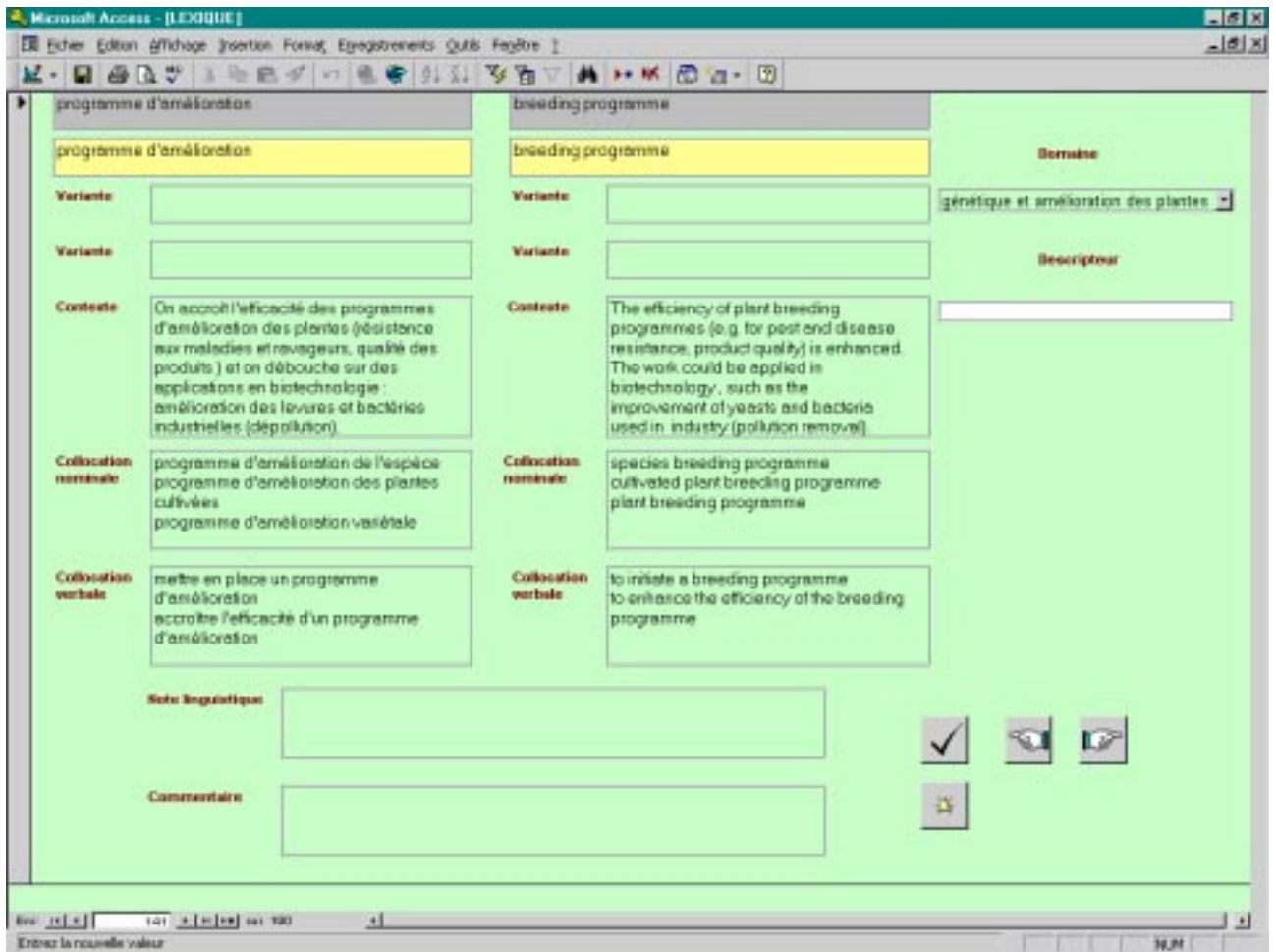Picht H 1987 Terms and their LSP environment – LSP phraseology. *Meta* 32(2) 149-155.

APPENDIX A: Hypertext interface for validation: TCs (NPS) by decreasing frequency

| freq tot | freq isol | freq non isol | prod tot | prod T | prod E | num Fam. | nb TCP | nb Corp | cat | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 78 | 17 | 61 | 39 | 18 | 21 | 0 | 4 | | SN | ⊙⊙⊙⊙⊙⊙ | matière organique |
| 78 | 46 | 32 | 21 | 8 | 13 | 0 | 6 | | SN | ⊙⊙⊙⊙⊙⊙ | facteur d'impact |
| 49 | 26 | 23 | 19 | 11 | 8 | 20 | 7 | | SN | ⊙⊙⊙⊙⊙⊙ | teneur en eau |
| 36 | 30 | 6 | 6 | 1 | 5 | 0 | 4 | | SN | ⊙⊙⊙⊙⊙⊙ | grande Alose |
| 33 | 9 | 24 | 16 | 3 | 13 | 0 | 7 | | SN | ⊙⊙⊙⊙⊙⊙ | bactérie lactique |
| 31 | 17 | 14 | 12 | 5 | 7 | 203 | 10 | | SN | ⊙⊙⊙⊙⊙⊙ | système de culture |
| 29 | 10 | 19 | 6 | 3 | 3 | 0 | 8 | | SN | ⊙⊙⊙⊙⊙⊙ | discipline scientifique |
| 29 | 15 | 14 | 8 | 1 | 7 | 0 | 7 | | SN | ⊙⊙⊙⊙⊙⊙ | cuivre extractible |

APPENDIX B: TCs and translation candidates in aligned contexts

matière organique — organic matter
matière organique — organic matter

Fermer

| 31-6-3-3_p7-4 | La quantité présente sous la forme 1 dépend essentiellement de la composition minéralogique du sol et de sa teneur en matière organique . | The amount present in form 1 depends mainly on the mineralogical composition of the soil and on its organic matter content . |
|---|---|---|
| 31-8-1_p2-1 | , la minéralisation de l'humus , qui résulte de la décomposition de la matière organique humifiée du sol . | , humus mineralisation , due to the decomposition of humified organic matter in the soil . |

APPENDIX C: Terminological entry



APPENDIX D: Head and expansion productivity

APPENDIX E: Verb entry