# HOW SPECIALIZED (OR POPULARIZED)?
# Terminological density as a clue to text specialization in the domain of food safety

ADRIANO FERRARESI
UNIVERSITY OF BOLOGNA

**Abstract** – Terminological density is often mentioned as one of the main lexical features that distinguish specialized forms of communication and their popularized counterpart, i.e. discourse conveying specialized knowledge to an audience of non-specialist readers. The underlying idea is that texts addressed to experts feature a high number of terms that are expected to be understood within the target discourse community, while in popularized texts, vocabulary is chiefly drawn from the general language, and hence density of terms is lower. Though intuitively appealing, the idea that terminological density can provide quantitative evidence as to the degree of specialization or popularization of texts has seldom been investigated empirically. Taking the domain of food safety as a case in point, the aim of this article is twofold. First, it aims to assess the extent to which terminological density reflects the distinction between more specialized and more popularized texts. To do so, it describes a method to operationalize and measure terminological density building on replicable corpus-based procedures and freely available tools for vocabulary profiling. Second, in a more descriptive perspective, it aims to relate quantitative findings to qualitative observations on the discursive strategies adopted in the popularized genres under consideration to target different audiences. The texts analyzed consist of scientific opinions addressed to experts published by the European Food Safety Authority, as well as their popularized versions (factsheets, Frequently Asked Questions and news), also produced by the Authority to disseminate knowledge to the wider public. Results provide evidence that terminological density as operationalized here not only reflects the difference between specialized and popularized texts, but can also point to more subtle differences related to popularized texts' varied audiences and discursive strategies.

**Keywords**: specialized corpora; popularization; terminological density; vocabulary profiling.

## 1. Introduction

The notion that scientific or technical texts are characterized by different degrees of specialization is a central one in several applied linguistics fields, from English for Special Purposes (ESP) to vocabulary research. Text specialization (or "technicality", see e.g. Krüger 2015) is closely linked to the

level of expertise in a subject matter displayed by text producers and expected in target audiences: communication between experts in a scientific or technical domain typically features a higher degree of specialization than communication between experts and non-experts. Hence the relevance of this parameter, e.g., in assessing the suitability of texts for teaching disciplinary contents (Chung, Nation 2003), or their intelligibility by ESP students (Ha, Hyland 2017).

Focusing specifically on the strategies for conveying specialized knowledge to non-experts, generally subsumed under the umbrella term of popularization strategies (Calsamiglia, van Dijck 2004), several authors have suggested that the very distinction between specialized and non-specialized texts is not a binary one, but rather forms a continuum (e.g. Myers 2003, Delavigne 2003). In this view, popularized texts occupy intermediate positions along the continuum, and may gravitate toward one pole or the other depending on the audience that they target. This can include a wide spectrum of subjects, from a public "with little or no specialized knowledge" to "highly informed laypersons who […] may already approach semi-expert status" (Krüger 2015, p. 74), and even to experts "as soon as they step outside their very limited specialism" (Myers 2003, p. 268). In the words of Hilgartner (1990, in Myers 2003, p. 270), popularization, just like specialization, is "a matter of degree".

A large body of work has been devoted to shedding light on the diverse linguistic/discursive features that characterize popularized texts, especially by comparing them to well-established genres of scientific and technical communication, such as research articles (Gotti 2014, p. 17). Such features include greater reliance on figurative language (e.g. Constantinou 2014) and on definitional/paraphrasing structures (e.g. Gotti 2011), as well as the use of less nuanced markers of modality to express certainty and probability (e.g. Bondi 2015).

The present contribution focuses on one such feature, which is often mentioned in the literature as distinguishing highly specialized forms of communication from their popularized counterparts, i.e. terminological density (see Section 2). The underlying idea is that texts addressed to experts display a high density of domain-specific terms which are expected to be understood by the target audience, while in popularized texts vocabulary is likely to be chiefly "drawn from the general language" (Gotti 2011, p. 207), which makes density of terms lower. Though intuitively appealing, the idea that terminological density can provide quantitative evidence as to the degree of specialization or popularization of texts has seldom been investigated systematically.

Taking the domain of food safety as a case in point, the aim of this article is twofold. First, it aims to assess the extent to which terminological density reflects the distinction between more specialized and more popularized texts.

To do so, it describes a method to operationalize and measure terminological density building on replicable corpus-based procedures and freely available tools for vocabulary profiling (Nation, Heatley 2002, Anthony 2014). Second, in a more descriptive perspective, it aims to relate quantitative findings to qualitative observations on the discursive strategies adopted in the popularized genres under consideration to target different audiences. The texts analyzed consist of scientific opinions addressed to experts published by the European Food Safety Authority (EFSA),[1] as well as their popularized versions (factsheets, Frequently Asked Questions and news), also produced by EFSA to disseminate knowledge to the wider public.

The remainder of the article is structured as follows. Section 2 reviews previous work focusing on terminological density as an indicator of text specialization, and expands on the theoretical and methodological underpinnings of the notion of term adopted in the study. Section 3 presents the study methodology, and specifically the steps taken to build the corpora of texts in the domain of food safety, the method employed to extract and classify terms, and the setup of the vocabulary profiling tools used in the analysis of EFSA texts. Section 4 presents the results of the quantitative and qualitative analysis, and Section 5 concludes by discussing implications of the findings for research on broadly conceived popularization discourse.

## 2. Background

### *2.1. Text specialization and terminological density*

The degree to which texts make use of domain-specific terminology is arguably among the most frequently adopted parameters to assess text specialization. Both Pearson (1998, p. 62) and Corpas Pastor and Seghiri (2010) use it to select appropriate scientific and technical texts for inclusion in specialized corpora, on the grounds that "[d]ocuments with a high level of technicality [display] 'terminological density', or in other words, a high number of units that convey specialized knowledge" (Corpas Pastor, Seghiri 2010, p. 124). Along similar lines, Cabré (2010, p. 358) suggests that density of terms can be used as a parameter to classify specialized texts, arguing that "terminological density of texts varies according to their level of specialization: [the] more specialized the text, the more terminology it will have". These studies, however, do not report on quantitative measures of term density, nor do they mention explicit criteria through which the notion could be operationalized.

Andersen (1996) and Skibitska (2015) present more quantitatively-oriented approaches to assess terminological density. The former focuses on

---

[1] http://www.efsa.europa.eu

operating manuals related to oil-drilling procedures distinguished by different degrees of technical complexity, while the latter compares tourist texts addressed to professionals of the tourism industry and to tourists. Both studies measure term density as the proportion of domain-specific terms to the overall number of words in texts, and find support for the hypothesis that more specialized texts, i.e. texts concerning more technical procedures in the case of Andersen (1996) or addressed to experts in the case of Skibitska (2015), are associated with a higher density of terms compared to less specialized ones. Though more systematic than the above-mentioned studies, neither Andersen nor Skibitska provide details as to the methodology adopted for term identification, raising issues as to the replicability of findings.

Perhaps the most explicit attempt at establishing a correlation between terminological density and text specialization is Piscopiello and Bertaccini's (2009) study comparing medical texts at different levels of specialization/popularization (texts addressed to experts, to medicine students and to the general public) in two languages (English and Italian). Their operationalization of terminological density corresponds to the type/token ratio measure (TTR), i.e. the proportion between the number of different words (types) and the total number of words (tokens) in a text. Results indicate that TTR only reflects a coarse-grained distinction between the most specialized and most popularized texts, and does so more reliably in Italian than in English. Going beyond the partial inconsistency of results, one wonders whether a measure like TTR, which is entirely agnostic to whether a word constitutes a term or not, can appropriately be defined as a measure of terminological density.

The method described in this study (Section 3) aims at overcoming these limitations by proposing replicable, quantitative procedures for term identification, based on present-day terminology theory and practice. This makes the object of the next Section.

## 2.2. Term specialization and corpus-based terminology research

The advent of descriptive approaches to terminology was a major breakthrough in terminology research (Cabré 2003). Contrary to the traditional view of prescriptive terminology, descriptive terminology holds that specialized units of knowledge (concepts) can be expressed by more than one linguistic sign (term), and that it is essential for terminology as a discipline to account for such variation. By observing the use of terms in specialized contexts, it has been shown that term variation is pervasive, and that it is motivated by cognitive, linguistic and situational factors (Freixa 2006).

The level of subject-field expertise of text producers and addressees features prominently among these factors. According to Delavigne (2003) popularized texts addressed to a lay audience may not only feature few

specialized terms (see also Section 2.1), but terms themselves may display a low degree of specialization: in her example, "rougeur de la peau" (*skin redness*) and its specialized counterpart "érythème cutané" (*skin erythema*) should be treated as equivalent terms distinguished by a different level of specialization, and hence of intelligibility by the audience (Delavigne 2003, p. 89). The point is taken up by Freixa (2006, p. 59), who states that:

> [o]ne of the most interesting conclusions one can draw from the analysis of the discourse of scientific popularization is that synonymy (including reformulations, paraphrases, etc.) is used to make a non-specialized audience understand the new concepts resulting from new discoveries. Different denominations are used, intended to help the reader to understand the limits of the concept […].

In order to account for such instances of functionally-motivated terminological variation, several authors have developed procedures to identify terms characterized by different degrees of specialization, notably by relying on corpora and corpus-related methods. Most of these procedures are based on comparisons between specialized and non-specialized (reference) corpora (see Marín 2014 for an overview), on the assumption that lexical units that are salient/typical in the specialized discourse under consideration tend to have significantly higher frequency in the former than in the latter. The resulting lists of lexical units are then subjected to further analysis for term classification.

The methods developed by Kwary (2011) and Ha and Hyland (2017), both concerned with financial terminology, involve the use of corpus-external resources. The former draws on input by a professional lexicographer and a subject-field expert, who are asked to classify lexical units as domain-specific financial terms, as terms belonging in other subject-fields, or as "non-terms" belonging in the general language (see also Chung, Nation 2003). The latter proposes a five-point technicality scale, ranging from "least technical" (e.g. "fee") to "most technical" (e.g. "subprime"), and assigns terms to one category or another based on manual analysis of their meaning as codified in general language and special language lexicographic resources, as well as general-language corpora. In both cases, (non-)availability of resources in the domain under investigation, be they experts or dictionaries/glossaries, may be a major hindrance, especially when working with under-resourced languages or domains.

A fully corpus-based, domain-independent method is the one developed by Marín (2016). The author aims to categorize "sub-technical" terms, i.e. terms "common to both the general and the specialized fields or amongst scientific disciplines" (Marín 2016, p. 84), along a continuous scale of specialization, with positive/negative numerical values indicating

*Lingue e Linguaggi*

greater/lower specialization. The classification method itself consists of an algorithm built around Williams' (1998) lexical network model, according to which "the more specialised a sub-technical term [is], the greater its capacity to attract other words in a specialised context, generating more populated networks around it and its collocates" (Marín 2016, p. 94). While firmly rooted in linguistic theory and returning promising results in terms of classification precision, Marín's method has the drawback of involving complex data processing procedures, which may be especially cumbersome for less technically-minded users.

The method for term extraction described in Section 3.3 is less sophisticated than Marin's (2016) or Ha and Hyland's (2017), and affords less precise observations on the status of terms as more/less specialized. However, it does offer the advantage of being easily implemented through freely available, user-friendly software tools. Furthermore, as will become clearer in Section 3.3, its objective is not that of maximising precision in the classification of single terms, but rather that of providing *lists* of terms that, when used in combination with vocabulary profiling tools, can provide reliable indications as to the status of whole texts as specialized/popularized.

# 3. Study setup

## 3.1. Research questions

As mentioned in Section 1, the present study aims to assess whether and how terminological density reflects the distinction between more specialized and more popularized texts. Unlike most previous studies (see Section 2.1), the notion is operationalized adopting replicable corpus-based procedures and freely available, user-friendly software tools both in the identification of domain-specific terms, and in measuring density of such terms in texts.
Taking the domain of food safety as case in point, our research questions are as follows: 1) do specialized and popularized texts display different levels of terminological density at the quantitative level? 2) Can quantitative differences (if any) be explained in terms of discursive strategies adopted in popularized texts to disseminate knowledge to a wide public?

## 3.2. Corpus resources

Two distinct corpus resources were built for the purposes of this study. The first one consists in the study corpus, i.e. a collection of specialized and popularized texts published by EFSA, which made the object of the quantitative and qualitative analysis presented in Section 4. The second one is a purpose-built benchmark corpus, which was used as a basis for extracting a list of terms in the domain of food safety (see Section 3.3). The two corpora are presented in what follows.

### 3.2.1. The EFSA corpus

The main goal of EFSA, as stated on its website, is to produce "scientific opinions and advice that form the basis for European policies and legislation"[2] in the field of food safety and related issues, such as animal health and plant protection. This research-driven, scientific role goes hand in hand with another key part of EFSA's mandate, which is to "communicate clearly […] to the public at large, to help bridge the gap between science and the consumer".[2]
The EFSA corpus includes samples of two types of texts, which reflect the two categories of stakeholders targeted by EFSA (see also Bondi 2015 on EFSA's documentation). The first are samples of scientific opinions,[3] which are "forms of expert-to-expert communication, [whose] intended addressees are mainly the scientific community and economic agents submitting applications for regulated products" (Bondi 2015, p. 95). The second are popularized versions of these texts aimed at a wider public, which can take three main forms:

- "News": these are the most frequent popularized texts on EFSA's website. They "act as key documents for users who are […] interested in checking on the final assessment provided by the Agency on the issues involved", but may also "act as 'trailers' for experts who might then want to go on and read the Opinion document" (Bondi 2015, p. 95).

- "Frequently Asked Questions" (FAQs): these texts have the typical structure of FAQs, i.e. they are lists of questions and answers on topics that are supposed to be of particular relevance to the public, such as chemicals in food or dietary reference values. Not all scientific opinions are accompanied by FAQs.

- "Factsheets": factsheets are summarized versions of scientific opinions targeted to the general public. Unlike FAQs and news, which are published

---

[2] http://www.efsa.europa.eu/en/aboutefsa

[3] EFSA also provides other types of scientific and technical documents, such as statements by the Agency or scientific and technical reports. These are labelled as "other scientific outputs" or "supporting publications", and are less central to EFSA's mission (see http://www.efsa.europa.eu/en/efsajournal/scdocdefinitions). For this reason, only scientific opinions were taken into account.

as text-only web pages, factsheets are published in PDF format, and include images and tables. Only few factsheets have been published, on a limited range of topics; no information is made available by EFSA on the criteria that were adopted for the selection of topics, which range from substances like caffeine or aspartame to medical issues related to food-borne diseases.

The point of departure for corpus compilation was the list of "topics" available on the EFSA website:[4] each document published by EFSA is associated with one such topic, and, in turn, each topic is associated with at least one macro-category. For each topic that was labelled as pertaining to the macro-category of food safety (as opposed to, e.g., biological hazards), a manual check was performed on the number and types of documents available. Only topics for which EFSA had published at least one scientific opinion and at least three popularized texts (belonging to one of the three categories mentioned above) were retained. All available FAQs and factsheets on these topics were downloaded, plus a maximum of two scientific opinions and ten news per topic. Since the measures of terminological density adopted (see Section 3.3) are sensitive to extreme differences in document length, news texts shorter than 200 words were discarded, while scientific opinions, which tend to be very lengthy documents, were cut to a maximum of 6,200 words, corresponding to 1.5 times the size of the longest popularized text. Table 1 provides information on the resulting corpus and its sub-corpora.

|                               | Factsheets | FAQs   | News    | Scientific opinions | TOTAL   |
|-------------------------------|-----------|--------|---------|---------------------|---------|
| Texts                         | 7         | 15     | 191     | 60                  | 273     |
| Tokens                        | 11,446    | 33,285 | 106,032 | 319,040             | 469,803 |
| Types                         | 1,778     | 3,153  | 5,841   | 14,532              | 25,304  |
| Text length: mean             | 1,638     | 2,221  | 556     | 5,385               | –       |
| Text length: std. deviation   | 400       | 1,240  | 279     | 412                 | –       |
| EFSA topics (examples)        | Acrylamide, Aspartame, Bisphenol A, Caffeine, Chemicals | | | | |

Table 1
Size information about the EFSA corpus.

## 3.2.2. The food safety benchmark corpus

The aim of the benchmark corpus was to provide an "external" resource (with respect to the study corpus) of specialized texts from which food safety terms could be extracted. Research articles from the *Journal of food safety*,[5] published by Wiley, and the *Journal of food protection*,[6] the official journal of the *International Association for food protection*, were included.

---

[4] http://www.efsa.europa.eu/en/topics-complete-list
[5] https://onlinelibrary.wiley.com/journal/17454565
[6] http://jfoodprotection.org/

To maximize comparability between the study corpus and the benchmark one, for each topic represented in the EFSA corpus a search was performed via Google in the databases of the two journals, using the name of the topic itself (see examples in Table 1) as a query term. A maximum of 10 articles per topic per journal were downloaded. Since texts were only used for term extraction and not for comparison with the EFSA texts, differences in text length were not considered as relevant, and hence whole texts were retained. Information on the resulting corpus can be found in Table 2.

| | |
|---|---:|
| Texts | 334 |
| Tokens | 2,359,388 |
| Types | 73,980 |
| Sources | *Journal of food safety, Journal of food protection* |

Table 2
Basic statistics on the food safety benchmark corpus.

## 3.3. Software tools and vocabulary lists

The method proposed in this study to measure terminological density is based on so-called vocabulary profiling (VP) software. Initially developed for language teaching purposes, and especially to assess intelligibility of texts by students of English for Academic Purposes, VP tools use lists of words matching pre-defined criteria, for instance featuring high-frequency words with which students are expected to be familiar vs. low-frequency words which are less likely to be known to them (e.g. Webb, Nation 2008). Based on such lists, VP tools can then measure the density of the (categories of) words of interest in the texts to be analyzed. Density itself is expressed as a proportion, e.g. of high-frequency vocabulary to the overall number of tokens or types in texts.

The study presented in Section 4 draws on a combination of existing and purpose-built vocabulary lists, representing general and specialized language. To assess the degree to which EFSA texts rely on general, high-frequency vocabulary, the *General Service List* (GSL, West 1953) was used. The GSL was complemented by a set of specialized term lists, and specifically: *a)* the *Academic Vocabulary List* (AVL, Gardner, Davies 2014), a genre-specific list of terms that are expected to be typical of scientific/academic writing in general, irrespective of differences among subject fields; and *b)* a domain-specific term list, featuring terms that are expected to be typical of the sub-

domain of food safety, created specifically for the study with the help of the CasualConc concordancer.[7]

Following common practice in terminology extraction (see Section 2.2), the list of food safety terms was created on the basis of a statistical comparison between the food safety benchmark corpus (see Section 3.2.2) and a general corpus of English, i.e. a 17-million-word subset of ukWaC (Baroni *et al.* 2009). Terms were defined as single-word lexical units occurring with a significantly higher frequency in the food safety corpus than in the general language corpus.[8]

Further, to reflect the view that domain-specific terms can differ in their degree of specialization (see Section 2.2), a set of parameters were taken into account to categorize them as "core" vs. "non-core" terms, corresponding to less vs. more specialized lexical units. Core terms were defined as those units that were distributed across at least 30% of the benchmark corpus texts (a parameter adopted, among others, by Coxhead (2000) and Marín (2016)), or whose frequency comparison with respect to the ukWaC corpus displayed a low effect size,[9] pointing to the fact that they tended to also be frequent in general language. By contrast, non-core terms were those which occurred in less than 30% of the benchmark corpus texts, or whose frequency comparison displayed a high effect size. Table 3 presents examples of the two types of terms.

| Core terms | Non-core terms |
| --- | --- |
| food, see, salmonella, coli, samples | campylobacter, inactivation, STEC, ND, jejuni |

Table 3
Top 5 core and non-core food safety terms, in decreasing frequency order.

Being based on single lexical units and simple frequency/distributional information, this method has obvious disadvantages: it does not take into account multi-word terms, and disregards refined semantic distinctions between specialized and non-specialized senses of words. However, as also argued in Section 2.2, it was not conceived as a full-fledged terminology classification method, but rather as a viable, easy-to-implement alternative to more sophisticated classification systems see Section 2.2. We leave it to further

---

[7] https://sites.google.com/site/casualconc/. CasualConc was the concordance of choice since it provides more advanced functionalities compared to other freely available concordancers. E.g. it computes word distribution values across corpus texts as well as other useful statistics, such as effect size measures for keywords (see Footnote 9).

[8] Specifically, terms were defined as lexical units for which the comparison returned a Log-Likelihood value equal or greater than 3.84, corresponding to a significance level of $p < 0.05$ (Rayson, Garside 2000).

[9] Effect size was assessed using the %DIFF measure (Gabrielatos, Marchi 2012). Low/high %DIFF values were defined as values below/above the median of %DIFF values for all terms.

work to experiment with the output of these systems as sources of data for terminological density measurements.

In the final stages of the data preparation phase, the four word lists (i.e. the GSL, AVL, and the core and non-core food safety term lists) were loaded in the VP tool AntWord Profiler (Anthony 2014). The four EFSA sub-corpora were then also loaded, one at a time, to measure their terminological density. This was calculated on a text-by-text basis, and defined as the proportion of specialized terms to the overall number of words per text; in order to reduce the effect of text length, analyses were carried out taking into account word types rather than tokens. Figure 1 displays the AntWord Profiler settings adopted.
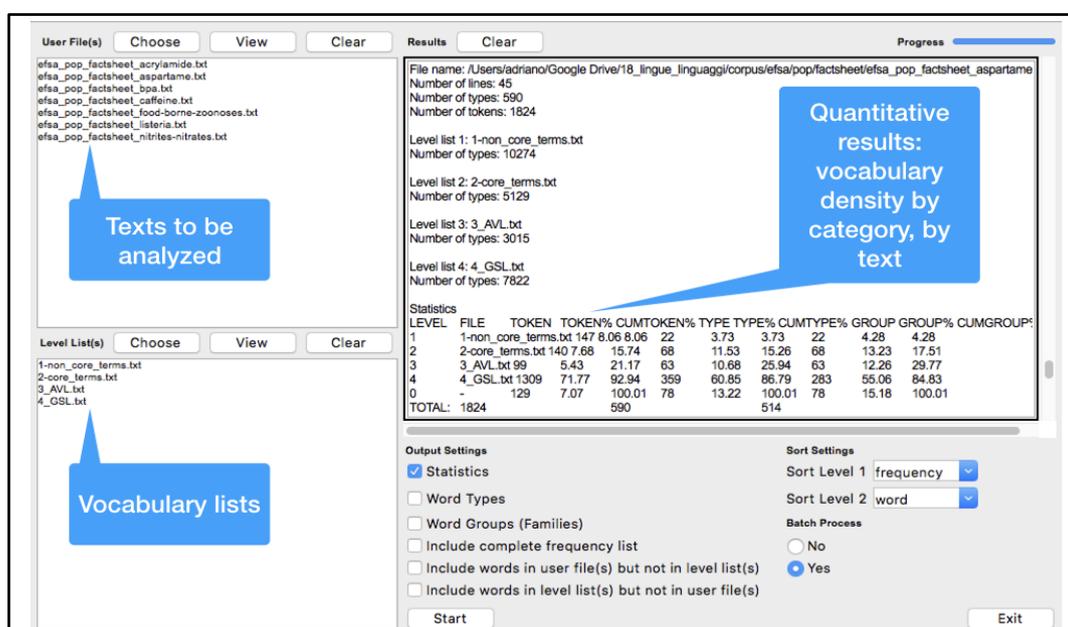


Figure 1
The AntWord Profiler interface.

## 3.4. Statistical procedures for the quantitative analysis

Results obtained from AntWord Profiler were exported to a spreadsheet, and statistical analyses were carried out to detect significant differences in terms of terminological density among sub-corpora. Analyses were performed in R.[10] Kruskal-Wallis tests were used to compare all sub-corpora. In the cases where they returned significant results, they were followed by post-hoc comparisons between pairs of sub-corpora using Wilcoxon rank sum tests with Holm correction. Results of these quantitative analyses are presented in Section 4.1.

---

[10] https://www.r-project.org

# 4. Results

## 4.1. Quantitative analysis: comparing terminological density across specialized and popularized texts

### 4.1.1. Overall analysis

To provide a general overview of the data, a mosaic plot is presented in Figure 2 showing the overall vocabulary profiles of the 4 EFSA sub-corpora. The height of the coloured bars indicates the density, expressed as a percentage, of three categories of vocabulary, i.e. *a)* specialized vocabulary, which corresponds to the density of genre- and domain-specific terms (see Section 3.3); *b)* general vocabulary, corresponding to the density of high-frequency words included in the GSL; and c) "other" vocabulary, i.e. a miscellaneous category of words which were not attested in any of the available lists, and which includes numbers, low-frequency words, words not typical of the domain, etc..
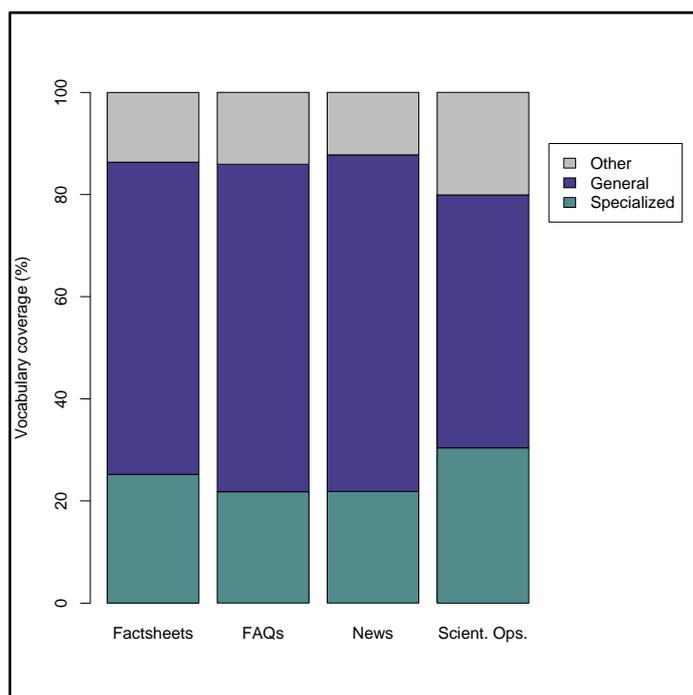


Figure 2
Overall vocabulary profiles of EFSA texts.

Focusing on the first two categories, FAQs and news are very similar in terms of density of both specialized and general vocabulary. Term density in factsheets is slightly higher, but, overall, the three popularized genres are more similar to each other than they are to scientific opinions, which display the highest density of terms and the lowest density of general vocabulary.

Zooming in on the parameter that is most central to the analysis, i.e. density of terms, the results of the statistical tests confirm the observations made on the basis of data visualization. A Kruskal-Wallis test reveals that differences among the four sub-corpora are significant (Kruskal-Wallis $\chi^2(3)=118.35$, $p < 0.001$), and results of post-hoc comparisons (see Tables 4a and 4b) show that terminological density is significantly higher in scientific opinions than in any of the popularized texts. By contrast, none of the comparisons between popularized texts display significant differences.

|  | Median (%) | SD |  | Factsheets | FAQs | News |
|---|---|---|---|---|---|---|
| Factsheets | 25.0 | 2.8 | FAQs | 0.051 | - | - |
| FAQs | 22.1 | 3.0 | News | 0.051 | 0.946 | - |
| News | 21.9 | 4.1 | Scient. Ops. | **0.002** | **<0.001** | **<0.001** |
| Scient. Ops. | 29.9 | 3.2 |  |  |  |  |

Table 4a; Table 4b
Terminological density; *p*-values of pairwise comparisons (significant results in bold).
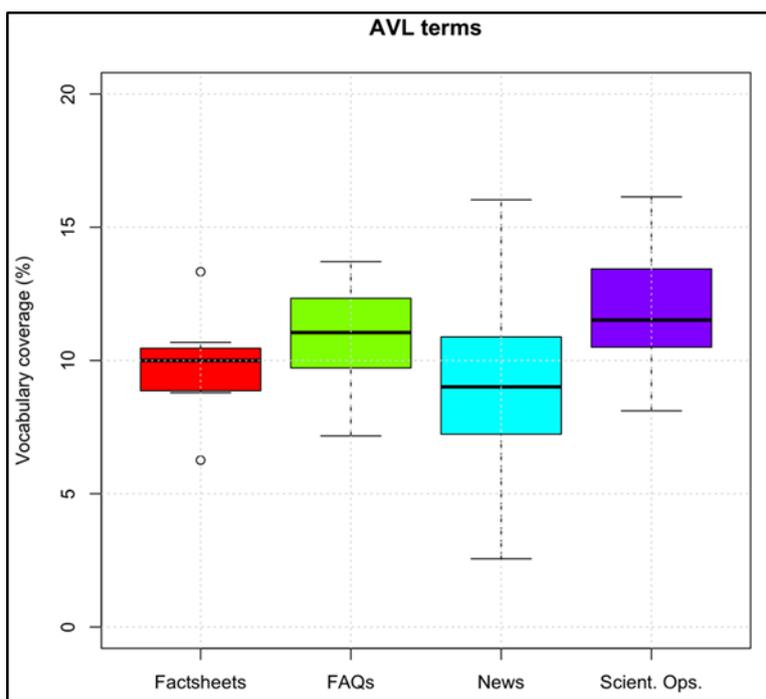


Figure 3
Density of academic terms in EFSA texts.

## *4.1.2. Detailed analysis*

As a follow-up to the overall analysis, the density of the different categories of terms was also considered. Graphical representations are provided in this case in the form of boxplots. For term density in the texts of each sub-corpus, these show the central tendency (median values, visualized as a thick black line cutting the "boxes" in two), the spread of the data around it (corresponding to the size of the boxes; the larger the box, the more different the texts within the sub-corpora), and the minimum and maximum values of the distribution (represented by the extreme points above or below the boxes).

Starting from the academic (AVL) terms (Figure 3) the comparison among sub-corpora returns significant results (Kruskal-Wallis $\chi^2(3)=54.2$, p < 0.001); post-hoc tests (see Tables 5a and 5b) reveal instead a different picture. As expected, scientific opinions display the highest density of academic terms. This is significantly higher than the density in factsheets and news, but not significantly higher than the density in FAQs, which in turn are significantly denser than news (but not factsheets). In this case, popularized texts do not form a homogenous group, and this is especially due to the high density of academic terms in FAQs.

|             | Median (%) | SD  |              | Factsheets | FAQs  | News   |
|-------------|-----------|-----|--------------|-----------|-------|--------|
| Factsheets  | 10.0      | 2.2 | FAQs         | 0.441     | -     | -      |
| FAQs        | 11.1      | 1.9 | News         | 0.517     | **0.036** | -      |
| News        | 9.0       | 2.5 | Scient. Ops. | **0.043** | 0.441 | **<0.001** |
| Scient. Ops.| 11.5      | 2.0 |              |           |       |        |

Table 5a; Table 5b
Density of academic terms; *p*-values
of pairwise comparisons (significant results in bold).

The results for the core terms analysis are displayed in Figure 4. Scientific opinions are, again, the terminologically densest texts. Unlike in the case of academic terms, here factsheets are the odd ones out among popularized texts. Their term density is not significantly different from that of scientific opinions, but is significantly higher than that of both FAQs and news (Kruskal-Wallis $\chi^2(3)=73.6$, p < 0.001; results of post-hoc tests in Tables 6a and 6b). In other words, when it comes to density of less specialized terms related to food safety, the EFSA texts cluster in two groups, with scientific opinions and factsheets being significantly denser than news and FAQs.
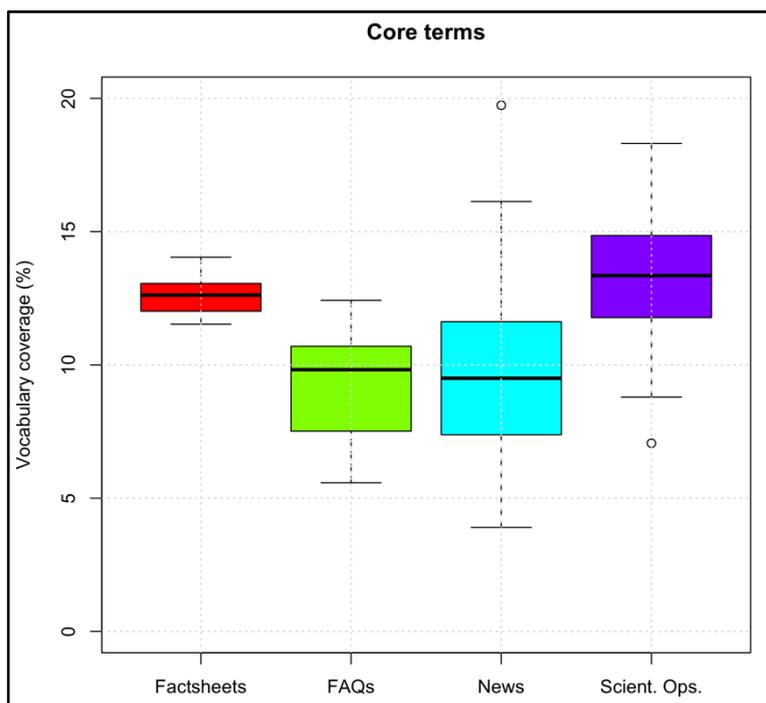
Figure 4
Density of core terms in EFSA texts.

| | Median (%) | SD | | Factsheets | FAQs | News |
|---|---|---|---|---|---|---|
| Factsheets | 12.6 | 0.9 | FAQs | **<0.001** | - | - |
| FAQs | 9.8 | 2.0 | News | **0.008** | 0.551 | - |
| News | 9.5 | 2.9 | Scient. Ops. | 0.430 | **<0.001** | **<0.001** |
| Scient. Ops. | 13.4 | 2.1 | | | | |

Table 6a; Table 6b
Density of core terms; *p*-values
of pairwise comparisons (significant results in bold).

We finally consider non-core terms (Figure 5), corresponding to the most specialized domain-specific terms. The pattern here resembles the overall one presented in Section 4.1.1, with scientific opinions being significantly denser than any of the popularized texts (Kruskal-Wallis $\chi^2(3)=53.8$, $p<0.001$; results of post-hoc tests in Tables 7a and 7b). Among these, however, the news sub-corpus is significantly denser than the sub-corpus of FAQs (but not denser than that of factsheets). This is probably due to the relatively large dispersion of data (see the large light blue box in Figure 5, as well as the SD value in Table 7a), pointing to the fact that news constitute a less homogeneous group than the rest of the popularized texts: some of the news texts have a density of non-core terms that is comparable to that of the densest scientific opinions. It should be noticed that such relative lack of homogeneity is also testified by the analyses of academic and core terms: SD values pertaining to news are

*Lingue e*
*Linguaggi*

consistently the highest ones among popularized texts.



Figure 5
Density of non-core terms in EFSA texts.

|  | Median (%) | SD |
|---|---|---|
| Factsheets | 2.4 | 1.2 |
| FAQs | 1.9 | 0.8 |
| News | 2.7 | 1.7 |
| Scient. Ops. | 4.7 | 2.1 |

|  | Factsheets | FAQs | News |
|---|---|---|---|
| FAQs | 0.132 | - | - |
| News | 0.976 | **0.028** | - |
| Scient. Ops. | **0.028** | **<0.001** | **<0.001** |

Table 7a; Table 7b
Density of non-core terms; *p*-values
of pairwise comparisons (significant results in bold).

Summing up, differences in terminological density were observed not only between specialized and popularized texts, but also between the different genres of popularization; such differences are related to the different types of terms considered, such that different genres are characterized by relative overuse of specific categories of terms. The next Section aims at explaining the quantitative trends observed here by adopting a more qualitative perspective, taking a closer look at discursive patterns surrounding the use of terms in popularized texts.
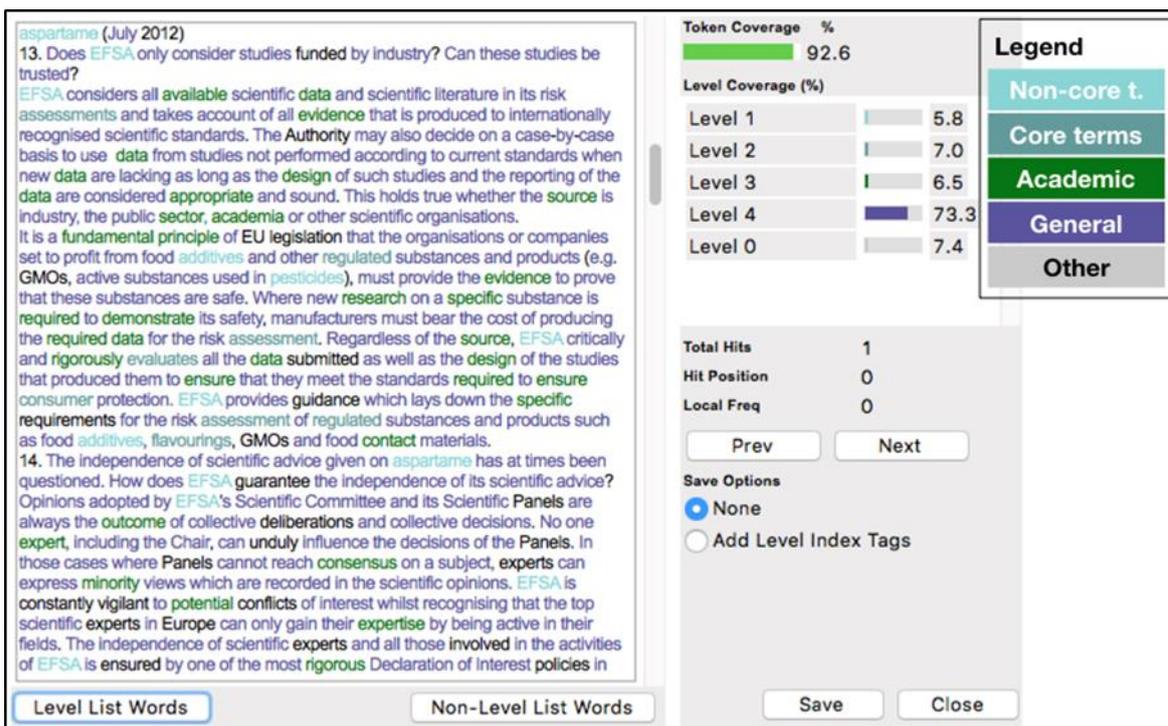
Figure 6
The text analysis function in AntWord Profiler.

## 4.2. Qualitative analysis: exploring terminological density in popularized texts

In order to look for possible explanations of the patterns of terminological density observed in popularized texts (see Section 4.1), we will consider here single texts belonging to different types of popularized genres and displaying especially high values of density of terms. To do so, we will rely on another feature of the AntWord Profiler tool (see Section 3.3), which makes it possible to visualize, through colours, which words were assigned to which vocabulary category in a text. Figure 6 illustrates this visualization function.

Starting from FAQs, a possible reason for their relatively high density of academic vocabulary is the emphasis they place not only on the results of EFSA's research initiatives on specific topics, but also on the methods and the wider scientific context of research itself. Examples 1 and 2 below, taken from FAQs on endocrine active substances and aspartame respectively, illustrate the point.[11] Underlined words correspond to terms classified in AVL as typical of academic registers:

---

[11] If not specified otherwise, examples are selected from the five texts displaying the highest density values for the category of terms under consideration.

1) On 4 July 2017 EU Member States voted in favour of the European Commission's proposal on scientific criteria to <u>identify</u> endocrine disruptors in the field of plant protection products. The criteria <u>identify</u> known and <u>presumed</u> endocrine disruptors. They also <u>specify</u> that the <u>identification</u> of an endocrine disruptor should be carried out by taking into account all <u>relevant</u> scientific <u>evidence</u> including animal, in vitro or in silico studies, and using a weight-of-<u>evidence</u> <u>approach</u>.

2) Where new <u>research</u> on a <u>specific</u> substance is <u>required</u> to <u>demonstrate</u> its safety, manufacturers must bear the cost of producing the <u>required</u> <u>data</u> for the risk assessment. Regardless of the <u>source</u>, EFSA critically and <u>rigorously</u> evaluates all the <u>data</u> submitted as well as the <u>design</u> of the studies that produced them [...].

In these examples, the aim of EFSA is not that of explaining whether or why endocrine active substances or aspartame may be dangerous to consumers' health; rather, it aims at describing the working methods and the guiding principles that were adopted in carrying out research, possibly addressing an audience of scientists or expert legislators (see the first sentence of Example 1). It should be noticed, however, that FAQs are composite texts made of several questions and answers, each potentially targeted to a different audience. In this light, the high density of academic terms, but not, say, of domain-specific terms (see Section 4.1.2), could be due to the fact that these popularized texts are targeted to a mixed audience, which may include academics and other experts in addition to a more "traditional" audience of lay persons.

Moving on to factsheets, it will be remembered that this sub–corpus displays a significantly higher density of core (less specialized) terms than the other popularized sub-corpora. Analysis of the top 5 texts by term density reveals that core terms are often either accompanied by definitions/paraphrases explaining them (see Example 3; core terms underlined), or may be part of the explanations themselves (Example 4).

3) <u>Listeria</u> is a group of <u>bacteria</u> that contains ten <u>species</u>. One of these, <u>Listeria monocytogenes</u>, causes the disease "listeriosis" in humans and animals.

4) <u>Exposure</u> – <u>concentration</u> or amount of a particular substance that is taken in by an individual [...] in a specific frequency over a certain amount of time. When experts assess <u>consumers'</u> <u>dietary</u> <u>exposure</u> to a <u>chemical</u> substance, they combine data on its <u>concentrations</u> in food with the quantity of those foods <u>consumed</u>.

In Example 3 and 4, the core terms "Listeria" and "exposure" are clearly thought to be in need of an explanation, which in the former case is provided through an in-text definition ("<u>Listeria</u> is a group of <u>bacteria</u> that contains ten <u>species</u>"); in the latter case, "exposure", a recurrent term in the text from which the example is taken, i.e. a factsheet on food nitrites and nitrates, is defined within a glossary (a strategy that is adopted for selected terms in all the factsheets considered). In turn, the definitions provided rely on other core terms. Factsheets thus seem to target a more homogeneous audience than

FAQs: the frequent use of terms characterized by a low degree of specialization or accompanied by explanations/paraphrases, paired with a non-significant overuse of the most specialized terms (see Section 4.1.2), seem to testify to the fact that the main audience of factsheets is mainly one of non-experts.

The same cannot be said of news, the last category of texts considered. In Section 4.1.2, the news sub-corpus was revealed as the one characterized by the highest degree of variation, especially in the analysis pertaining to non-core, highly specialized terms. Hence, part of the news texts display a degree of term density that is comparable to that of other popularized genres; other news texts, however, display peaks of density that make them more similar to specialized texts than to popularized ones. Closer inspection of the former set of texts shows that, like in the case of factsheets, highly specialized terms, if present, are accompanied by definitions or paraphrases (see Example 5; non-core terms underlined).

5) At the end of last year the European Food Safety Authority (EFSA) published its first annual report[1] on infectious diseases <u>transmissible</u> from animals to humans (<u>zoonoses</u>) which affect over 380,000 EU citizens per year.

By contrast, the news texts displaying high term density values seem to reproduce the style and content of the scientific opinions that they should contribute to disseminate (see Examples 5 and 6).

6) Bacteria that produce <u>ESBL</u> <u>enzymes</u> show multi-drug resistance to β-lactam <u>antibiotics</u>, which include <u>penicillin</u> <u>derivatives</u> and <u>cephalosporins</u>. The prevalence of <u>ESBL</u>-producing E. coli varied across countries, from low to very high.

7) The report also provides data on other important <u>zoonotic</u> diseases such as Escherichia coli (<u>verocytotoxin</u> and <u>Shiga</u> toxin producing E. coli, collectively termed <u>VTEC</u>), <u>Mycobacterium</u> bovis, <u>Brucella</u>, <u>Yersinia</u>, <u>Trichinella</u>, Echinococcus, and <u>Toxoplasma</u>.

This is not to say that no attempt is made in these examples at paraphrasing or expanding on the most technical notions (e.g. "β-lactam antibiotics […] include penicillin derivatives and cephalosporins", and "verocytotoxin and Shiga toxin producing E. coli, collectively termed VTEC"). Nonetheless, the degree of domain expertise expected on the part the reader is arguably higher than in factsheets (e.g. readers are expected to understand what "cephalosporins" are). The accumulation of several highly-specialized terms ("Brucella, Yersinia, Trichinella") further contributes to the impression that these news items are targeted to an expert, rather than lay, audience (see also Bondi 2015, p. 95).

The next Section summarizes and discusses these results.

## 5. Conclusion

This article set out to test empirically the hypothesis that terminological density can be used as a measurable indicator to tell apart specialized and popularized texts, and that quantitative differences in terms of this parameter can be related to differences in terms of discursive strategies and target audiences of the texts considered. Focusing on the documentation published by EFSA on its website, it compared specialized scientific opinions and three different popularized genres, namely factsheets, FAQs and news. To do so, it proposed a method to operationalize and measure terminological density that is entirely based on replicable procedures, and which can be implemented using freely-available, user-friendly software tools.

Results of the quantitative analyses provided evidence that terminological density as operationalized here does reflect the difference between specialized and popularized texts, and can further point to more subtle differences related to popularized texts' varied audiences and discursive strategies. Analyses showed that scientific opinions generally display a higher density of terms than their popularized counterparts. At the same time, popularized texts were revealed as a non-homogeneous group, providing support for the view that popularization, on a par with specialization, should be conceived as a multi-faceted and gradable phenomenon (see Section 1).

Specifically, FAQs and news were found to display varying degrees of specialization (or popularization, depending on the perspective adopted), by combining popularizing strategies and a relative lack thereof: the former make very limited use of domain-specific terms and instead rely extensively on academic vocabulary (at least in certain text parts); the latter include texts where highly specialized food safety terms are accompanied by definitions and paraphrases, and other texts which closely reproduce the contents and style of specialized scientific opinions. Such differences, it was argued, might be related to the mixed audiences of both lay persons and (semi-)experts to which these texts are addressed. By contrast, factsheets emerged as the most "prototypically" popularized genre, positioning themselves halfway between specialized texts and the other popularized genres in terms of two out of three parameters of term density. The only category of terms for which factsheets are significantly denser than FAQs and news, and as dense as scientific opinions, is that of core terms: these, however, are systematically defined either in texts or within glossaries, possibly in an attempt to "[inform] and […] extend the reader's knowledge" (Gotti 2014, p. 17), the primary function of popularized texts.

To conclude, we would like to go back to the idea that vocabulary in popularized texts is chiefly "drawn from the general language" (Gotti 2011, p. 207; see Section 1). While the overall analysis of the vocabulary profiles of

EFSA texts provided support for the claim, the role of terminology in popularization discourse should not be underestimated: all of EFSA popularized texts were shown to display at least a degree of terminological density, and in this sense to display what Delavigne (2003, p. 84; my translation) calls the "discursive dilution of technical jargon". Somewhat paradoxically, popularized texts need to either explain or paraphrase technical terms in order to be understood by the audience; at the same time, however, they need these terms to gain legitimacy and authoritativeness in the eyes of the audience itself. Striking a balance between these two poles, as the analysis has hopefully suggested, is one of the biggest challenges for (successful) popularization.

**Bionote:** Adriano Ferraresi is Senior assistant professor at the Department of Interpreting and Translation of the University of Bologna, where he teaches a Master's module on information mining and terminology. His main research interests are in the area of quantitative, corpus-informed methods for the study of specialized varieties of English and in translation studies. He has given several invited seminars on these topics at various Italian and European universities. He recently published the monograph *Collocations in institutional academic English. Corpus and experimental perspectives* (Aracne).

**Author's address:** adriano.ferraresi@unibo.it

# References

Andersen Ø. 1996, *Degree of specialization, network density and status - a text typological approach*, in Budin G. (ed.), *Multilingualism in specialist communication*, International Network for Terminology, Vienna, pp. 85-93.

Anthony L. 2014, *AntWordProfiler* [Software]. http://www.laurenceanthony.net/ (24.04.2019).

Baroni M., Bernardini S., Ferraresi A. and Zanchetta E. 2009, *The WaCky Wide Web: A collection of very large linguistically processed web-crawled corpora*, in "Language Resources and Evaluation" 43 [3], pp. 209-226.

Bondi M. 2015, *Language policy in web-mediated scientific knowledge dissemination: A case study of risk communication across genres and languages*, in Plo Alastrué R. and Pérez-Llantada C. (eds.), *English as a scientific and research language*, Vol. 2, De Gruyter, Berlin, pp. 85-111.

Cabré M.T. 2003, *Theories of terminology*, in "Terminology" 9 [2], pp. 163-199.

Cabré M.T. 2010, *Terminology and translation*, in Gambier Y. and van Doorslaer L. (eds.), *Handbook of Translation Studies*, Vol. 1, Benjamins, Amsterdam, pp. 356-365.

Calsamiglia H. and Van Dijk, T.A. 2004, *Popularization discourse and knowledge about the genome*, in "Discourse & Society" 15 [4], pp. 369-389.

Chung, T.M. and Nation, P. 2003, *Technical vocabulary in specialised texts*, in "Reading in a Foreign Language" 15 [2], pp. 103-116.

Constantinou M. 2014, *Conceptual metaphors of anger in popularized scientific texts. A contrastive (English/Greek/French) cognitive-discursive approach*, in Baider F.H. and Cislaru G. (eds.), *Linguistic approaches to emotions in context,* Benjamins, Amsterdam, pp. 159-188.

Corpas Pastor G. and Seghiri M. 2010, *Size matters: A quantitative approach to corpus representativeness*, in Rabadan R. (ed.), *Language, translation, reception. To honor Julio César Santoyo*, Universidad de Leon, Leon, pp. 111-145.

Coxhead A. 2000, *A new academic word list*, in "TESOL Quarterly" 34 [2], pp. 213-238.

Delavigne V. 2003, *Quand le terme entre en vulgarization*, in "Terminologie et Intelligence artificielle", pp. 80-91.

Freixa J. 2006, *Causes of denominative variation in terminology*, in "Terminology" 12 [1], pp. 51-77.

Gabrielatos C. and Marchi A. 2012, *Keyness: Appropriate metrics and practical issues*, paper presented at *Critical approaches to discourse studies 2012*, Bologna, 14 September 2012.

Gardner D. and Davies M. 2014, *A New Academic Vocabulary List*, in "Applied Linguistics" 35 [3], pp. 305-327.

Gotti M. 2011, *Investigating Specialized Discourse*, 3rd ed., Peter Lang, Bern.

Gotti M. 2014, *Reformulation and recontextualization in popularization discourse*, in "Ibérica" 27, pp. 15-34.

Ha A.Y.H. and Hyland K. 2017, *What is technicality? A Technicality Analysis Model for EAP vocabulary*, in "Journal of English for Academic Purposes" 28, pp. 35-49.

Krüger R. 2015, *The interface between scientific and technical translation studies and cognitive linguistics*, Frank &Timme, Berlin.

Kwary D.A. 2011, *A hybrid method for determining technical vocabulary*, in "System" 39, pp. 175-185.

Marín M.J. 2014. *Evaluation of five single-word term recognition methods on a legal English corpus*, in "Corpora" 9 [1], pp. 83-107.

Marín M.J. 2016, *Measuring the degree of specialisation of sub-technical legal terms through corpus comparison. A domain-independent method*, in "Terminology" 22 [1], pp. 80-102.

Myers G. 2003, *Discourse studies of scientific popularization: Questioning the boundaries*, in "Discourse Studies" 5, pp. 265-279.

Nation I.S.P. and Heatley A. 2002, *Range: A program for the analysis of vocabulary in texts* [Software]. http://www.vuw.ac.nz/lals/staff/paul-nation/nation.aspx (24.04.2019).

Pearson J. 1998, *Terms in Context*, Benjamins, Amsterdam.

Piscopiello M.G. and Bertaccini F. 2009, *Densità e ricchezza terminologica: generi testuali a confronto*, in "MediAzioni" 7. http://mediazioni.sitlec.unibo.it (24.04.2019).

Rayson P. and Garside R. 2000, *Comparing corpora using frequency profiling*, in *Proceedings of Workshop on Comparing Corpora*, held at the ACL conference 2000, pp. 1-6.

Skibitska O. 2015. *The language of tourism: Translating terms in tourist texts*, in "Translation journal", October issue. http://translationjournal.net/October-2015/the-language-of-tourism-translating-terms-in-tourist-texts.html (24.04.2019).

Webb S. and Nation P. 2008, *Evaluating the vocabulary load of written text*, in "TESOLANZ Journal" 16, pp. 1-10.

West M. 1953, *A General Service List of English Words*, Longman, London.

Williams G. 1998, *Collocational networks: Interlocking patterns of lexis in a corpus of plant biology research articles*, in "International Journal of Corpus Linguistics" 3 [1], pp. 151-171.