

# Terminological Consistency in Croatian Translations of EU Legislation: A Corpus-Based Study

---

Marić, Laura

Master's thesis / Diplomski rad

2024

*Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj:* **University of Zagreb, Faculty of Humanities and Social Sciences / Sveučilište u Zagrebu, Filozofski fakultet**

*Permanent link / Trajna poveznica:* <https://urn.nsk.hr/urn:nbn:hr:131:592331>

*Rights / Prava:* [In copyright](#) / [Zaštićeno autorskim pravom.](#)

*Download date / Datum preuzimanja:* **2024-07-29**



Sveučilište u Zagrebu  
Filozofski fakultet  
University of Zagreb  
Faculty of Humanities  
and Social Sciences

*Repository / Repozitorij:*

[ODRAZ - open repository of the University of Zagreb  
Faculty of Humanities and Social Sciences](#)



University of Zagreb  
Faculty of Humanities and Social Sciences  
Department of English

**Terminological Consistency in Croatian Translations of EU Legislation:  
A Corpus-Based Study**

Master's thesis

Laura Marić

Supervisor: Nataša Pavlović, Ph.D.

Zagreb, 2024

Sveučilište u Zagrebu

Filozofski fakultet

Odsjek za anglistiku

**Terminološka dosljednost u hrvatskim prijevodima zakonodavstva EU-a:  
korpusno istraživanje**

Diplomski rad

Laura Marić

Mentor: prof. dr. sc. Nataša Pavlović

Zagreb, 2024.

## Table of Contents

1. Introduction.....	1
2. Previous research.....	2
3. Key concepts .....	3
3.1. <i>Terminology and term</i> .....	3
3.2. <i>Terminological consistency and variation</i> .....	5
4. Aims and hypotheses.....	6
5. Methodology .....	7
5.1. Corpus compilation .....	7
5.2. Term list.....	9
5.3. Analysis .....	11
6. Results .....	15
6.1. Descriptive results .....	15
6.1.1. <i>Frequency</i> .....	15
6.1.2. <i>Number of words</i> .....	15
6.1.3. <i>Herfindahl Hirschman Index score</i> .....	16
6.2. Correlation between frequency and terminological consistency .....	16
6.3. Correlation between structure and terminological consistency .....	17
7. Discussion .....	18
7.1. Terminological (in) consistency in Croatian translations.....	18
7.2. Correlation between frequency and terminological consistency .....	18
7.3. Correlation between structure and terminological consistency .....	20
7.4. Limitations and relevance.....	21
8. Conclusion.....	21
References.....	23
Appendix A.....	25
Appendix B.....	26
Appendix C.....	28

## **Abstract**

The EU produces thousands of translations every year since all legally-binding documents have to be available in all 24 official languages. A particular emphasis is placed on terminology, or rather terminological consistency due to its impact on clarity and legal certainty. However, a question still arises about how consistently the terminology is really used considering the fast pace the translations are produced at. The aim of this study is, therefore, to check for inconsistencies in the use of trade-related terminology in Croatian translations of EU legislation, and to determine if there is a correlation between the frequency and structure of terms and their consistency. Terminological consistency is measured using the HHI, following the example of Itagaki et al. (2007). Trade-related terms are taken from IATE and analysed on a corpus compiling English-Croatian translation memories from 2020 provided by the European Commission's DGT. The results confirm the presence of terminological inconsistency and show a weak positive correlation between consistency and frequency, and no correlation between consistency and structure. The main contribution of the study is to show the usefulness of EU materials and the HHI for linguistic and terminological research and create a space for further discussions about terminological consistency.

Key words: terminology, terminological consistency, European Union, EU translation, Herfindahl–Hirschman Index

## 1. Introduction

The European Union has become the most multilingual body of institutions in the world. According to data provided by European Commission (2023), more than 2.5 million pages get translated every year, with the legislation making up 54% of that amount in 2023. This includes translations into all 24 official languages of the EU as per its language policy<sup>1</sup>. To ensure terminological consistency across all these languages and documents, significant effort has been put into terminological work, with one of the biggest milestones being the introduction of the EU's multilingual terminology database, *Interactive Terminology for Europe* or IATE in 2004. However, some research suggests that there might still be terminological inconsistencies, especially in Croatian translations, as the *acquis* was infamously translated “under pressure” and “by numerous translators of various degrees of expertise and experience” (Bratanić & Lončar, 2016, p. 210). That, together with the fact-paced, “multitranslator” environment in which the translations are produced raises questions about how consistently terminology is really used and translated. The aim of this study is therefore to check if there are inconsistencies in the use of terminology in Croatian translations of EU legislation, and furthermore to determine if there is a correlation between the frequency and structure of a term and its terminological consistency where inconsistencies are detected. The scope has been reduced to trade-related terminology, due to the limited scope of a Master’s thesis. The study was conducted on a corpus made up of translation memories comprising texts from 2020, made available by the European Commission's Directorate-General for Translation, and the analysed terms were extracted from IATE. Finally, the method used to measure terminological consistency includes the use of the Herfindahl–Hirschman Index, or HHI, as proposed by Itagaki et al. (2007).

The paper is organized as follows: the next section provides an overview of previous research related to the topic, and it is followed by the *Key terms* section in which central terms and concepts of the study are defined. The fourth section presents the aims and hypotheses, while the fifth section gives a detailed explanation of the used methodology and resources. Finally, in sections *Results* and *Discussion* the study’s findings are reported and further discussed. Various data sets, as analysed and reported on in the study, can be found in the *Appendices* section at the end.

---

<sup>1</sup> <https://www.europarl.europa.eu/factsheets/en/sheet/142/language-policy>

## 2. Previous research

In Croatia, terminology and terminological consistency in Croatian legal translations gained traction in the scholarly sphere when Croatia was preparing for its accession to the EU. One of the works outlining and discussing the pre-accession state of Croatian legal terminology was *Hrvatski jezik na putu u EU* edited by Maja Bratanić (2011). It comprised articles dealing with various terminological issues, from term formation to term standardisation. In one contribution, Bajčić and Stepanić (2011) highlighted the inconsistent use of competition law terminology. Their study, conducted on a number of English-Croatian legal translations and terminological resources, such as *Eurovoc* and *Euroterm*, found that several competition law terms were inconsistently translated. The authors underlined the importance of terminological consistency in legal translation, advocating for more cooperation between domain experts, terminologists and translators in the then upcoming translation of EU legislation. Although the resources used in the study are less relevant today, because Croatia has since joined the EU and its terminology and translation framework, its results illustrate that the problem of terminological inconsistency was present in Croatian legal translation even before the accession, and probably influenced the consistency in the later translation of EU law.

Bratanić and Lončar (2016) also examined the myth of terminological consistency in the EU on the example of the Croatian translation of the *acquis*, documents that constitute the body of EU law. The study exemplified instances of terminological inconsistency, aiming to provide an overview of both linguistic and extra-linguistic reasons behind it. The linguistic reasons mostly included the lack of clarity of definition of terms and their relationship in the national and EU legal systems, while the extra-linguistic ones were related to the relatively adverse circumstances in which the *acquis* was translated.

A more international perspective on terminological work and consistency in the EU is given in the works of Stefaniak (2017) and Pozzo (2020). Stefaniak (2017) gave a detailed overview of the terminological process in the EU and described most common terminological problems in translation. She pointed out the importance of terminological consistency and the consequences inconsistency can have in the legal context of the EU. Pozzo (2020) analysed the impact of multilingualism on the harmonisation process of European private law. She found that despite the EU's efforts to harmonise private law terminology, inconsistency was still present, both at a monolingual level and in translations. Although there is more research concerning the terminology process in the EU, very few studies observe consistency using an example-based, empirical approach. Moreover, all of the aforementioned studies observed

terminological inconsistency using a qualitative approach.

However, Gašpar (2013) and later Gašpar et al. (2022) proved that, by using the HHI method first introduced by Itagaki et al. (2007), terminological consistency can also be quantitatively assessed. The method, which utilizes the Herfindahl–Hirschman Index, commonly used to measure market concentration is more thoroughly explained in the *Key concepts* and *Methodology* sections. The first study (Gašpar, 2013) assessed terminological consistency in a Croatian-English legal parallel corpus which showed an expected low index of consistency in Croatian-English legal translations and supported the implementation of the HHI for the measurement of terminological consistency. In the second study, Gašpar et al. (2022) employed the same method, but on a bigger corpus. The corpus consisted of three types of legal subcorpora, Croatian-English parallel corpus (1991–2009), Latin-English and Latin-Croatian versions of the Code of Canon Law (1983), and the English and Croatian versions of the EU legislation (2013– ). The results confirmed the presence of inconsistency in all corpora, with the consistency index being higher for the English-Croatian language pairs. They also showed a diachronic increase in consistency and supported the implementation of HHI for assessment of terminological consistency on the Croatian-English language pair. However, the study had certain limitations due to the small terminology dataset. Only the most frequent terms from each corpus were analysed, with the numbers of terms per corpus being the following: Croatian-English parallel corpus, 100; Canon Law corpus, 25; EU legislation corpus, 15. The limitations leave the question of how representative the consistency index was as a whole, since there wasn't a great distribution in terms of frequency. That is a question this study tries to bring more insight into. It should also provide a more recent assessment of state of terminological consistency in Croatian translations of EU legislation.

### **3. Key concepts**

#### **3.1. Terminology and term**

Both terminology as a discipline and the term *term* have been defined and redefined by many scholars since the publication of Wüster's seminal work *General Theory of Terminology* in 1979, which set the groundwork for the development of the discipline. Traditionally, terminology has been defined as “the study of and the field of activity concerned with the collection, description, processing and presentation of terms” (Sager, 1990, p. 2). It is also often used to refer to “internally consistent and coherent set of terms belonging to a single subject field” (Sager, 1990, p. 3). It plays a crucial role in specialized fields, including law, as



it provides a framework that establishes connections between specific *concepts*, or ideas, and their lexical counterparts, or *terms*, as well as maps relationships between those concepts. That connection is formed primarily by means of a definition which gives a precise description and reference of the concept, and it is further strengthened by its relationship to other concepts belonging to the same domain. The definition can also be complemented by other morphological, syntactic and pragmatic specifications, such as information about context and usage (Sager, 1990, pp. 21-40). As follows, *terms* can be defined as “items which are characterised by special reference within a discipline” (Sager, 1990, p. 19), as opposed to *words*, which are linguistic units that have a general reference in a language. Since the meaning and function of terms is defined and confined by the domain they exist in, they can also be described as “a functional class of lexical units” (Sager, 1998, as cited in Kockaert & Steurs, 2015, p. 48). In translation, as explained by Fischer (2022), terms can also be evaluated in a broader sense, encompassing any lexical unit which is expected to be translated in a specific way, i.e. any word, phrase or sentence that restricts the translator’s freedom. She also insists that such approach should be followed in all discussions surrounding EU translation, since EU’s terminology work is largely influenced by translation (Fischer, 2022). Moreover, the application of these principles can be noticed in the criteria for inclusion of terms in IATE, as it, in addition to terms, comprises special expressions that are useful in translation, but that would not be considered terms in the traditional sense (Fischer, 2010). Furthermore, one of the main characteristics of terms in traditional terminology is univocity, designating both monosemy (one concept per term) and mononymy (one term per concept) (Temmerman, 2000, p. 10). That means that each concept can be signified by only one term, and in turn that term cannot be used to refer to another concept. This implies the elimination of term synonymy and polysemy as well as terminological variation. This view has, however, been challenged by a more recent sociocognitive terminology theory developed by Rita Temmerman. Temmerman (2000) based her approach on frameworks proposed by sociolinguistics and cognitive linguistics. She describes terms as prototypical in nature, with their categorization being based on similarity, not defining characteristics, and their structure being less delineated and more malleable (Temmerman, 2020, pp. 63-66). Likewise, concepts are replaced by less restrictive and also mostly prototypical *units of understanding*, emphasizing the cognitive dimension they function in. In this view, synonymy, polysemy and terminological variation are consequently accepted as functional, as they express changes in meaning and show different perspectives; “[c]ategories evolve, terms change in meaning, understanding develops” (Temmerman, 2020, p. 16). In the context of EU translation, the

sociocognitive approach can nevertheless be hard to accept because of its implications for law interpretation and legal certainty<sup>2</sup>, even though it might provide a better insight into the ways in which we process and understand terminology than the traditional one. As Bratanić and Lončar (2016) point out, due to the multilingual, multicultural and “multilegal” context that EU translation is situated in translators already have to face many challenges to ensure the uniformity of law interpretation across all languages, as misinterpretations can have serious legal consequences. That uniformity hinges on unambiguity and clarity which could be even harder to achieve in the presence of terminological variation, which is why in such cases “a more traditional approach to term harmonisation and standardisation should still be at the forefront of the discussion” (Bratanić & Lončar, 2016, p. 217). Following these arguments, this study draws on the traditional definitions of terminology and terms and gives precedence to terminological consistency over terminological variation.

### **3.2. Terminological consistency and variation**

Being the focus of this study, terminological consistency also needs to be defined. Although it is a characteristic of both monolingual and translated specialized texts, it is more common in translations and, in this study, it will only be observed from the perspective of translation. *Terminological consistency* can, therefore, be defined as the use of one and the same translation equivalent for a given source term (Gašpar, 2013, p. 1). Additionally, consistency does not only refer to the use of “the same term for the same referent throughout a particular communication“ (Rogers, 2008, p. 107), but also “throughout all communications within a particular organisation if a terminology policy is in place” (Rogers, 2008, p. 107), which is certainly the case in the EU. On the other hand, *terminological variation* is the result of term polysemy and synonymy, i.e. one source term having two or more *terminological variants*, which are all connected to the same concept, and have co-referential status (Gašpar, 2013, p. 16). In the context in which an ideal of consistency is present, terminological variation is replaced by *terminological inconsistency* which is defined as “the use of two or more translation variants for a given source term” (Gašpar et al., 2022, p. 2).

Finally, terminological consistency can be analysed both quantitatively and qualitatively, and in this study the former approach was employed. The method was first proposed by Itagaki et al. (2007) in their study *Automatic Validation of Terminology Translation Consistency with Statistical Method*, where they utilized the Herfindahl–Hirschman Index or HHI, commonly

---

<sup>2</sup> Legal certainty is a principle that “rules should be clear and precise, so that individuals may be able to ascertain unequivocally what their rights and obligations are and may take steps accordingly” (Craig, 2012, p. 549).

used to measure market concentration to measure terminological consistency in a given text. The HHI is usually calculated using the following formula:

$$HHI = \sum_{i=1}^n s_i^2$$

where  $S$  indicates the market share of a firm in the market, and  $n$  is the number of firms. If the index is 10000, or  $100^2$  that means one firm dominates the market. When applied to translation, Itagaki et al. (2007) explain that “ $S$  becomes the ratio of each translation ( $i$ ) to the total number of translations ( $n$ ) within a product” (p. 5). The index was later applied by Gašpar (2013) and Gašpar et al. (2022) in studies on Croatian-English and English-Croatian corpora, using the following adapted formula:

$$C_t = \sum_{n=1}^n \left( \frac{f}{k} \times 100 \right)_i^2$$

The calculation and application of the HHI index in this study is further explained in the *Methodology* section.

#### 4. Aims and hypotheses

The main aim of this study was to check if there are inconsistencies in the use of terminology in Croatian translations of EU legislation. Due to its small scope, the observed terminology was limited to trade-related terms only. Furthermore, where inconsistencies were detected, the second aim was to determine if there is a correlation between a term's frequency and its structure (i.e., the number of words it consists of) on the one hand and how consistently it is translated on the other. The following hypotheses were tested:

H1 There are inconsistencies in the use of trade terminology in Croatian translations of EU legislation.

H2 Terms that have a higher frequency, i.e. those that are used more often, are translated more consistently.

H3 Longer terms, i.e. those that consist of more words are translated less consistently.

The basis for the first hypothesis was the previously discussed expectation that terminology might not be consistently translated based on the context of EU translation and the suggestions of previous studies (e.g. Gašpar, 2022). That hypothesis also lays the foundation for the other two hypotheses: they can be tested only if it is accepted, i.e. if inconsistencies are found.

Furthermore, it is expected that more frequent terms will be translated more consistently because they have more established terminological equivalents due to being used more often. Another aspect that also plays a role in this correlation is the use of, and reliance on translation memories (TMs). All EU translators work with *Euramis*, the EU's central translation memory, which automatically retrieves similar segments that it recognizes as useful for the translation of the new document. This not only allows for the translations to be produced at a faster rate, but also helps ensure consistency across all documents, especially when translating delegated or implementing acts. In 2016, *Euramis* contained over 1 billion segments across all official EU languages (European Commission, 2016). Owing to this fact, it is presumed that terms that are more frequent will be present more often in the TMs and therefore automatically be translated more consistently, as opposed to less frequent terms which might not be present in the TMs and for which the translator will have to do extra research to find or create the appropriate equivalent. More frequent terms are usually also better known, so the translator might know their equivalents even without double-checking the term base or the TM.

The third hypothesis takes into account the number of words that make up a term, because it is expected that that aspect might pose a challenge to translators in some contexts, especially when it comes to longer terms. If there are inconsistencies, maybe they arose because the intended equivalent made the sentence less readable and understandable, or the translator did not recognize the whole phrase as a term and translated it only partially consistently.

Initially, a fourth hypothesis was to be tested that was meant to test if terminological consistency was related to the part of speech a term belongs to. However, upon analysis of a random sample, it was observed that only one term was a verb, one term was an adjective, and the rest were all either nouns or noun phrases. This hypothesis was then dismissed, as the data to test it was insufficient.

## **5. Methodology<sup>3</sup>**

### **5.1. Corpus compilation**

Since this is a corpus-based study, the first step regarding data collection was finding or compiling a relevant corpus. When it comes to EU legislation, there are several pre-existing corpora that were considered first. The most recent public corpus is *EUR-Lex 2/2016 parallel*,

---

<sup>3</sup> I would like to thank prof. Stanojević for his invaluable advice and guidance on the methodological approach, particularly on corpus tools and analysis. I would also like to extend my gratitude to prof. Tonković for her helpful advice regarding the statistical analysis of the data.

a parallel corpus with multilingual subcorpora in all official languages of the European Union. It can be accessed through Sketch Engine, and it compiles European Union law and other public documents up until 2016 that are available in EUR-Lex, the online database of EU legal documents.<sup>4</sup> *The Digital Corpus of the European Parliament*<sup>5</sup> or *DCEP* is another publicly available corpus comprising documents, including legislative documents, published on the European Parliament's official website between 2001 and 2012. Although these corpora are the most recent corpora of EU legislation that are publicly accessible, their data, i.e. texts published up until 2012 and 2016, were deemed to be too dated for this study, especially since Croatia joined the EU in 2013. The final corpus that was considered is *DGT Translation Memory parallel corpus* in 24 official EU languages, including Croatian. It is also available on Sketch Engine; however, there is no clear indication of when it was created, which texts it compiles and from which time period. Therefore, a new, more relevant corpus had to be created. This was done using translation memories made public by The European Commission's Directorate-General for Translation on their website<sup>6</sup>. These TMs consist of parallel texts from the *acquis communautaire*, as well as some other texts, in all 24 official EU languages. The *acquis* comprises all treaties, regulations and directives adopted by the European Union, or in other words the EU legislation. According to the European Commission (n.d.), the texts were aligned in accordance with the DGT's segmentation rules and were pre-processed "to reduce the number of entries of low value for the translators (short sentences, long sentences, obvious mismatches, etc.)". There are several versions of the TMs depending on the year they were released. Since this study has a smaller scope, the corpus compiles only the most recently released TMs that consist of parallel texts from 2020. As mentioned, the TMs are available in all 24 official EU languages, so to acquire the ones for the English-Croatian language pair, a bilingual extraction was performed using the extraction tool TMXtract, made available by the DGT. Finally, TMs consisting of only English-Croatian parallel texts from 2020 were used to compile the corpus. This was done using Sketch Engine, a web-based corpus tool which offers features like corpus creation, automatic lemmatization and tagging for parts of speech, term extraction etc., which was useful for this study. The resulting corpus consists of two parallel corpora, with the English one compiling 7,170,658 words and the Croatian one 6,445,680 words.

---

<sup>4</sup><https://eur-lex.europa.eu/homepage.html>

<sup>5</sup>[https://joint-research-centre.ec.europa.eu/language-technology-resources/dcep-digital-corpus-european-parliament\\_en](https://joint-research-centre.ec.europa.eu/language-technology-resources/dcep-digital-corpus-european-parliament_en)

<sup>6</sup>[https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory\\_en#dgt-memory](https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory_en#dgt-memory)

## 5.2. Term list

The second step was compiling a term list which would contain terms to be examined in the study. EU legislation covers a wide range of domains and areas of human activity, which has translated into the EU developing very diverse terminology. To improve terminology management, the EU has, as mentioned, created IATE, its own term base which is available as an online tool, and can also be downloaded in multiple formats and all official EU languages. On IATE there is a filtering option that lists all 22 domains, or fields of knowledge, the EU terminology can be categorized into, ranging from law and economics to industry and energy. Since this is a small-scale study, the scope was reduced to analysing only trade-related EU terminology. Trade is one of the main domains of EU legislation, so it is the subject of many legislative documents and therefore has a well-developed and comprehensive terminology, which is why it was chosen for this study. It should be noted though that the domain variable could also be an interesting point in further research regarding terminological consistency, since the nature of the domain and the way its terminology is dealt with could influence its consistency in translation. However, this correlation will not be explored in this study.

As mentioned, terms in IATE are categorized by domain, which made it easier to compile a list of only trade-related terms for the analysis. In general, IATE comprises the majority of English terms used in EU documents and provides a number of translation equivalents for them in official languages. It is concept-based, meaning that each entry should correspond to a single concept. Furthermore, since its main purpose is to “facilitate multilingual drafting and translating of EU legal texts” (Stefaniak, 2017, p. 111), it does not only contain terminology, but also proper names, titles of documents and agreements, abbreviations and a number of phrases that could not necessarily be considered terms, but that occur often in EU texts and should be uniformly translated. However, IATE still has a number of downfalls, as pointed out by Stefaniak (2017) and Bratanić and Lončar (2016), such as the varying quality of entries, with some containing little to no information or many entries having low reliability as evaluated by IATE's own reliability system. The Croatian terminological network on IATE can especially be lacking at times, since there remains a significant number of English terms that do not have a listed Croatian equivalent. Nonetheless, this does not affect the results of this study since it is concerned only with the translation of terms in the actual texts and does not use IATE's Croatian term base. This is also why the term list is monolingual, i.e. contains only terms in English.

The base of the term list was created by downloading the English IATE term base for the trade domain in .xsl format. That list was then manually edited so that it is accessible and ready for further processing, which included separating entries that were initially in the same cell in Excel, so that each respective entry can be properly recognized by Excel, as well as deleting terms that are abbreviations and Latin terms. Those entries are not of interest for this study since they are always translated in the same way, so there is little to no room for possible inconsistencies apart from translator's sheer lack of attention. Since the list comprised thousands of terms which cannot all be guaranteed to be found in the study's corpus, the next step was to do an automatic term extraction from the English part of the parallel corpus in Sketch Engine. The extraction resulted in a list of over 150,000 single- and multi-word units that Sketch Engine recognized as terms and as expected, there were, as Bowker (2015) calls it "instances of both noise (non-pertinent items identified) and silence (relevant terms missed)" (p. 310). For this reason, she emphasizes the importance of manually editing the list, usually done by a domain expert (Bowker, 2015). In this study, the IATE term list was used as reference for the validity of the recognized terms instead. The automatically generated list of terms was compared to the IATE term list in Excel to determine which of the IATE terms are present in the corpus. The final term list consisted of 909 English terms (see Appendix A), and then a random sample of one hundred terms was taken using Excel. The random sample was supposed to provide a more even distribution of terms with different frequencies and structures, and possible inconsistencies.

Due to IATE's aforementioned varying reliability and quality of listed entries, there are limitations to this study, as the validity of certain terms that were analysed can be questioned. The reliability values on IATE are reflected through a stars rating system, with entries having 3 or 4 stars being manually verified and considered reliable and very reliable, respectively, and 1 or 2 stars indicating unverified and low reliability. Some terms also have no listed definition or other pragmatic information. That might weaken their reliability because, as explained in the *Key terms* section, a term's meaning, or its connection to the concept it describes is established through a definition (Sager, 1990, p. 21). The lack of definition can also lead to different interpretations of the term's reference, especially if observed in the context of different national legal systems (Ferrari, 2010, as cited in Pozzo, 2020). It is also interesting to note that many terms are listed in IATE as separate entries under different domains, and some of them have a listed definition under only one of the entries. For example, *replacement certificate* has five separate entries, one under finance, one under

international relations and three under trade, with two of them having low reliability value, and none of them having a definition. Moreover, some terms cannot be found on IATE's online search tool and only exist in the downloadable term base. Out of 100 terms analysed in this study, 22 terms have a satisfactory reliability value, but no listed definition, six terms have an unsatisfactory reliability value, and one is not listed on the online database. The remaining 71 terms fit all the reliability criteria, having both a listed definition and a high reliability score.

### **5.3. Analysis**

The process of gathering data for the analysis was corpus-based. In other terms, it involved searching for the occurrences of the terms in the corpus and then recording the relevant data, which included the frequency of a term in the corpus, the number of words a term consists of, and its translation equivalents present in the corpus, i.e. its terminological variants in the target language. When it comes to frequency, two types of frequency could have been used: a term's relative frequency in the corpus, or its absolute frequency. Relative frequency shows the relation between the number of occurrences of a term and the total number of tokens, or words, in the corpus and it is usually used to compare frequencies between corpora of different sizes. Absolute frequency is just the number of individual occurrences, or hits, in the corpus, which is why it is also referred to as raw frequency. For example, the relative frequency of *trade committee* is 0.001878%, or 18.78 per million tokens, while its absolute frequency is 165. This study observes and uses only the absolute frequency for the following reasons. Firstly, since the study is conducted on only one corpus, its size is not an important factor that would have had to be taken into consideration if this were a multi-corpus study, in which case the relative frequency would have been a better representation of a term's frequency. Secondly, not all hits in the corpus contained only the relevant term. In other words, the absolute frequency of certain terms was lower than indicated in the corpus, which meant the relative frequency would also have to be recalculated. This calculation was essentially much harder to do in comparison to the simple manual adjustment of the absolute frequency. Finally, absolute frequency was also needed to calculate the consistency index, or the HHI, so it was decided that only the absolute frequency of terms and its correlation to terminological consistency would be analysed in this study.

As mentioned, the corpus search doesn't always yield results containing only relevant term occurrences. Therefore, a set of criteria for the exclusion of certain results was laid down. New terms are often created by principle of recursion, i.e. by taking established terms and



combining them into new phrases with different meanings, as exemplified by terms *trade* → *trade policy* → *Trade Policy Committee*. Because of this phenomenon, when searching for certain terms in the corpus, especially single-word ones, it generates results that essentially contain a different, expanded term. In the context of the EU legislation, that also frequently happens with official documents or agreements that contain terms in their names, and are considered separate terms themselves. All of these instances had to be observed and taken note of with respect to the absolute frequency of terms, since any instances where the resulting example from the corpus contained what could be considered a separate term had to be eliminated from the analysis. For example, 43 instances of the term *international trade* were eliminated due to the term being a part of other terms like *international trade rules*, or *Convention on International Trade in Endangered Species of Wild Fauna and Flora*. The other exclusion criteria concerned the translation of the term, namely the instances where the term was transposed, i.e. replaced with a different word class, or even excluded from the translation. The exclusion criterion was especially important because of the calculation of the HHI. The HHI for each variant is a ratio of its absolute frequency and the absolute frequency of all terminological variants found in the translation. If the examples where the term was excluded in translation were counted towards a term's absolute frequency, the ratio would be skewed, and the HHI would, therefore, be inaccurate. In addition to data about overall frequency of the term, the number of its terminological variants was recorded, as well as the variants' form and frequency. For the purpose of testing the third hypothesis, each term was analysed on the level of structure, i.e. how many words it is made up of, and whether it is a single- or a multi-word unit. As previously discussed, there was an intention to analyse the correlation between the type of phrase or part of speech a term belongs to, and its consistency, but the random sample contained only one adjectival term and one verbal term, with the rest of them being nouns or noun phrases. That number was insufficient to draw any conclusions about the correlation, so they were replaced by the next two random sample terms that were nouns or noun phrases. Furthermore, that eliminated the part of speech as a possible intervening variable.

The last step of the analysis was to calculate the consistency index for each of the terms. As explained in the *Key terms* section, Herfindahl Hirschman Index or HHI is a commonly used measure of market concentration in economics. However, it was introduced into terminological research by Itagaki et al. (2007) as a way to measure and automatically validate terminological consistency which was up until then evaluated only qualitatively.

Their study was focused on terminology in localized materials, like manuals, and in training of example-based and statistical MT systems, since the training can be hindered by terminological inconsistency. The authors adapted the HHI formula to fit the context of translation and terminology which resulted in the following formula:

$$C_t = \frac{\sum_{n=1}^p \sum_{i=1}^n \left(\frac{f}{k} \times 100\right)_i^2}{p}$$

$C$  is the consistency index for a specific term ( $t$ ),  $p$  is the number of texts that contain the term,  $f$  is the absolute frequency of a particular translation variant, and  $k$  is the total number of occurrences of the term, or a sum of absolute frequencies of all variants, within a text or corpus (Itagaki et al. 2007, p. 5). This formula, especially the  $p$  variable, was pertinent to their methodology and aims, because the study analysed terminological consistency across multiple groups of texts belonging to different products. When Gašpar (2013) applied this method to assess the terminological consistency of translated terms in a Croatian-English parallel corpus of legislative texts, she further adapted the formula by removing the  $p$  variable to calculate the HHI score for individual terms. The adapted formula was as follows:

$$C_t = \sum_{n=1}^n \left(\frac{f}{k} \times 100\right)_i^2$$

It can be said that the consistency index of a particular term is the sum of the consistency indexes of all its respective variants found in a text or corpus. The “frequency share” for each variant is calculated as a ratio of its absolute frequency and the total occurrence of all variants of that term. Gašpar et al. (2022) applied this formula again on an expanded range of corpora, this time including also Latin-English and Latin-Croatian versions of the Code of Canon Law (1983), and the English and Croatian versions of the EU legislation (2013-). Both studies confirmed that HHI as a measure of terminological consistency can be successfully applied to Croatian-English and English-Croatian legal translations. Therefore, the aforementioned formula to calculate the HHI, and consequently assess terminological consistency, will be used in this study as well. Table 1 shows an example of an HHI calculation. The final values of the HHI, seen in the far right column of Table 1, were normalized to a range of 0-100, with 0 marking complete terminological inconsistency, i.e. the term being translated differently in every instance, and 100 marking complete terminological consistency, i.e. the term being translated using the same terminological variant in every instance.

SOURCE TERM	VARIANTS	FREQUENCY	$\left(\frac{f}{k} \times 100\right)$	$\sum_{n=1}^n \left(\frac{f}{k} \times 100\right)_i^2$
price suppression	sprečavanje rasta cijena	15	4253.31	4593.56
	smanjenje cijena	2	7.61	
	pritisak na cijene	3	170.13	
	pad cijena	1	18.90	
	sniženje cijena	2	7.61	

Table 1: HHI calculation for the term *price suppression*

Although the number of terminological variants indicates the presence of terminological inconsistency, it is actually the ratio of these variants' frequency that affects the index the most. For example *subtotal* had two terminological variants, *međuzbroj* and *ukupno*, and its HHI was 71.18, while *import price* had four terminological variants *uvozna cijena*, *cijena uvoza*, *izvozna cijena* and *obujam uvoza* and its HHI was higher, at 87.16. This is because the two variants for *subtotal* were relatively evenly distributed, while *uvozna cijena* was the clearly dominant variant for *import price*. It follows that the consistency index evaluates terms with multiple variants out of which one is dominant as “more consistent” than terms with fewer variants that are equally distributed.

Finally, due to the distribution of the HHI scores being skewed, another value was added to each term in the spreadsheet for statistical analysis purposes. Their HHI scores were ranked on an ordinal scale as shown in Table 2, and each term got assigned a number from 1 to 5. For the final spreadsheet with all data for all one hundred terms, see Appendix B.

HHI SCORE	ORDINAL VALUE
0.00-19.99	1
20.00-39.99	2
40.00-59.99	3
60.00-79.99	4
80.00-100.00	5

Table 2: HHI scores with assigned ordinal values

Statistical analysis and correlation tests were performed using JASP (version 0.17.2.0), an open-source program for statistical analysis. Owing to the fact that the distributions of all three variables (frequency, number of words, HHI score) were skewed, as seen in Figures 1, 2 and 3 in the *Results* section, Spearman's rank correlation, a non-parametric correlation test,

was employed to test the hypotheses about the correlation between consistency, and frequency and term structure, respectively.

## 6. Results

The following sections present the results of the statistical analysis. The first section gives an overview of descriptive results for each of the variables used in the correlation tests. The second and third section report the results of correlation tests regarding terminological consistency and frequency, and term structure, respectively.

### 6.1. Descriptive results

#### 6.1.1. Frequency

First variable that was analysed was frequency. As previously discussed, only absolute frequencies of terms were recorded. The frequency distribution was right-skewed (see Figure 1). The median frequency was 11.5 (IQR = 5-29). The term with the highest frequency was *consignment* with 1645 occurrences, followed by *trade* (uncountable noun) with 1304 occurrences. *Middle-value contract* had the lowest frequency, occurring only twice in the corpus.

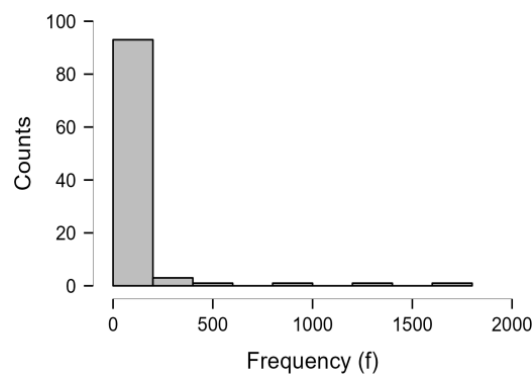


Figure 1: Frequency distribution

#### 6.1.2. Number of words

The distribution of the number of words a term consists of was slightly right-skewed as well, as shown in Figure 2. Out of 100 terms, four terms were single-word units and the remaining 96 were multi-word units; sixteen terms consisted of three words; eight of four words and two of five words. The most terms, 70 of them, consisted of two words. The median was 2 (IQR = 2-3).

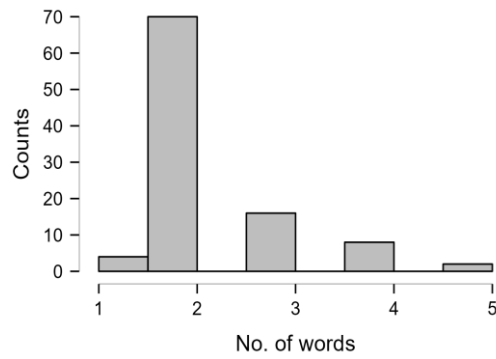


Figure 2: Number of words distribution

### 6.1.3. Herfindahl Hirschman Index score

The distribution of the HHI scores was left-skewed (see Figure 3), with 51 terms having optimal scores of 100. This means that the overall consistency was fairly high, with the median score being 100.00 (IQR = 63.31-100.00). The term with the lowest consistency of 33.33 was *specific contract*, followed by *corporate entity* ( $C_t = 33.56$ ) and *supply contract* ( $C_t = 42.15$ ). However, the terms with the most terminological variants (5) were *price suppression* ( $C_t = 45.94$ ) and *trade* ( $C_t = 96.23$ ).

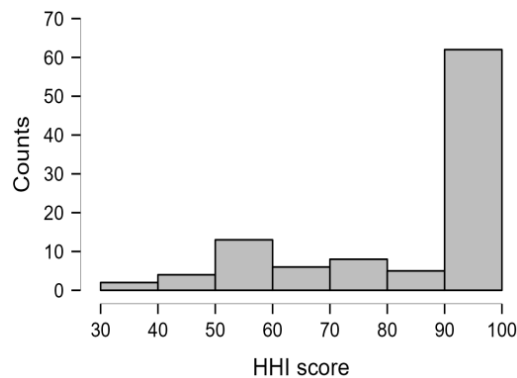
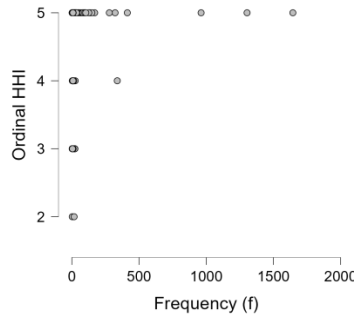


Figure 3: HHI score distribution

## 6.2. Correlation between frequency and terminological consistency

As explained in the *Methodology* section, due to the skewness of distribution across all variables, Spearman's rank correlation coefficient, or Spearman's rho ( $\rho$ ) was computed to assess the correlations. Additionally, instead of the HHI score variable, the aforementioned ordinal HHI variable was used as a measure of terminological consistency for both correlations. The results of the correlation test, as seen in Figure 4, show that there is a weak positive correlation between frequency and ordinal HHI score,  $\rho(98) = .32$ ,  $p < .001$ . This means that more frequent terms tend to have a higher HHI score, or rather are translated more consistently. The effect size, as measured by Spearman's rho, indicates a medium effect (Goss-Sampson, 2022, 41).



## **7. Discussion**

### **7.1. Terminological consistency in Croatian translations**

The main aim of this study was to determine if there are terminological inconsistencies in the Croatian translations of the EU legislation. Terminological consistency is considered one of the most essential features of specialized texts, and in turn translations. It enhances readability and information transfer, as well as reduces the possibility of misunderstanding and ambiguity (Gašpar et al, 2022, p. 2). In the context of the EU, consistency is even more insisted upon because of the implications of EU legislation, as well as the legal dimension it exists in. Since EU legislation is implemented into the law of every Member State, terminological consistency ensures there are no ambiguities or difficulties in its interpretation, and consequently enhances the harmonisation of laws between Member States. Legal uncertainty can have serious consequences at both national and EU level, resulting in misinterpretation of rights and obligations, or even legal disputes (Stefaniak, 2017). Due to those circumstances, terminology work in drafting and translation processes has been brought to the forefront of the EU's language services tasks. Nonetheless, the results of the analysis confirm that there are terminological inconsistencies in Croatian translations, i.e. that translators rendered certain terms using more than one translation equivalent, or terminological variant. However, even though inconsistency is present, it can be said that the overall consistency of the analysed terms in this corpus was relatively high, as indicated by 51 of the 100 examined terms having the optimal HHI score of 100.

### **7.2. Correlation between frequency and terminological consistency**

Furthermore, since inconsistencies were detected, the next aim of this study was to test if the (in) consistency might correlate with a term's frequency, or more precisely, if higher frequency positively correlated to higher consistency. The sample consisted of many more infrequent terms than frequent ones, resulting in a skewed distribution. This, however, was expected due to the skewness of frequency being "a design feature of language" (Taylor, 2012, p. 180). As Taylor (2012) explains, this means that normally "a small number of very common words make up the bulk of a text, a fair number of moderately frequent words constitute somewhat smaller proportion, while a very large number of infrequent words account for only a tiny amount of a text" (p. 156). In view of this, the skewed distribution was accounted for in the analysis by using Spearman's non-parametric correlation test. The results of the test found a weak positive correlation between frequency and terminological consistency, confirming the second hypothesis. This means that the more frequent the term is,

the more consistently it tends to be translated. At the same time, it is important to remember that consistency, as measured by HHI, does not account as much for the number of variants, as for their ratio across all occurrences. As a result, a term like *Union producer* that has four terminological variants, still has a high HHI,  $Ct = 95.71$ , while a term with two variants like *special fiscal territory* has a low HHI,  $Ct = 50.78$ , because the variants are relatively equally distributed. These findings could be partially explained by the process of EU translation in which translation memories play a central role. As previously explained, the EU has a central translation memory called *Euramis*, which is automatically integrated into the translation process. It offers the translator already translated segments stored in the memory that are similar to the ones being translated, and the translator can choose to copy them, retain them with alterations or ignore them, depending on how similar the retrieved and new segments are. The TMs can also be used to look up words or phrases in older documents to get an overview of the context they occur in. As follows, the highly-frequent terms are more likely to occur in a greater number of segments in the TM, and therefore more likely to be included in the segments recommended to the translator, as well as be found when the translator manually looks them up. On the other hand, less frequent terms might not be as present in the TM. In such cases the translator might have to use other resources, like IATE which can, as mentioned, be lacking, to find a term's right translation equivalent, and if none can be found, they might have to create a new one. However, because of their dependence on older translations, TMs can both enhance and reduce terminological consistency. If the stored TMs consist of translations with multiple terminological variants for one term, the TM system might suggest two segments containing different variants to two different translators, depending on the context the term occurs in. If the translator does not check which one is the established term translation in their language but rather automatically copies the unsuitable term, its frequency in the TM might increase, perpetuating the process. It is possible that such line of events is the reason behind certain terms like *subtotal* or *special fiscal territory* having a rather equal distribution of variants. For this reason, all language departments in the DGT have so-called sentence managers, whose main task is to update translation memories, as well as update IATE in cooperation with the terminologists (European Commission, 2012, p. 25).

Additionally, another aspect that comes into play in this correlation is that translators are more likely to learn, and recognize more frequent terms, and consequently, know how they should be translated. Even if they do not know the correct translation equivalent, they would be aware that it is a term they should look up to retain consistency. On the other hand, less



frequent terms might go unrecognized and therefore be translated differently depending on the context. For instance, *specific contract* had a low absolute frequency of three in this corpus,  $f = 3$ , and was translated differently each time. According to IATE, the term denotes a “contract specifying details of a particular task based on the previously signed framework contract or agreement, dynamic purchasing system or qualification system”, and its Croatian equivalent is *posebni ugovor*. However, two translations contained variants *pojedini* and *pojedinačni ugovor*, pointing to the fact that the translators probably failed to recognize the phrase as a term, thus did not look it up, and translated it using a phrase with a more general meaning which seemingly fit the context, but did not retain the true meaning of the term.

As previously discussed, one of the study’s limitations regarding term extraction was possible unreliability of IATE’s entries. For this reason, the correlation test was also performed on data which excluded possible unreliable terms from the study’s term list. The criteria of unreliability included the term not existing in IATE’s online base, the term’s reliability score on IATE being low, and the term having no listed definition on IATE. The correlation coefficient for data excluding the seven terms with a low reliability score didn’t change significantly. However, when the test was performed on data excluding additional 22 terms with no listed definition, the correlation coefficient increased significantly, indicating a stronger or more precisely moderate correlation between frequency and terminological consistency. It is hard to say if the coefficient’s value changed due to the smaller sample size, or if IATE’s unreliability and its possible effect on the translation process presented itself as an intervening variable. Nonetheless, the results of the study indicate that a correlation between terminological consistency and frequency does exist, but they should be taken as preliminary due to certain methodological limitations and the lack of other analogous research in the area.

### **7.3. Correlation between structure and terminological consistency**

Conversely, the same cannot be said for the correlation between terminological consistency and the number of words a term consists of. The results of the correlation test showed no correlation between those two variables, rejecting the third hypothesis. The basis of the hypothesis was the expectation that longer terms might be translated differently in certain contexts where the readability would decrease with the use of their intended translation equivalent. It was also expected that, similarly to infrequent terms, translators might not recognize the whole phrase as a term in case of longer terms which would result in partially consistent translation. No such cases were recorded, however, which could indicate that

consistency is indeed given precedence over possible decreased readability. Moreover, the length of a term might not even be a particularly relevant factor in the translation process, especially since the distribution of terms indicates that most terms consist of two words (see Figure 2).

#### **7.4. Limitations and relevance**

The main limitations of this study, as pointed out in the *Methodology* section, come from the use of IATE as a reference for the validity of extracted terms from the corpus. The term list might have been more credible had it been evaluated by a domain expert. However, the applied methodological approach does give an insight into the current state and practical use of IATE, which presents itself as an additional contribution. Other limitations stem from the use of nonparametric correlation tests which are less precise than parametric tests and need a larger sample size to show sufficient results (Eddington, 2015, p. 37). The results should consequently be taken as preliminary, since the scope of the study was rather small, focusing only on trade-related terminology. It is also hard to account for the intervening variable of the translator's lack of attention or skill which might have influenced the translation of some of the entries. Furthermore, despite the general discourse about terminological (in)consistency in the EU, there have not been many studies exploring the issue using a more quantitative or statistical approach, so the findings cannot be compared or interpreted in a broader context of the research area. However, they do provide an insight into the state of terminological consistency in recent Croatian translations of EU legislation. Their main relevance is therefore in creating a space for further research on terminological consistency in EU translation, not only regarding the Croatian translations, but other official languages as well. Finally, the study expands on the implementation of the Herfindahl Hirschman Index for measuring terminological consistency and supports its use on the English-Croatian language pair, as well as in EU legislation.

### **8. Conclusion**

Due to its impact on clarity and interpretation, terminology is one of the crucial aspects of translation in the EU. Its harmonisation presents itself as a challenge both because of the interplay between the EU legal system and Member States' diverse legal systems, and because of the fast-paced environment in which terminological work and translation are done. Integration of IATE, the EU's multilingual terminology database, was supposed to facilitate

this process and help ensure terminological consistency across all documents and languages. The aim of this study was therefore to check if there are still inconsistencies in the use of terminology in Croatian translations of EU legislation, and to determine if there is a correlation between the frequency and structure of a term and its terminological consistency where inconsistencies were found. The consistency of 100 trade-related terms taken from IATE was analysed in a corpus compiling English-Croatian EU translation memories from 2020. It was measured using the Herfindahl Hirschman Index, a method innovated by Itagaki et al. (2007). The results have found terminological inconsistency to be present, with the overall consistency still being relatively high. Furthermore, statistical analysis reported a weak positive correlation between consistency and frequency, and no correlation between consistency and structure. However, due to the limitations of the study, notably the small dataset and lack of analogous research, the results are to be taken as preliminary.

The study's main contribution is thus to provide an insight into the current state of terminological consistency in Croatian translations of EU legislation and encourage further discussions and research on terminological consistency in EU translation, not only in Croatian translations, but also regarding other official languages. Additionally, it provides an overview of several EU language resources, such as IATE and DGT's TMs, and illustrates how useful they can be for linguistic and terminological research. Future studies could focus on a diachronic analysis of EU translations, to see if there has been any change in quality and consistency over time. Likewise, the domain-specificity of terms in the context of consistency could also be an interesting variable to assess. The correlation between frequency and terminological consistency and the implications that translation memories have for it could also be further analysed, as these aspects might have practical outcomes.

Lastly, this study continued the work of Itagaki et al. (2007), Gašpar (2013) and Gašpar et al. (2022) by implementing the Herfindahl Hirschman Index to measure terminological consistency. The method was proven to be a suitable tool for quantitative terminological research, with the results supporting its use on the English-Croatian language pair and in EU legislation.

## References

- Bajčić, M., & Stepanić, M. (2011). Nedosljednost u prevođenju pojmova iz prava tržišnog natjecanja Europske unije. In Maja Bratanić (Ed.), *Hrvatski na putu u EU* (pp. 133-151). Institut za hrvatski jezik i jezikoslovlje, Hrvatska sveučilišna naklada.
- Bratanić, M., & Lončar M. (2016). The Myth of EU Terminology Harmonization on National and EU Level. In Susan Šarčević (Ed.), *Language and Culture in EU Law* (pp. 207-218). Routledge.
- Bowker, L. (2015). Terminology and translation. In Kockaert, H., & Steurs, F. (Eds.), *Handbook of terminology.: Volume 1* (pp. 304-323). John Benjamins Publishing Company. <https://doi.org/10.1075/hot.1>
- Craig, P. (2012). *EU Administrative Law* (2nd ed.). Oxford Univ. Press. <https://doi.org/10.1093/acprof:oso/9780199568628.003.0018>
- Eddington, D. (2015). *Statistics for Linguists: A Step-by-Step Guide for Novices*. Cambridge Scholars Publishing.
- European Commission, Directorate-General for Translation. (2012). *Quantifying quality costs and the cost of poor quality in translation : quality efforts and the consequences of poor quality in the European Commission's Directorate-General for Translation*. Publications Office of the European Union. <https://data.europa.eu/doi/10.2782/44381>
- European Commission, Directorate-General for Translation, (2016). *Translation tools and workflow*, Publications Office of the European Union. <https://data.europa.eu/doi/10.2782/703257>
- European Commission, Directorate-General for Translation. (2023). *Translation in figures 2023* [Brochure]. Publications Office of the European Union. <https://data.europa.eu/doi/10.2782/438164>
- European Commission. (n.d.). *DGT-Translation Memory*. EU Science Hub. [https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory\\_en#dgt-memory](https://joint-research-centre.ec.europa.eu/language-technology-resources/dgt-translation-memory_en#dgt-memory) (Accessed 26/1/24)
- Fischer, M. (2010). Language (policy), translation and terminology in the European Union. In Marcel Thelen & Frieda Steurs (Eds.), *Terminology on everyday life: Terminology and Lexicography Research and Practice 13* (pp. 21-34). John Benjamins.

- Fischer, M. (2022). Horizontal and vertical terminology work in the context of EU translations. *Porta Lingual*, 7–15. <https://doi.org/10.48040/pl.2022.1.1>
- Gašpar, A. (2013). Računalno potpomognuta provjera terminološke dosljednosti prijevoda hrvatskog zakonodavstva na engleski jezik [Doctoral dissertation, Faculty of Humanities and Social Sciences]. Hrvatski nacionalni skupni katalog.
- Gašpar, A., Seljan, S., & Kučič, V. (2022). Measuring Terminology Consistency in Translated Corpora: Implementation of the Herfindahl Hirshman Index. *Information* 13(2), 43. <https://doi.org/10.3390/info13020043>
- Goss-Sampson, M. (2022). *Statistical analysis in JASP: a guide for students*. <https://jasp-stats.org/jasp-materials/statistical-analysis-in-jasp-a-students-guide-v16/>
- Itagaki, M., Aikawa, T., & He, X. (2007). Automatic validation of terminology translation consistency with statistical method. Machine Translation Summit. <https://aclanthology.org/2007.mtsummit-papers.36.pdf>
- Kockaert, H., & Steurs, F. (Eds.). (2015). *Handbook of terminology.: Volume 1*. John Benjamins Publishing Company. <https://doi.org/10.1075/hot.1>
- Pozzo, B. (2020). Looking for a consistent terminology in European contract law. *Lingue Culture Mediazioni – Languages Cultures Mediation (LCM Journal)*, 7(1), 103–126. <https://doi.org/10.7358/lcm-2020-001-pozz>
- Rogers, M. (2008). Consistency in Terminological Choice: Holy Grail or False Prophet? In Ingrid Simmons (Ed.), *SYNAPS 21: Festschrift for Magnar Brekke* (pp. 107-113). NNH.
- Sager, J. C. (1990). *Practical Course in Terminology Processing*. John Benjamins Publishing.
- Stefaniak, K. (2017). Terminology work in the European Commission: Ensuring high-quality translation in a multilingual environment. In Tomáš Svoboda, Lucja Biel & Krzysztof Łoboda (Eds.), *Quality aspects in institutional translation* (pp. 109–121). Language Science Press. DOI:10.5281/zenodo.1048192
- Taylor, J. R. (2012). *The Mental Corpus: How Language Is Represented in the Mind*. Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780199290802.001.0001>
- Temmerman, R. (2000). *Towards New Ways of Terminology Description: The sociocognitive approach*. John Benjamins Publishing Co. <https://doi.org/10.1075/tlrp.3>

## **Appendix A**

List of 909 IATE terms present in the corpus

The appendix is attached to the Master's thesis in the ODRAZ repository in the Digital Academic Archives and Repositories (DABAR).

## Appendix B

Table containing data for all 100 sampled terms used for the statistical analysis

TERM	FREQUENCY (f)	S1 or MW2	No. of words	No. of variants	VARIANTS	Variants f	HHI variants	HHI term	HHI term 100	HHI ordinal
capital fund	7	2	2	1	fond kapitala	7	10000	10000	100	5
product control number	13	2	3	1	kontrolni broj proizvoda	13	10000	10000	100	5
economic area	9	2	2	1	gospodarsko područje	9	10000	10000	100	5
consignment	1645	1	1	1	pošiljka	1645	10000	10000	100	5
customs territory	170	2	2	1	carinsko područje	170	10000	10000	100	5
net mass	28	2	2	2	neto masa	27	9298.47	9311.23	93.11	5
					neto težina	1	12.76			
capacity calculation	22	2	2	2	izračun kapaciteta	2	82.64	8347.1	83.47	5
					proračun kapaciteta	20	8264.46			
household washing machine	32	2	3	2	perilica	1	9.77	9394.54	93.95	5
					kućanska perilica rublja	31	9384.77			
works contract	8	2	2	2	ugovor o radovima	7	7656.25	7812.5	78.13	4
					ugovor radova	1	156.25			
identity of means of transport	3	2	5	1	identitet prijevoznog sredstva	3	10000	10000	100	5
certificate of inspection	18	2	3	1	potvrda o inspekciji	18	10000	10000	100	5
exclusive right	66	2	2	1	isključivo pravo	66	10000	10000	100	5
security exception	6	2	2	2	iznimka u vezi sa sigurnošću	2	1111.11	5555.55	55.55	3
					iznimka za sigurnost	4	4444.44			
joint proposal	13	2	2	1	zajednički prijedlog	13	10000	10000	100	5
provisional disclosure	50	2	2	2	privremena objava	49	9604	9608	96.08	5
					privremeno objavljivanje	1	4			
subsidised import	90	2	2	1	subvencionirani uvoz	90	10000	10000	100	5
international trade	78	2	2	1	međunarodna trgovina	78	10000	10000	100	5
arbitration	27	1	1	1	arbitraža	27	10000	10000	100	5
certificate of approval	5	2	3	3	potvrda o homologaciji	3	3600	4400	44	3
					certifikat o homologaciji	1	400			
					potvrda o odobrenju	1	400			
electronic transport document	6	2	3	1	elektronička prijevozna isprava	6	10000	10000	100	5
immediate packing	54	2	2	1	neposredno pakiranje	54	10000	10000	100	5
price of natural gas	5	2	4	1	cijena prirodnog plina	5	10000	10000	100	5
corporate income tax	20	2	3	2	porez na dobit	5	625	6250	62.5	4
					porez na dobit trgovačkih društava	15	5625			
military equipment	28	2	2	1	vojna oprema	28	10000	10000	100	5
request for a review	4	2	4	1	zahtjev za reviziju	4	10000	10000	100	5
actual value	11	2	2	2	stvarna vrijednost	9	6694.21	7024.79	70.25	4
					trenutna vrijednost	2	330.58			
economic operator	413	2	2	1	gospodarski subjekt	412	9951.63	9951.69	99.52	5
					subjekt u tržišnom gospodarstvu	1	0.06			
passive means of transport	6	2	4	1	pasivno prijevozno sredstvo	6	10000	10000	100	5
subtotal	338	1	1	2	međuzbroj	279	6813.57	7118.26	71.18	4
					ukupno	59	304.69			
unpaid balance	3	2	2	1	nepodmireni iznos	3	10000	10000	100	5
amortisation of intangible assets	3	2	4	1	amortizacija nematerijalne imovine	3	10000	10000	100	5
trade agreement	15	2	2	1	trgovinski sporazum	15	10000	10000	100	5
					oslobođenje od polaganja					
guarantee waiver	6	2	2	3	osiguranja	1	277.78	5000	50	3
					odricanje od osiguranja	4	4444.44			
					odricanje od jamstva	1	277.78			
manufacture of goods	21	2	3	3	proizvodnja robe	14	4444.44	5102.03	51.02	3
					proizvodnja proizvoda	5	566.89			
					povezivanje ili zamatanje robe	2	90.7			
cocoa butter	16	2	2	1	kakao maslac	16	10000	10000	100	5
domestic sale	134	2	2	2	domaća prodaja	133	9851.3	9851.86	98.52	5
					prodaja na domaćem tržištu	1	0.56			
					nacionalni portal za gospodarske					
national trader portal	32	2	3	1	subjekte	32	10000	10000	100	5
import price	322	2	2	4	uvozna cijena	300	8680.22	8715.53	87.16	5
					cijena uvoza	19	34.82			
					izvozna cijena	1	0.1			
					obujam uvoza	2	0.39			
country of preferential origin	6	2	4	1	zemlja povlaštenog podrijetla	6	10000	10000	100	5
specific contract	3	2	2	3	pojedinačni ugovor	1	1111.11	3333.33	33.33	2
					posebni ugovor	1	1111.11			
					pojedini ugovor	1	1111.11			
malt extract	28	2	2	1	sladni ekstrakt	28	10000	10000	100	5
price suppression	23	2	2	5	sprečavanje rasta cijena	15	4253.31	4593.56	45.94	3
					smanjenje cijena	2	75.61			
					pritisak na cijene	3	170.13			
					pad cijena	1	18.9			
					sniženje cijena	2	75.61			
trade committee	149	2	2	1	Odbor za trgovinu	149	10000	10000	100	5
express request	3	2	2	2	izričit zahtjev	2	4444.44	5555.55	55.55	3
					izričit poziv	1	1111.11			

TERM	FREQUENCY (f)	S1 or MW2	No. of words	No. of variants	VARIANTS	Variants f	HHI variants	HHI term	HHI term 100	HHI ordinal
approved exporter	99	2	2	2	ovlaštenu izvoznik	98	9799	9800.02	98	5
					odobreni izvoznik	1	1.02			
customs decisions system	4	2	3	2	sustav Carinske odluke	3	5625	6250	62.5	4
					sustav za carinske odluke	1	625			
diversion notification	3	2	2	1	obavijest o preusmjeravanju	3	10000	10000	100	5
special account	3	2	2	1	poseban račun	3	10000	10000	100	5
arbitrary discrimination	28	2	2	1	proizvoljna diskriminacija	28	10000	10000	100	5
business register	21	2	2	1	poslovni registar	21	10000	10000	100	5
information package	6	2	2	2	opisna dokumentacija	5	6944.44	7222.22	72.22	4
					informatijski paket	1	277.78			
conventional duty	3	2	2	2	konvencionalna carina	2	4444.44	5555.55	55.56	3
					uobičajena carina	1	1111.11			
trade (uncountable)	1304	1	1	5	trgovina	1279	9620.24	9622.87	96.23	5
					trgovanje	21	2.59			
					trgovačka djelatnost	2	0.02			
					trgovinska razmjena	1	0.01			
					promet	1	0.01			
award procedure	47	2	2	1	postupak dodjele	47	10000	10000	100	5
medicinal substance	3	2	2	2	medicinska tvar	1	1111.11	5555.55	55.56	3
					ljekovita tvar	2	4444.44			
common customs tariff	126	2	3	2	Zajednička carinska tarifa	119	8919.75	8950.61	89.51	5
					zajednička carinska tarifa	7	30.86			
special fiscal territory	16	2	3	2	posebno fiskalno područje	9	3164.06	5078.12	50.78	3
					posebno porezno područje	7	1914.06			
purchase option	7	2	2	3	pravo kupnje	5	5102.04	5510.2	55.1	3
					pravo na kupnju	1	204.08			
					opcija kupnje	1	204.08			
nature of the goods	9	2	4	2	priroda robe	8	7901.23	8024.69	80.25	5
					narav robe	1	123.46			
investment aid	17	2	2	1	potpora za ulaganje	17	10000	10000	100	5
competent customs office	8	2	3	2	nadležan carinski ured	7	7656.25	7812.5	78.13	4
					ovlašten carinski ured	1	156.25			
wheat gluten	14	2	2	1	pšenični gluten	14	10000	10000	100	5
Union producer	962	2	2	4	proizvođač iz Unije	941	9568.17	9570.88	95.71	5
					proizvođač u Uniji	5	0.27			
					proizvođač Unije	15	2.43			
					Unijin proizvođač	1	0.01			
occasional service	4	2	2	2	povremeni prijevoz	3	5625	6250	62.5	4
					povremena usluga	1	625			
supply contract	22	2	2	4	ugovor o nabavi	3	185.95	4214.88	42.15	3
					ugovor o nabavi robe	5	516.53			
					ugovor o isporuci	1	20.66			
					ugovor o opskrbi	13	3491.74			
debt management	7	2	2	1	upravljanje dugom	7	10000	10000	100	5
corporate entity	17	2	2	4	korporativni subjekt	4	553.63	3356.39	33.56	2
					poduzeće	8	2214.53			
					gospodarski subjekt	1	34.6			
					pravna osoba	4	553.63			
motor-driven fan	8	2	2	2	motorni ventilator	2	625	6250	62.5	4
					ventilator na motorni pogon	6	5625			
clearing house	9	2	2	2	klirinška kuća	5	3086.12	5061.43	50.61	3
					mehanizam za razmjenu	4	1975.31			
marketable quality	5	2	2	1	tržišna kvaliteta	5	10000	10000	100	5
sales agreement	5	2	2	2	ugovor o prodaji	2	1600	5200	52	3
					sporazum o prodaji	3	3600			
special stamp	3	2	2	1	posebni pečat	3	10000	10000	100	5
sole distributor	4	2	2	1	jedini distributer	4	10000	10000	100	5
transaction value	9	2	2	2	vrijednost transakcije	3	1111.11	5555.55	55.56	3
					transakcijska vrijednost	6	4444.44			
domestic price	83	2	2	2	domaća cijena	82	9760.49	9761.94	97.62	5
					cijena na domaćem tržištu	1	1.45			
domestic procedure	6	2	2	1	domaći postupak	6	10000	10000	100	5
late payment	25	2	2	2	zakašnjelo plaćanje	22	7744	7888	78.88	4
					kašnjenje u plaćanju	3	144			
non-automatic import licensin	3	2	3	1	neautomatski postupci					
distinct market	3	2	2	2	izdavanja uvoznih dozvola	3	10000	10000	100	5
					zasebno tržište	2	4444.44	5555.55	55.55	3
					različito tržište	1	1111.11			
underselling margin	38	2	2	2	marža sniženja ciljnih cijena	37	9480.61	9487.54	94.88	5
					marža nelojalnog sniženja ciljnih cijena	1	6.93			



TERM	FREQUENCY (f)	S1 or MW2	No. of words	No. of variants	VARIANTS	Variants f	HHI variants	HHI term	HHI term 100	HHI ordinal
market surveillance	96	2	2	4	nadzor tržišta	87	8212.89	8250.88	82.51	5
					nadzor nad tržištem	5	27.13			
					nadziranje tržišta	3	9.77			
					tržišni nadzor	1	1.09			
memo item	17	2	2	1	bilješka	17	10000	10000	100	5
award of the contract	4	2	4	2	odjela ugovora	3	5625	6250	62.5	4
					odjeljivanje ugovora	1	625			
autonomous tariff suspension	11	2	3	2	autonomna tarifna suspenzija	9	6694.21	7024.79	70.25	4
					autonomna carinska suspenzija	2	330.58			
originating product	104	2	2	1	proizvod s podrijetlom	104	10000	10000	100	5
country of last known destination	3	2	5	2	zemlja posljednjeg poznatog odr	2	4444.44	5555.55	55.55	3
					zemlja zadnjeg poznatog odrediš	1	1111.11			
middle-value contract	2	2	2	1	ugovor srednje vrijednosti	1	10000	10000	100	5
free movement	279	2	2	3	slobodno kretanje	273	9574.52	9576.84	95.75	5
					slobodni protok	3	1.16			
					sloboda kretanja	3	1.16			
central government authority	5	2	3	1	tijelo središnje državne uprave	5	10000	10000	100	5
export refund	28	2	2	1	izvozna subvencija	28	10000	10000	100	5
oleaginous fruit	14	2	2	2	uljani plod	12	7346.94	7551.02	75.51	4
					uljni plod	2	204.08			
replacement certificate	5	2	2	2	zamjenska potvrda	2	1600	5200	52	3
					zamjenski certifikat	3	3600			
auto loan	12	2	2	1	kredit za kupnju automobila	12	10000	10000	100	5
preferential origin of goods	4	2	4	1	povlašteno podrijetlo robe	4	10000	10000	100	5
level of trade	35	2	3	1	razina trgovine	35	10000	10000	100	5
border crossing	8	2	2	3	prelazak granice	6	5626	6250	62.5	4
					prijelaz granice	2	625			
fallback procedure	3	2	2	1	rezervni postupak	3	10000	10000	100	5
harmful substance	10	2	2	1	štetna tvar	10	10000	10000	100	5
nominal mass	7	2	2	1	nazivna masa	7	10000	10000	100	5
exemption certificate	8	2	2	1	potvrda o oslobođenju	8	10000	10000	100	5

## Appendix C

Table containing seven terms with low reliability score on IATE (Row 1) and 22 terms with no listed definition on IATE (Rows 2&3)

LOW RELIABILITY SCORE	NO LISTED DEFINITION	NO LISTED DEFINITION
manufacture of goods	capital fund	level of trade
express request	identity of means of transport	nominal mass
occasional service	certificate of approval	arbitrary discrimination
motor-driven fan	immediate packing	nature of the goods
special stamp	price of natural gas	competent customs office
domestic procedure	request for a review	Union producer
central government authority	unpaid balance	marketable quality
	domestic sale	sole distributor
	country of preferential origin	domestic price
	diversion notification	distinct market
	special account	preferential origin of goods