# PROCESSING TERMINOLOGY FOR TERMINOGRAPHY OUTPUTS

## JOZEF ŠTEFČÍK

*University of Economics in Bratislava, Faculty of Applied Languages, Dolnozemská cesta 1, 852 35 Bratislava, Slovensko*
*Email: jozef.stefcik@euba.sk*

**Abstract**

The paper presents the terminology work and terminography's concepts and principles for creating terminology glossaries, datasets or termbases used as parts of more complex terminology databases. It introduces the structure of a terminography record, points out preferences in creating the structure of multilingual documents, and briefly introduces frequently used online terminography platforms. Also, it demonstrates what tools and corpora can be applied to terminology processing. Furthermore, the paper intends to reflect on setting up terminology glossaries and what elementary mistakes should be avoided when preparing them.

**Keywords**

terminology work, terminography, terminology database, glossary, dictionary, online platform

## Introduction

Terminography is a part of terminology that is primarily concerned with describing terms or concepts and creating terminology products such as various dictionaries, databases, glossaries, thesauri, etc. Like Terminology and Terminology Work, terminography is a team and interdisciplinary activity based on multiple technical, formal, and procedural recommendations that are subject to international harmonization. It describes professional language in multilingual communication and creates specialized glossaries with specialized lexis (Jurčacková, 2002). Terminography targets the attention mainly of the professional public, using a large number of electronic products. The following text introduces terminography basics based on terminology work principles (published in the article *Terminologická práca ako súčasť odborno-jazykového vzdelávania* (Štefčík, 2020).

## 1 Terminology work

Terminology work is based on teamwork. Terminology is deeply rooted in logic, which helps to understand the ontologies between concepts (Ogden, Richards; 1921). Some particularities require knowledge of the area processed as terms and the methodology with standards of terminology work (ISO 704 Terminology and terminology work). However, there is a vast difference between the work of a translator and a terminologist. Simply put, terminologists deal with names, concepts, and designations out of the concept and definitions, and translators look for the correct application regarding the context (Arntz-Picht-Mayer; 2002). As terminology is in continuous development, terminologists and terminographers need to be experts, register new changes or the emergence of new terms and expand their vocabulary in the field of terminology. They must be familiar with the subject area in which they prepare the outputs and also master the terminology system of the language for target users. They should use the established terminology means and deviate from them only in justified cases, such as neologisms. In addition, terminologists must be aware of all term-formation procedures to handle even more complex issues and not make inadequate shifts in meaning. Unlike translators, who must look for words and their equivalents in various encyclopaedias or good-quality dictionaries, terminologists and terminographers have to use more resources. Finally, every terminology project brings multiple problems that must be resolved with experts in the particular domain. The second, equally important area is linguistics, which refers to creating designations. Terminology work is a system of steps to excerpt, identify and analyze concepts in the design of the particular domain as validated (and validatable) names (Levická-Zumrík, 2019). Terminology work consists of several processes to be carried out systematically and complementarily (Wüster, 1991). Any terminology work necessarily covers pieces including discussions with several experts for the terminology domain. This is the most critical in any professional or academic work in terminology. Then, the outcome includes processing several terminology entries in a closely researched domain. As the primary goal of the terminology is to collect terms that are as unambiguous as possible, it is necessary to include them in the specific area of users - professionals. To do this, terminologists must precisely define the field whose terminology is being processed. This is very important, especially from the point of view of further elaboration of the concepts (elaboration of the conceptual system), as

well as determining initial names and terms that belong to the thematic area with the target group of users of the terminology being processed.

## 1.1 Main objectives of terminography

Terminography attempts to describe lexical units as terms in a dictionary form. It is a part of terminology work carried out after excerption (selecting terminology data such as names, definitions, and context) and harmonization (analysis of conceptual systems and subject area according to standards, differences, and distinctions between concepts). Terminography tries to solve three cardinal issues:

- defining the word "term" for a specific field;

- incorporating the term into a single unit;

- defining a specific thematic area.

Based on specific criteria, the linguist-terminologist attempts to determine, whether particular terms should be included in the project's terminography output. This is challenging because each domain, terms, and non-terms are very fluid. Therefore, in a particular project, words from the common vocabulary may also appear in the list of terms, but they may be identified as a term for the specific domain. The terms include mainly nouns and multi-word names.

In the terminography project, the output also includes terms that have only been coined during the terminography project, and thus, they are not terms belonging to a specific subject area.

These terms become part of the terminology through the work of terminographers, who often add new terms to dictionaries and terminology databases that are less familiar even to experts in the particular field. They were created as equivalents of foreign language terms and were not used in any official text.

However, the biggest problem tends to be defining the boundaries of the thematic area. In practice, it is common for disciplines including terminologies to overlap. That is why one term can often have its place in different areas. That is why the terminographers should be experts in the field in which they prepare the terminology dictionary. Terminographers approach the project deliberately as experts in their field.

The result of their whole terminography project is a specific type of terminology dictionary.

There are several terminography outputs in print or digital form:

-monolingual terminology dictionaries;

-bilingual terminology dictionaries with different sources and target languages;

-multilingual terminology dictionaries;

-translation of bilingual terminology dictionaries with different sources and target languages;

-translation of multilingual terminology dictionaries with different sources and target languages.

Currently, terminography outputs are mainly produced as primary and secondary resources for specific user groups:

**Glossaries and terminology databases**

**Glossaries** are created for highly specialized professionals who need up-to-date terminology information. They are not only characterized by their limited content but also by their precise targeting of a small group of users. They contain highly specialized up-to-date terminology that has appeared in various specialized sources such as journals, monographs, and various scholarly works. They often occur in smaller units, companies, and multiple manuals as basic instrumentation without meeting exact terminography criteria.

**Terminology databases** are collections of specialized terminologies that include nomenclatures, standardized terms, and lexicalized syntagma, together with the identifying information of their source, which can be used at any time as a monolingual or multilingual reference dictionary, as a basis for the creation of dictionaries, as a means of checking the correct use of terms and their product, and also as a complementary tool of information and documentation technologies. (Sager, 1990)

### 1.2 Terminology records

A terminology record is a written record of information from previous steps of terminology work. There are different approaches to the creation of terminology records. Most of these are based on the recommendations of terminology standards (ISO 12620:1995 and ISO740:000). The structure of the terminology record in each language should be represented by the following data (Křečková 2000):

- Terminology data related to the term (basic form of the term, abbreviated forms of the term, spelling variants of the term, synonyms, foreign language equivalents, context, etc.).

- Data related to the term (definition, graphical representation, antonyms, etc.).

- Administrative data such as concept identifier, language symbol, record creation date, sources, etc.

The International Terminology Standards (Creation and Editing) recommend three essential pieces of information that a record should contain (ISO 10241):

- a record number,

- a recommended term representing the concept,

- definition of the term.

It also gives the possibility to add optional information in rare cases:

- pronunciation

- short form

- grammatical information

- thematic area

- links to resources

- non-recommended terms (permitted, prohibited, obsolete, superseded)

- references related to other records

- an example of term usage

- note and author of the record

- equivalent terms in other languages

- date of record processing (Cabré, 1999).

When compiling some national terminology databases, the following term selection qualifiers are selected according to terminology standards (e.g., Slovak Terminology Database):

*Table 1 Qualifiers for term selection (ISO 10241)*

| | |
|---|---|
| recommended | The terminology committee of the relevant discipline or other relevant institution has recommended the term. |
| standardized | The term is defined by the standard(s) |
| legislative | The term is defined by law or decree. |
| proposed | The relevant terminology committee submits the term for consideration by the professional authority |
| neologism | The term is newly created or newly adopted, the form and content of which may still be subject to change. |
| incorrect | The term is recommended not to be used due to its incorrect form or content. |
| obsolete | The term is recommended not to be used due to its archaic nature. |

## 2 Entries in terminology records

The most critical information in a terminology record is the keyword, definition, and context. The data categories should be recorded uniformly (Horecký, 1971; Masár, 1991).

Adjectives and verbs are less frequently used as headwords, with nouns or abbreviated words being the most common.

The definition or interpretation of the meaning of a term is the most critical information for the user in the terminography output. The terminography product gives only the meaning of the definition of a particular subject area; it does not give the meanings of other subject areas. In a professional terminography project, the meaning of an interpretation is not formulated by the terminographer. Still, it is always taken from the specific professional source to which the terminographer refers. There are also cases where verbal definitions are unnecessary, and an illustration (or video) is provided instead of a definition.

Context is a valuable piece of information that provides an example of the usage of a given term in a text or sentence. Unlike a definition, where the term must not be provided, the term must be given in the case of a contextual indication (Temmerman, 2000).

As an optional entry in the terminology record, synonyms are primarily recorded in explanatory terminography dictionaries. Synonyms have a unique position in terminographies because they are rare, even undesirable, in many projects. This is because terminologies attempt to be as precise and unambiguous as possible in technical vocabulary. However, synonyms tend to be assigned to new fields that gradually build up their conceptual apparatus (e.g., the area of computing, etc.). Nevertheless, terminographers strive to standardize terminology by leading to the use of one term for one concept. It is obvious, but almost impossible, to avoid synonyms, which is why synonyms have become a standard part of terminology glossaries and databases in various fields. Therefore, it is essential that in the case of terminology notation, a uniform structure is maintained so that all synonyms are on the same plane at the level of the concept.

Grammatical entries are rare in terminography projects, and from the point of view of the function, they are considered undesirable (Cíbiková, 2011). The terminography advocates the principle that grammatical entries are given only in the case of exceptional grammatical phenomena.

## 3 Internationalization of terminology and linking terminography projects

The result of any terminography project is the enormous challenge of connecting the work of the linguist and terminologist with that of the IT professional (Drewer-Schmitz, 2017).

This should be both a goal and an opportunity to include terminology management in terminology work with a collaborative model. Terminology management represents ways of managing terminology databases regarding the dynamic evolution of the language. Developing the structure of the terminology record is one of the last and vital steps of terminology work. This should include terminology management with a collaborative model that presents ways of managing the terminology database created concerning the dynamic development of the language. However, the problem is that terminology data needs to be more cohesive across different IT systems and platforms, making it difficult to access and use. In addition, many terminology resources are unavailable online and exist as separate outputs within institutions. Many are stored in basic formats such as Word, Excel, TXT, DOCS, and OTF. In terminology management, the word "termbase" means a database consisting of conceptually oriented terminology records and related information, usually in multilingual format and linked to a specific subject area (terminology database). Terminology databases are also part of CAT systems (CAT=Computer-Assisted Translation), such as MemoQ, SDL, etc. Some CATs exist as stand-alone systems, e.g., the Multiterm product from SDL (RWS group). They use specific translation-related software to help manage the wealth of data in terms of translation and localization outputs. Their internal structure in term entries may vary, but they can be adapted by integrating and (re) different programming APIs. Eurotermbank was a breakthrough in this area.

*EuroTermBank is the largest centralized online terminology database for the European Union and Icelandic languages, linked to other terminology banks and resources. EuroTermBank facilitates the exchange of terminology data with existing national and European terminology databases by establishing cooperative relationships, harmonizing methods and standards, and designing and implementing data exchange mechanisms and procedures* (www.eurotermbank.com).

It functions as a set of terminology resources for the new EU Member States languages, linked to other terminology databases, covering more than 2 million terms and 27 languages.

## 3.1 International terminography projects

The terminology process is a complex set of tasks that use aggregated tools (e.g., corpora) to process large volumes of data.

The following figure proposed by EuroTermBank (www.eurotermbank.com) is one of the terminology processes composed of sub-tasks.
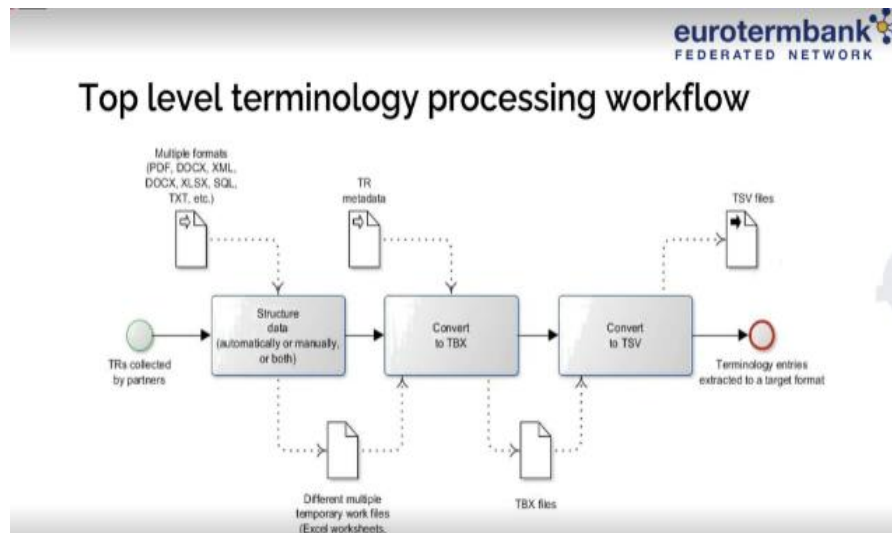


*Figure 1 (www.eurotermbank.com)*

Important terminology projects funded by the EU are IATE or Eurotermbank. It was created as a single access point for terminology search, which has three advantages:

- easy handling and adaptability;

- data sharing;

- wide availability.

IATE (late.europa.eu) is a terminology database. EU institutions and agencies have used it since 2004 to collect, disseminate and manage EU-specific terminology. The Translation Centre works for the Bodies of the European Union in Luxembourg. IATE provides web infrastructure and has undergone many modifications and improvements since its inception. IATE operates in a user-friendly search engine mainly used by translators and interpreters but also by students, teachers, and other professionals. The search is adapted to a semi-professional setup for matching functions, term type search, field search, and domain filtering, which make the terminology a full-fledged database containing e.g. necessary metadata information, and confidence level in all EU languages. Therefore, the database is difficult to manage and maintain, requiring the involvement of a sufficient number of terminologists.

Additional terminology-based resources can be searched via the European Parliament's TermCord platform (DG Trad). This platform includes various terminology projects (e.g. Inmyownterms) and terminology tools. Out of their large number, the following media can be mentioned as an example:

- *Tilde*, which offers terminology extraction and cloud search (https://term.tilde.com/).

- *TermWiki.com* is an online terminology portal that allows users to search, upload, translate and share terms and definitions with other users (https://pro.termwiki.com/).

- *EuroVoc* is a multilingual and multidisciplinary thesaurus containing terms in EU languages (http://publications.europa.eu/eurovoc/index_en.htm).

- *UNTerm* is the United Nations terminology database. It contains technical and specialized terminology in each of the six official UN languages (https://unterm.un.org/unterm/portal/welcome).

- *FAO* Terminology Portal - United Nations Food and Agriculture Organization terminology portal in 6 UN languages (https://www.fao.org/faoterm/en/).

- *Fodina* (TermCatch) - a web platform for creating and maintaining high-quality, consistent terminology.

Other terminology database projects include UNESCOTERM, UNHCR, WTOTERM, OECD Terminology, MultiTes (World Bank Thesaurus), and many others. As already mentioned, the problem with these databases is their fragmented and inconsistent design, which needs more data, various languages, timeliness, and technicality. Given the complexity of terminology work, conventional approaches to terminology work still need to be improved. Traditional terminology work focuses on a formalized approach to defining terms, open data, and semantic search.

Terminology trends include open data proliferation, interconnection, and multilingualism. They target different target groups:

- language community;

- specialized expert groups;

- ordinary users.

Collaboration between the academy and the business community will be essential, as terminology is integral to any operational digital system. A good example is FedTerm, the platform growing and accepting terminology portals from other countries (Federated eTranslation TermBank Network).

## 4 Procedures for creating terminology database

Creating a terminology database is divided into several steps (see Štefčík – Gašová, 2022. Creating a German-Slovak Hunting Terminology Database):

- a conceptual apparatus of the relevant partial professional field;

- searching for and extracting relevant terms and their definitions;

- searching for possible equivalents in the target language(s);

- searching for and extracting relevant technical terms and their definitions in reliable and available resources;

- comparing the definitions related to the source term;

- verifying the meaning and frequency of the term in source and target resources;

- comparing the meaning of terms in source and target resources;

- searching and documenting suitable contexts for a given technical term from authentic sources;

- consulting the terminology record with a professional authority and updating it, if necessary.

Finally, terminology records are directly exported from the Excel sheets to the relevant platform, such as Termweb, XTMCloud, TermWiki, computer-aided CAT translation tools or software for multi-lingual terminology management tlTerm, or term extraction tools for maintaining terminology consistency through the terminology web-based platform.

## 5 Terminology platforms

The everyday basis for any terminology project is the Excel sheets, in which the term entries are structured and processed. There are several advantages to using spreadsheets through Excel, most notably that Excel is compatible with several terminology platforms and CAT tools.

One of the standard terminology platforms is *TermBase eXchange*, an international standard for representing structured concept-oriented terminology data, co-published by ISO and the Localization Industry Standards Association.

However, using Excel sheets also has several disadvantages, especially when processing terminology entries. The main problems include in particular:

- structure = uniform structure of the terminology entries;

- synonyms and data mixing;

- compliance with the principles of uniformity in terminology entries;

- avoiding double work;

- importing into databases of CAT tools – automatic search function aligned with CAT and other applications.

Similar term tools are listed as follows:

-*DatCatInfo* is a Category Repository (DCR) that replaces the former ISOcat. It was developed according to the ISO 12620:2019 series of standards.

-*Termweb4,* which is software for terminology management.

-*TBX info,* the International standard for representing and exchanging information from databases. It includes:

- archiving of information in the terminology database;

- support for software changes;

- exchange of information between systems;

- sending information from the terminology database to the content creation tool;

- translation;

- data mining (exporting most/all information from the terminology database for analysis using XML).

**6 Creating a database or glossary**

Today, the price and availability of different software allow even smaller institutions with limited budgets to create glossaries. One can automatically set up a termbase in the translation process (e.g., translation memory) or manually (e.g., SDL MultiTerm, MultiTerm Extractor, etc.). There are also commercial or freeware tools for the creation and management of translation databases and glossaries:

*Helium* (Microsoft) - transferring volumes of data.

*LogoPort* - translation tool for real-time collaboration, enables virtual workgroups.

*LocStudio* – localization of software products

*Alchemy Catalyst* - visual localization technology

*Multilizer* - multilingual translator of documents

*Toolbox* - managing lexical data for analysis and interlinearization.

There are terminology projects not for open community sharing, e.g., highly specialized and narrowly used professional language within companies. To create smaller glossaries or terminology datasets, the terminologists could use corpora (e.g. national corpus). A popular corpus tool is SketchEngine which has embedded functionalities of term extraction, glossary creation, or creation of specialized corpus from the text on the internet or creating subcorpora.

Sketch Engine supports monolingual and bilingual term extraction. Keyword and term extraction are used to:

- *extract terminology for translation and interpreting;*

- *extract single-word and multiword units typical of a corpus/document/text or which define its content or topic;*

- *compare two corpora/documents/texts by identifying what is unique in the first corpus compared to the second one.*

Sketchengine uses an innovative interface with a responsive design that enables access from different tools (smartphones, laptops). The frequently used functions are:

- *Concordance Search* in one or two languages according to the corpus and pre-defined search criteria (e.g., lemma, phrase, word, character, CQL) for examples of use in context.

- Comparing collocations of two words by using the *Word Sketch Difference* function.

- *Thesaurus* displays the function of words occurring in similar contexts and demonstrates differences in collocation profiles.

- *Creating a frequency list* by using the function of *Wordlist*.

- *Parallel concordance* with parallel corpora for translation search.

- *Keywords* for bilingual terminology extraction.

- *One-click Dictionary* enables automatic dictionary drafts exported into the Lexonomy writing system by creating an account and an API key.

- *Create a corpus* to build a private corpus from the web or personal texts (documents). A translation-useful function of the private corpus is extracting terms by creating a frequency list out of the private corpus and comparing it with a frequency list of a larger reference corpus.

From the practical perspective, each terminology project struggles with questions that should be asked and resolved initially.

Those questions include the following issues:

- What lexemes should be registered in the new database as terms? The content applies not only to parts of speech (nouns and verbs are mainly used as terms) but also to those with the same form but different meanings in various fields and domains. Typically, standard terms may also be included in the database once they have a specific meaning in the domain being processed and tailored to target users.

- Application of terms for purposes other than translation. Each termbase glossary creator should remember that not only translators might use the glossary. Therefore, the focus should be on the content rather than grammar information.

- More equivalents for one term. This should be perfectly all right if the termbase or glossary is multilingual and targets the translation domain.

- Multiple meanings of one lexeme. In a terminology database, this can be solved by domains and their categorization. However, a glossary should be polysemy avoided by providing explicit definitions of concepts.

- One term belongs to one line. This rule should prevent synonymy, homonymy, and orientation confusion in the glossary.

- Terms should be standard and uniform (e.g., singular, nouns, - infinitive, lowercase letters, omit "ing" forms in English, etc.).

- Usage of alternative forms of terms should be avoided. Alternative forms should be included only exceptionally in specific cases.

- Grammatical information about the term must be structured and differentiate parts of speech.

- In the case of two terms, such as regular term and its acronym, they should be placed in two separate lines (cells).

- Synonyms and lexemes with different conceptual apparatus must be distinguished.

- Context information is always recommended and should be included in every terminology project.

- Definition with context must be clear and strictly distinguished.

- Notes in term entries are essential to add information, not in other term records.

- Information about the author of the record, the date, and the applicability of the term is relevant information that should be included.

- Part of speech with the type of term should be clear (e.g., full term and its masculine form).

**Conclusion**

Terminography outputs should not be perceived only as a single-purpose tool for removing barriers in professional communication, but, as a tool for expanding the professional vocabulary in the cooperation of several experts in a specific subject area. Moreover, our efforts should be directed toward finding solutions for multilingual platforms by unifying the form and structure of terminography records so that they are easily accessible at any time and from

any device. There are significant international projects involving teams from the commercial sector to address these challenges. Slovak, as a less widely spoken language, should, within its capacities, develop terminology competence and terminography projects that would become part of the large family of other languages included in comprehensive terminology databases.

There are also several other technology instruments to be used for term extraction, including CAT tools, e.g., SDL Trados Multiterm, MemoQ, with translation memory software, translation memory editors, terminology management software, review software, etc., used by translators. Also, several other possibilities do exist in terms of terminology processing. Additionally, concerning the near future perspective, besides standard tools such as CAT or TMS, other disruptive technologies like *blockchain* (see How Blockchain Technology Can Reshape the Language Industry) may be considered in the following research stage of the terminology project with a focus on terminography.

## References

ARNTZ, R. PICHT, H. MAYER, F. 2002. *Einführung in die Terminologiearbeit*, Olms.

CABRÉ, M. T. 1999. Terminology. *Theory, methods, and applications*. John Benjamins Publishing Company.

CÍBIKOVÁ, I. 2011. *Terminologický manažment verejnoprávnej tematickej oblasti*. Vydavateľstvo Edis.

HORECKÝ, J. 1956. *Základy slovenskej terminológie*. Bratislava: SAV, 1956.

HORECKÝ, J. 1974. *Obsah a forma termínu*. In: Kultúra slova, 1974, roč. 8, č. 10.

JURČACKOVÁ, Z. 2002. *Terminológia. Základné zásady, metódy a ich aplikácia*. CVTI Bratislava.

KREČKOVÁ, V. *Terminológia, terminologická práca a medzinárodné normy*. In: Tlmočenie – preklad, roč. 11, č. 53, s. 13.

JUAN C. SAGER. 1990 *A Practical Course in Terminology Processing*. Amsterdam; Philadelphia: John Benjamins.

LEVICKÁ, J., ZUMRÍK, M. 2019. *Terminologické inšpirácie profesora Jána Horeckého*. VEDA: SAV, 2019.

MASÁR, I. 1991. *Príručka slovenskej terminológie*. Bratislava.

MASÁR, I. 2000. *Ako pomenúvame v slovenčine.* Bratislava.

OGDEN, C. K., RICHARDS, I. A. 1923 *The meaning of meaning*. [online] 2020 [cit. 2020-08-10] Dostupné na internete:

https://courses.media.mit.edu/2004spring/mas966/Ogden%20Richards%201923.pdf

SCHMITZ, K-D.-DREWER, P. 2017. *Terminologiemanagement*. Springer Verlag

ŠTEFČÍK, J. 2020. *Terminologická práca ako súčasť odborno-jazykového vzdelávania* In: Aplikované jazyky v univerzitnom kontexte VII. Technická Univerzita vo Zvolene.

ŠTEFČÍK, J., GAŠOVÁ, Z. 2022. *Creating a German-Slovak Hunting Terminology Database*. In Folia Linguistica et Litteraria : Časopis za nauku o jeziku i književnosti. - Nikšić : Institute for Language and Literature.

TEMMERMAN, R. *2000 Towards new ways of terminology description. The sociocognitive-approach*. John Benjamins Publishing Company.

WREDE, O., ŠTEFČÍK, J., DRLÍK, M. 2016. *Úvod do terminológie a terminologickej práce*. UKF Nitra.

WÜSTER, E.1991. *Einführung in die allgemeine Terminologielehre und terminologische Lexikographie*. Bonn

https://termcoord.eu/terminology-websites/

https://xtm.cloud/

https://sk.termwiki.com/

https://interverbumtech.com/

SketchEngine (https://www.sketchengine.eu/)

Slovenská Terminologická Databáza: terminologickyportal.sk

Transterm:http://kger.web2v.ukf.sk/transterm/home

Federated eTranslation TermBank Network (https://cst.ku.dk/english/projects/fedterm/)

STN ISO 704: *Terminology. Terminology work*

STN ISO 10241: *International terminology standards.* Creation and modification