# TESTING CHATGPT ON TERMINOLOGY GENERATION, DEFINITIONS: TRANSLATION, AND ONTOLOGY CREATION IN GERMAN, ENGLISH AND POLISH

*BARBARA LEWANDOWSKA-TOMASZCZYK*
University of Applied Sciences in Konin
barbara.lewandowska-tomaszczyk@konin.edu.pl

*GRZEGORZ PAWŁOWSKI*
University of Applied Sciences in Konin
grzegorz.pawlowski@konin.edu.pl

**Abstract**
The paper presents the generation process and analysis of Chat GPT-produced terminology, definitions, translation and ontology creation of content derived from a restricted domain of electrotechnology on the basis of German manufacturing instruction (*Fertigungsvorschrift)* for the assembly and functioning of thermal switches. The first part of the study presents a brief historical sketch of corpus linguistics towards the development of LLMs and their applications. Research steps include the generation of terminology and a basic restricted domain ontology in German and their equivalents in English and Polish as identified in general web resources, followed by the same tasks based on the analysed tool instructions. The tests were first performed on earlier ChatGPTPro, followed by its recent version ChatGPT4, also with regard to their English and Polish equivalents, for different tasks, i.a., relevant thesauri building. An Assistant called SLA (Special Language Assistant) was created for the purpose of analysing prompts, recognising context and intent, and processing data using language models.

The tests have been carried out in terms of 4 prompts. First, specialist terms typical of the product and its parts were identified in German and the translation of these terms and of other relevant specialist phrases into English and Polish was performed. Finally, the ontological categorisation and its visualisation were generated. Results indicate areas of fair correspondences with manual intervention needed for the term definition refinement and ontology specification. The presentation emphasizes in the conclusions the effects of AI tasks as a contribution to lexicography, translation and foreign language education. The study can also serve as a reference for NLP researchers to improve the functioning of LLM tools.

**Keywords:** ChatGPT, cluster equivalence, corpus, LLM, ontology, ontoterminology, prompt, term, translation

# 1. Introduction

The main tenet of the present paper refers to the evaluation of LLMs, most notably ChatGPT 4, in generating data for restricted domain thesauri. The general theme discussed here, often researched at present (e.g., Veseli et al. 2023, Vrolijk 2022), refers to the automatic terminology generation and ontology creation in German, English, and Polish to serve as a tool towards building an ontoterminological dicionary/thesaurus of a restricted domain (compare Radszuweit & Spalier 1982, Li et al. 2004, Dornseiff 2003, Hallig & Wartburg 1963, Sierra 2000, Schryver 2023). The main reason for such an enterprise is the fact that manually-curated background knowledge for many fields is a rather scarce resource Roget (1852, 1984), while the possible use of automatic tools would make this part of the product preparation easier and faster.

The development of corpora and the maintenance of restricted domain thesauri, crucial for information retrieval, knowledge organization, and semantic interoperability in specialized fields, have traditionally been a resource-intensive endeavour. Carefully curated controlled vocabularies demand expert human input to define precise relationships between terms, ensuring accuracy and consistency within a narrow subject area. However, the emergence of Large Language Models (LLMs) presents a compelling opportunity to automate and accelerate aspects of this process, particularly in generating initial data and suggesting new entries or relationships (Rees & Lew 2023). The application of language corpora is fundamentally connected to checking the functioning and evaluating the performance of Large Language Models (LLMs). Language corpora function as a source of knowledge and training data for such tools. This study evaluates the potential of LLMs in generating data for restricted domain thesauri, exploring their capabilities, inherent limitations, and the critical considerations for their effective and responsible application (Lewandowska-Tomaszczyk 2022/2023).

LLMs, with their remarkable ability to generate human-like text, can leverage vast amounts of general and, when fine-tuned, domain-specific text data to infer semantic connections and propose terms relevant to a particular field. Their capacity to identify synonyms, broader and narrower terms, and related concepts from unstructured text, could significantly reduce the manual effort involved in the initial stages of thesaurus construction. This synthetic data generation can be particularly valuable in ow-resource domains where existing structured data is scarce, offering a rapid way to bootstrap the thesaurus creation process.

Despite these promising capabilities, LLMs possess inherent limitations that necessitate careful evaluation when applied to restricted domain thesauri. Firstly, their knowledge is a reflection of their training data, meaning they may struggle with novel concepts, highly specialized jargon not adequately represented in their corpus, or real-time changes within a dynamic domain. This can lead to "hallucinations" – the generation of plausible but factually incorrect or nonsensical terms and relationships – which are unacceptable in the precision-driven world of controlled vocabularies. Secondly, LLMs lack true reasoning and understanding. They identify patterns and predict sequences based on statistical correlations, not a deep comprehension

of the underlying domain. This can result in errors in hierarchical or associative relationships, misinterpretations of polysemous terms within a specific context, or an inability to discern subtle nuances that are critical in a restricted domain.

Furthermore, the nature of many LLMs makes it challenging to trace the provenance of generated data or understand the rationale behind proposed relationships, hindering validation and quality control. The contextual window limitations also mean that LLMs can only process a finite amount of information at a time, making it difficult for them to grasp the full scope and intricate interconnections of a large, complex thesaurus. Therefore, while LLMs can serve as powerful data generation engines, their outputs demand rigorous human expert review and validation. Their role is best envisioned as an intelligent assistant, automating the heavy lifting of initial data discovery and suggesting potential entries, rather than a fully autonomous thesaurus creator. The ultimate success of LLM integration into restricted domain thesauri development hinges on a synergistic approach that combines the generative power of AI with the indispensable domain expertise and critical oversight of human thesaurus builders. In other words, the study presented in this paper refers to the area of restricted domains, terms and terminology.

## 2. Restricted domains, terms and terminology

The field of information organization and retrieval is profoundly shaped by the concepts of *restricted domain*, *terms*, and *terminology*. While seemingly straightforward, their precise definitions are crucial for understanding the complexities of specialized knowledge representation, particularly in the context of thesauri, ontologies, and controlled vocabularies.

A *restricted domain* (also known as a specialized, specific, or narrow domain) refers to a delimited subject area or field of knowledge characterized by a specific set of concepts, relationships, and often, a specialized vocabulary. Unlike general domains that encompass broad areas of human knowledge, a restricted domain focuses on a particular discipline, industry, or activity. Examples include medicine, law, aerospace engineering, art history, or specific sub-disciplines within these, such as cardiology within medicine, or intellectual property law within law. The defining characteristic of a restricted domain is its inherent boundaries and the homogeneity of its subject matter, which allows for a high degree of precision and detail in its knowledge representation. The information within a restricted domain is often highly structured, interdependent, and requires specialized knowledge to interpret accurately. This focus enables the development of highly effective, domain-specific tools for information management, as the ambiguity inherent in general language can be significantly reduced.

*Terms* are the linguistic units, typically words or phrases, used to represent concepts within a given domain. In a restricted domain, terms are not merely everyday words; they are often *technical terms* or *domain-specific terms* that carry precise and unambiguous meanings within that particular field. For instance, while "cell" in a general context could refer to a prison cell, a battery, or a phone,

in the restricted domain of biology, "cell" specifically denotes the basic structural and functional unit of all known organisms. Terms can be single words (e.g., "algorithm," "neuron"), multi-word expressions (e.g., "artificial intelligence," "deoxyribonucleic acid"), or even acronyms and abbreviations commonly used within the domain (e.g., "MRI" in medicine). The selection and standardization of terms are fundamental to ensuring clarity and consistency in communication and information organization within a restricted domain. The goal is to avoid synonymy where a single concept is represented by multiple terms, and homonymy where a single term represents multiple concepts, as both introduce ambiguity. Thus *terminology*, in its broadest sense, refers to the systematic collection of terms used in a specific domain. More precisely, it encompasses the study and systematic collection of specialized words and phrases (terms) and their usage within a particular field. Terminology is not just a list of words; it involves the analysis of concepts, their relationships, and the linguistic expressions used to denote them within a restricted domain. It is a structured and organized body of knowledge that defines the vocabulary of a specific subject. Developing a robust terminology for a restricted domain involves identifying key concepts, defining them precisely, and establishing the preferred terms to represent them. This systematic approach ensures that experts and practitioners within the domain communicate effectively and unambiguously, facilitating knowledge sharing, education, and research. In essence, terminology provides the linguistic framework and conceptual clarity necessary for efficient operation within any specialized field. In other words,, terms are the lexical units assigned a unique sense and identified with reference to only one conceptual system or a restricted knowledge domain as e.g., the concept of *light* in the wave theory or the concept of *grammar* in the linguistic transformational-generative theory or in a cognitive grammar. Both the understanding of terms in a SL as well as their TL equivalents presuppose the familiarity with the unique conceptual system underlying a theory a given concept is a part of.

## 3. Problems: synonymy, polysemy and cluster equivalence

Some of the problems with the identification and generation of terms reside in the fact that they are either synonymous with other terms or else that they are part of a larger polysemy network. These facts are particularly troublesome in the case of translation and the question of (the range of) their translational equivalence. Let's take as an example the Polish form *umowa.* In the case of general language use, this word can have the meaning of an agreement, but at the same time, it can refer only to some act or action in the future, close in its sense to a 'planned proposal'. As a legal term though, it would correspond rather to the English phrase 'a legal contract', i.e., a legal act which is binding for the two co-operating parties.

However, there can also be the form 'umowa' used with reference to so-called 'unnamed contract' in Polish, i.e., *umowy nienazwane*, i.e., *contractus innominate*, which, in their definition, come closer to Pol. *obietnica* – Eng. *promise*[1].

## 4. Problems: synonymy and translational non-equivalence

Connected with this fact is the issue of synonymy in such cases as *contract - promise* (USA), the latter being an instance when a promise lacks the formal elements of a contract but can be considered an instance of a *promissory estoppel*:

1. *Promissory Estoppel:* promisoryjne ograniczenie

   1. Even when a promise lacks the formal elements of a contract, courts may enforce it under the doctrine of promissory estoppel. This doctrine prevents a promisor from going back on their promise if the promisee reasonably relied on it to their detriment.

   2. For example, if someone promises to donate money to a charity, and the charity relies on that promise to take specific actions, the court may enforce it.

The problem here resides in the fact that, for example, in about half of the states in the United States, a promise to marry is considered to be legally enforceable as long as that promise meets all the basic requirements of a valid contract. In Polish legislature though, such a promise would not be legally enforceable even though there exist some circumstances to enforce it. Therefore, the analysis of such monolingual synonymy of a cross-language type of equivalence or meaning differences is best investigated in terms of the contextual uses in large corpus datasets.

## 4. Corpora, corpus linguistics and Large Language Models (LLMs)

The fields of corpus collection and corpus linguistics have undergone a transformative journey, evolving from nascent theoretical concepts into indispensable tools for linguistic research, language teaching, and natural language processing. This evolution, spanning several decades, reflects advancements in computational power, an increasing understanding of linguistic phenomena, and a growing appreciation for empirical, data-driven approaches to language study. Corpus linguistics is the study of language based on large collections of real-world text and speech and the development of corpora is

---

[1]   *umowa: legal contract*

   1. *Umowa to czynność cywilno-prawna,* której dokonuje się w oparciu o oświadczenie woli obu stron umowy. W prawie cywilnym jest to porozumienie, które określa wzajemne prawa i obowiązki, np. dwóch kontrahentów podejmujących współpracę.

   2. *Umowami nienazwanymi (contractus innominati – unnamed contract)* są te umowy niemające w ogóle odrębnego, normatywnego odniesienia, nie zostały one uregulowane ani w kodeksie cywilnym, ani w innych aktach prawnych porządku krajowego vs obietnica.

the foundational step, as their quality, size, and representativeness directly impact the validity and generalizability of any linguistic insights derived.

The earliest forays into corpus creation were largely manual and driven by a desire to capture language in its natural use, moving beyond introspective or anecdotal evidence. A pivotal early milestone was the *Brown University Standard Corpus of Present-Day American English*, compiled in the 1960s by Henry Kučera and W. Nelson Francis. This million-word corpus, meticulously hand-tagged for part-of-speech information, was revolutionary. It demonstrated the feasibility and immense value of systematically collecting and analyzing a substantial body of text. Its British English counterpart, the *Lancaster-Oslo/Bergen (LOB) Corpus*, followed, paving the way for comparative studies of different English varieties. These early corpora, though small by today's standards, established fundamental principles: the importance of representativeness (sampling different genres and domains), the need for standardized annotation, and the potential for quantitative analysis of linguistic patterns.

The 1980s and 1990s witnessed a significant acceleration in corpus development, largely fueled by advancements in computing technology. The advent of personal computers and improved storage capabilities made it possible to handle much larger datasets. This era saw the creation of monumental projects like the *Bank of English* (later incorporated into the *Collins Word Web*), a multi-billion-word corpus that truly showcased the power of scale. Similarly, the *British National Corpus (BNC)*, a 100-million-word balanced corpus of written and spoken British English, became a gold standard for its size, careful design, and comprehensive annotation. These larger corpora moved beyond simple frequency counts, enabling researchers to explore more subtle grammatical patterns, lexical collocations, and register variations. The development of sophisticated corpus management software, like the Wordsmith Tools and the BNC's own query system, made these vast datasets accessible to a wider research community.

The turn of the millennium brought forth an explosion in corpus diversity and size, driven by the proliferation of digital text and the demands of emerging fields like natural language processing (NLP). The internet became an unparalleled source of linguistic data, leading to the creation of web-as-corpus initiatives. While offering immense size, web-derived corpora presented new challenges related to data cleaning, noise reduction, and ensuring representativeness across different online domains. Simultaneously, specialized corpora began to emerge – parallel corpora for machine translation, learner corpora for second language acquisition research, historical corpora for diachronic studies, and spoken corpora capturing the nuances of everyday conversation. The development of robust annotation tools and guidelines, often building on shared standards like the *Text Encoding Initiative* (TEI), became crucial for making sense of this increasingly complex data.

The current landscape of corpus development is characterized by several key trends. Firstly, the drive for ever-larger corpora continues, with multi-billion and even trillion-word corpora becoming increasingly common, particularly in the context of training large language models (LLMs). Secondly, there is a growing

emphasis on multimodal corpora, integrating not just text but also audio, video, and even physiological data to capture a more holistic view of language in use. Thirdly, annotation is becoming more sophisticated and multi-layered, moving beyond basic part-of-speech tagging to include syntactic parsing, semantic roles, discourse structures, and even emotional states. This rich annotation allows for more nuanced linguistic analysis and facilitates advanced NLP applications. Finally, accessibility and sustainability are paramount, with many major corpora being made publicly available and initiatives like the CLARIN European Research Infrastructure Consortium (ERIC) aiming to provide long-term access to linguistic resources.

The impact of corpus development on corpus linguistics has been profound and symbiotic. The availability of large, diverse, and well-annotated corpora has transformed linguistic research from a predominantly introspective and prescriptive endeavor into an empirical and descriptive science. Corpus linguistics has enabled a data-driven understanding of language in several key areas:

- Lexicography and Lexicology: Corpora are the backbone of modern dictionary making. By analyzing vast amounts of text, lexicographers can identify the most frequent words, their typical collocations, their various senses, and their usage patterns. This has led to more accurate, up-to-date, and usage-based dictionaries. Corpus analysis also facilitates the study of semantic change, neologisms, and the specialized vocabulary of different domains.

- Grammar and Syntax: Corpus data allows researchers to identify actual grammatical patterns, challenging prescriptive rules that may not reflect real-world usage. It has revealed the frequency and distribution of different grammatical constructions, shedding light on their functional motivations. For example, corpus studies have illuminated the nuances of passive voice usage, the patterns of complementation, and the preferred structures in different genres.

- Discourse and Pragmatics: By examining extended stretches of text or transcribed speech, corpus linguistics provides insights into how language is used in context, beyond the sentence level. This includes the study of cohesive devices, turn-taking in conversation, speech acts, and the linguistic features that characterize different discourse types (e.g., academic writing, news reports, casual conversation).

- Sociolinguistics and Variation: Corpora with metadata about speakers (e.g., age, gender, region, social class) are invaluable for studying language variation and change. Researchers can track the spread of linguistic innovations, analyze dialectal differences, and investigate how social factors influence language use. Diachronic corpora, specifically designed to capture language over time, are crucial for historical linguistics.

- Language Teaching and Second Language Acquisition (SLA): Learner corpora, composed of texts written or spoken by language learners, provide empirical evidence of common errors, developmental stages, and areas of difficulty. This information is vital for designing effective language teaching materials and developing data-driven pedagogical approaches. Pedagogical grammars and dictionaries can be informed by what is truly frequent and useful for learners, rather than relying solely on intuition.

- Natural Language Processing (NLP): It is one of the most significant beneficiaries of corpus development is NLP. Corpora are the raw material for training machine learning models in tasks such as machine translation, speech recognition, sentiment analysis, text summarization, and information extraction. The availability of massive, often unannotated, text datasets has been a key factor in the recent successes of deep learning models in NLP, though richly annotated corpora remain essential for fine-tuning and evaluating these Larhe Language Models (LLMs).

## 5. Translational cluster equivalence (Lewandowska-Tomaszczyk 2017)

As discussed in the previous sections, corpus development and corpus linguistics are subject to continuous innovation and expansion. From the pioneering efforts of the Brown Corpus to the vast, multi-modal datasets of today, the field has consistently pushed the boundaries of how we collect, process, and understand linguistic data, i.a., towards the lexical and sentential embeddings as well as automatically generated ontologies. Translation memories and large data processing models.

One of such tools is *Sketch Engine (https://www.sketchengine.eu/)*, which can help to identify synonym and polysemy referred to in Section 3 of the present paper, crucial for the preparation of the tools for ontoterminological materials in thesauri and dictionary building such as e.g., in  in distinguishing close senses such as offensive vs insulting () and a development of parallel – translational corpora and translational equivalents identification.
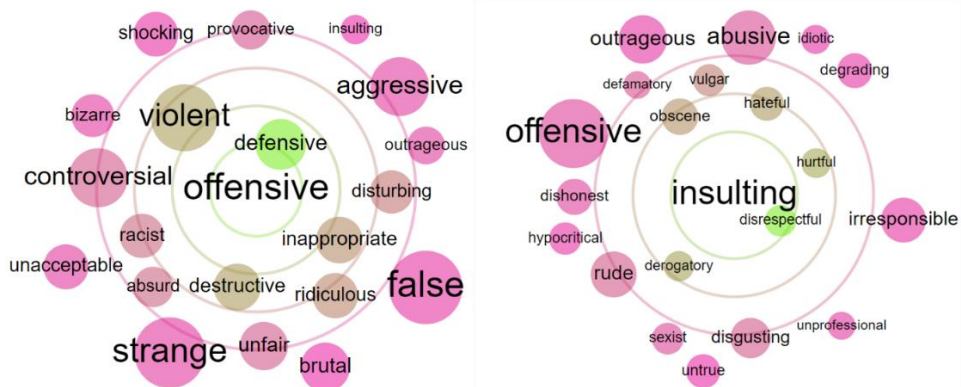
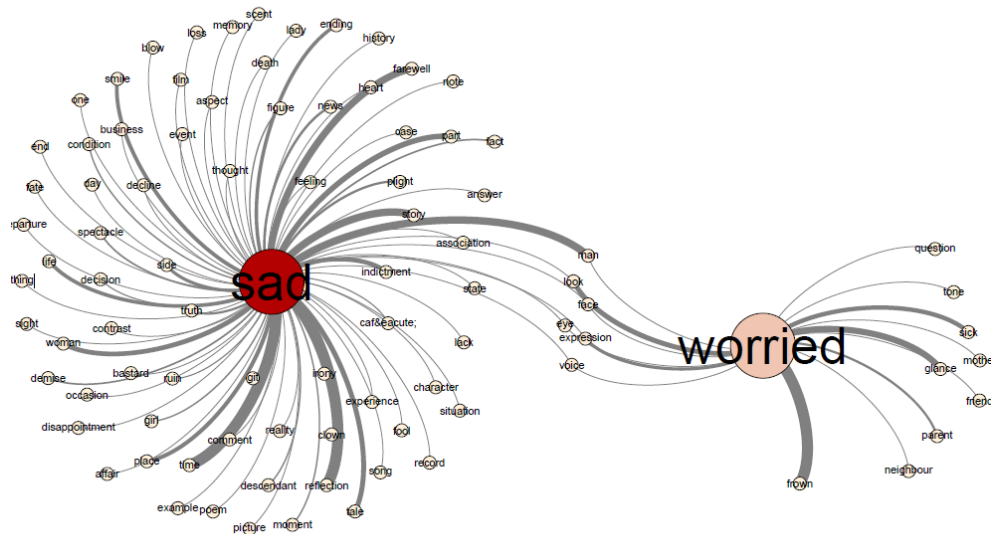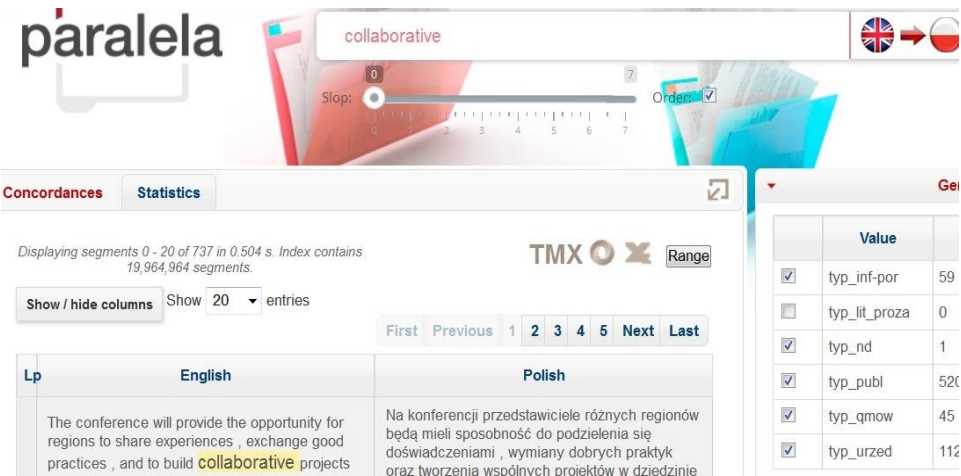**Figure 1:** Comparison of synonyms Sketch Engine (https://www.sketchengine.eu/



**Figure 2:** Colosaurus (Pęzik 2014): collocational differences

## 6. Paralela: translation corpus



*Paralela* (Pęzik 2016) English to Polish parallel corpus

The translational corpora help us identify *specific wording* in the TL text and *subtle and yet meaningful) differences* between what looks like putative synonymy among the translation variants, It was precisely the move to translational corpora that paved the way to the development and identification of the concept of the *translational cluster equivalence* (Lewandowska-Tomaszczyk 2017), which covers the (partial) synonymity and polysemy relations between the Source Language and Target Language corresponding items in a particular domain of their uses.

## 7. Ontologies of restricted domain

Originally developed for all natural terms in the surrounding reality by Aristotle, nowadays, each linguistic variety characteristic of a given restricted domain (sublanguages) is associated with its own *ontology*, i.e., a hierarchical system of terms representing a relevant part of the outside reality, such as e.g. contract law, social sciences, petroleum industry or biochemistry. Every concept in a restricted domain can be classified into different levels with different depths according to the purpose of ontology construction, existing ontology foundation, as well as to the level of ontology complexity, and application demands.

The figures below (3,4), for example, show part of the classification hierarchy of Petro-Onto and has been developed by Tang et al. (2023) and Li et al. (2024).

## 7.1 Ontology visualisation

The Petro-Onto ontology presents petroleum exploration and development - basin oil and gas resource prediction and evaluation, together with the basin simulation, basin simulation and other parameters as well as more dependent material features such as fluid property - viscosity, and others.



**Figure 3:** Petro onto-terminology (Tang et al. 2023)

In the below figure (4) petroleum industry repositories per country, crucial for petroleum industry are represented:



**Figure 4:** Petroleum Industry repositories per country (Tanng et al. 2023).

In the forthcoming sections we will scrutinize the application of ChatsGPT in defining, translation and ontology forming tasks in a restricted electrotechnical domain in German, English, and Polish.

## 8. Analysis of electrotechnological data in a manufacturing instruction manual

In this part of the paper we will refer to a manual ontoterminological analysis of electrotechnological data in a manufacturing instruction manual presented in Pawłowski (2023) and that generated by Chat GPT's comparative responses from available sources of the same domain across the three languages: the German original and its English and Polish translations and ontology formation with a request addressed to the LLM to explicitly account for each of the steps in proceeding of these tasks.

## 9. Research

### 9.1. Background

In Pawłowski (2023) the author presented a *manual analysis* of a German manufacturing instruction (*Fertigungsvorschrift)* terminology, for the assembly and functioning of thermal switches. We make reference to this study as a source of the original manual analysis of the content and reference to the current LLMs- based extended version, geared for more complex tasks.

### 9.2 Major research steps in the current study

It needs to be recalled that the major research aim in this study is to present results of testing of ChatGPT functions for the terminological translation capacity and ontology building for relevant dictionary/thesaurus preparation.

The tests were performed on two tools independently – an earlier version of LLMs – ChatGPTPro, and on recent ChatGPT4 version. They included two larger batches: prompts referring to the manual data, and prompts asked to scrutinize all data available to the system: (i) generating sets of relevant terminology in German, then in Polish, and in English  (ii) providing definitions of the terms (words and phrases ) in German, (iii), translating them into Polish and English, followed by (iv) to generate, and (v) to visualise the ontoterminological data in German, Polish, and English, and (vi) provide the data description and interpretation.

In the first batch, with prompts based on the manual, ChatGPT was asked (vii) to additionally provide – in the three languages – *informal professional equivalents* of all the terms generated in the previous testing steps

## 10. Extraction and ontological analysis of professional terminology

This practical section presents an example of ontological extraction and analysis for the translation of technical documentation from the electrical engineering

industry from German into English/Polish using appropriately prepared AI tools. Each stage of the analysis corresponds to a separate research task involving the extraction of terms, their semantic interpretation, classification in an ontological structure, assignment to categories, visualisation of relationships and evaluation of the usefulness of the resulting model for translators and linguists.

The whole process illustrates how linguistic and ontological tools, and above all the use of AI tools, can be used to systematically describe domain knowledge and facilitate and accelerate the translation of technical documentation.

The documentation concerns products from the electrical engineering industry - thermal switches, also known as bimetal thermal switches, which are an important protective element against overheating in household appliances, in the automotive industry and in other industrial devices. Their operation is based on the deformation of a bimetallic disc in response to temperature, which results in the opening or closing of an electrical circuit. Such systems require precise manufacturing and quality control. The technical documentation contains specialised terms describing components, tools, process parameters and individual assembly steps.

To analyse the documentation of the above-mentioned products, we used the aforementioned LLM ChatGPT 4.0, in which an assistant called SLA (Special Language Assistant) was created for this purpose. The AI assistant works on the basis of natural language processing (NLP) and machine learning, analyses entered commands (prompts), recognises context and intent, and processes data using language models trained on huge data sets to predict the most accurate response or action. Finally, it generates responses in a natural way that is appropriate for the objectives set. An important feature of the AI assistant is its ability to learn through interaction.

As part of the assistant's configuration, its role has been defined as follows:

*"You are a linguist and specialist in professional language analysis and translation. You have been tasked with linguistic analysis of technical documentation in electrical engineering as in the uploaded files. This is a setup manual for thermal switches, which are used as protection for several household appliances. You are analysing the professional terms, phrases and constructions which are specific to electric and electronic fields and these products."*

Next, the original technical documentation for the production of the thermal switches described above was loaded in the form of 16 PDF files. These files contain lists of parts, descriptions of assembly operations, tools required for production, safety equipment, quality control instructions, etc. All instructions are written in German using specialist terminology specific to the electrical engineering industry.

## 10.1 Instructions

The analysis performed using the AI assistant defined above consists of several stages, described by individual project instructions (prompts).

(1) The first stage involves extracting concepts using the following instruction (prompt 1):

*"First you are analysing the language of the documentation and then find all technical terms in the form of (1) nouns, (2) adjectives + nouns, (3) noun + verb, (4) adverb-verb combinations that are typical for the language of technology, this product, and its assembly processes described in the uploaded documentation. Find and list all of the German terms in the uploaded documentation one by one".*

The aim of this stage is to identify all technical terms appearing in the thermal circuit breaker assembly instructions. The search has been narrowed down to four dictionary structures: (1) independent nouns, (2) adjective-noun compounds, (3) verb-noun-phrases, and (4) adverb-verb phrases. This definition of terminology allows for the capture of the entire lexical layer relevant to the description of all assembly processes, selection of production means, and quality control.

First, technical terms appearing in the documentation were identified, including nouns (e.g. "Bimetallscheibe", "Kraftmessgerät"), etc., resulting in list no. 1. List No. 2 presents adjective-noun combinations (e.g. "mechanische Belastung"), noun-verb expressions (e.g. "Kappe bördeln") were collected in list No. 3, and adverb-verb combinations (e.g. "sorgfältig verschließen") in list No. 4. The lists of terms resulting from the search of the instructions contain several dozen key words describing parts, tools, assembly operations, parameters and control processes, thus creating a glossary of terms that guarantees consistency and uniformity of terminology in the translation.

(2) The next step is semantic analysis aimed at standardising concepts and their ontological classification, which is described in the instruction (prompt 2):

*"Conduct a detailed analysis of the meaning of all above terms and propose some possible and common ontological categories for all these terms, the production of thermal switches and other processes".*

In this way, all terms from the above four groups of concepts were analysed in terms of their meaning and role in the production process, so that their understanding and classification as, for example, tools, parts or means of production left no room for doubt. This is particularly important when the translation is carried out by a team of translators with varying levels of experience, subject knowledge and industry expertise.

As a result of the above prompt, we obtain a detailed semantic analysis of the technical terms extracted from the thermal switch production documentation which reveals a highly specialised lexicon that serves a structured and procedural environment. These terms support precise communication in manufacturing, quality assurance, and safety processes for electro-technical devices. Below is an analysis of the meaning dimensions involved, followed by a proposal of ontological categories suitable for organising such terminology in an industrial knowledge model.

To structure the lexicon for technical documentation, a domain-specific ontology might include the following top-level and mid-level classes:

*1. Component*
- Physical parts and materials used in the device.
- Examples: *Socket*, *bimetal disc*, *cap*, *spring*, *stranded wire*.

*2. Tools and devices*
- Tools, machines, fixtures used in production or testing.
- Examples: *Hydraulic press*, *welding machine*, *dosing device*.

*3. Process Step*
- Actions performed during manufacturing.
- Examples: *flanging*, *measuring*, *gluing*, *bending*.

*4. Process Parameter*
- Measurable quantities that influence or define processes.
- Examples: *working pressure*, *welding current*, *test voltage*, *torque*.

*5. Measurement and Evaluation*
- Terms for control, inspection, and qualification.
- Examples: *visual inspection*, *continuity test*, *force measurement*, *test certificate*.

*6. Defect and Deviation*
- Descriptions of flaws, errors, and tolerance violations.
- Examples: *deformed*, *not OK*, *exceeded*, *rejected*.

*7. Operator and Role*
- Human agents involved in the process and their responsibilities.
- Examples: *worker*, *foreman*, *setter*, *quality control*.

*8. Material and Substance*
- Substances and materials applied or joined.
- Examples: *adhesive*, *solder*, *sealing ring*, *casting compound*.

*9. Protective Equipment and Safety*
- Personal or procedural safety devices and rules.
- Examples: *safety goggles*, *protective mask*, *safety regulations*.

*10. Fixture and Mounting Support*
- Non-final parts used for holding, aligning, supporting during production.
- Examples: *assembly pallet*, *fixing adapter*, *gripping container*.

*11. Data and Documentation Artifact*
- Records and digital/analogue artefacts of production.
- Examples: *measurement report*, *construction specifications*, *production order*.

The above eleven categories seem to fully reflect the ontological structure of the domain, and the proposal prepared by SLA for dividing the obtained terminology into ontological categories allows for indicating the level of abstraction and relationships between concepts, as well as for assigning all previously extracted terms to one of the above categories. It should also be noted that this categorisation is based on

ontological engineering practices and the needs of translators, taking into account specific substantive aspects (e.g. "Sockel" as a housing component), process-related aspects (e.g. "nietować" as an act of joining), measurement aspects (e.g. "Prüfspannung"), tools (e.g. "Hydraulikpresse") or work safety (e.g. "Schutzkleidung").

(3) The consequence of the above stage is the process of assigning the concepts obtained using the first prompt to the eleven ontological categories provided by the SLA and generating results for better clarity in the form of a spreadsheet table. This is done using the instruction (prompt 3):

*"Now assign absolutely all technical terms from the attached documentation listed above to the above 11 categories. Complete this task in the form of 11 tables, one for each category. The category name is in the header of each table. The first column contains the ordinal number, the second column contains the technical term in German, and the third and fourth columns contain the translation of this term into English and Polish."*

The result is a table of the previously listed terms assigned to the above categories in the form of a table:

| Nr. | Professional Word (DE) | Professional Word (EN) | Professional Word (PL) |
|---|---|---|---|
| 1 | Sockel | base (housing) | Cokół |
| 2 | Bimetallscheibe | bimetal disc | bimetal |
| 3 | Kappe | cap | Obudowa |
| 4 | Stößel | plunger | Popychacz |
| 5 | Kontaktträger | contact carrier | sprężyna nośna |
| 6 | Feder | spring | sprężyna |
| 7 | Anschluss | terminal | styk |
| 8 | Rohrniet | tubular rivet | nit rurkowy |
| 9 | Niet | rivet | nit |
| 10 | Passivscheibe | passive disc | podkładka |
| 11 | Scheibenträger | disc carrier | sprężyna okrągła |

| 12 | Flachanschluss | flat terminal | styk płaski |
| 13 | Schweißanschluss | weld terminal | styk zgrzewany |
| 14 | Lötanschluss | solder terminal | styk lutowany |
| 15 | Schraubanschluss | screw terminal | styk na wkręt |

**Table 1.** Terms assigned to categories

(4) The aim of the next stage was to prepare a graphica l representation of the ontology separately for each category in order to visually show the relationships between the category and its constituent terms. The instruction (prompt 4) given to the AI to generate the visualisation was as follows:

*"Present the above ontological classification of all terms from the above tables in graphical form separately for each category. Show the relationships."*

The ontology visualisation was generated based on the above classification using the DALEE module, which is also a component of ChatGPT. For each of the eleven categories, a separate diagram was created in the form of a simple graph, in which the central node represented the name of the category and the surrounding nodes represented the assigned terms. Arrows indicated the membership relationship. This made it possible to quickly view the terminological structure of each group. Visualisations facilitate the understanding of relationships and the membership of terms, emphasising the role of hierarchy and categorisation in ontology, which corresponds to the postulates of systems theory and computer science.



**Figure 5:** Example of the ontological structure of "Tool and Device" developed in our experiment by ChatGPT

(5) The fifth stage completes the ontological structures of each category described above and presents a combined visualisation of the ontology of all concepts:

The aim was to combine all categories into a single comprehensive graphical diagram and show the relationships between categories. The requirement specified that each category and its terms were to be presented in a different colour according to the instructions (prompt 6):

*"Create a graphic visualisation for all 11 categories and their subcategories. Each category and its subcategories should be displayed in a different colour."*

This approach is particularly valuable for visual users and supports the memorisation and understanding of complex terminological structures.

The summary diagram included all categories and their assigned terms, using a separate colour scheme for each cluster. The colourful visualisation made it easier to distinguish between groups and showed the full picture of the ontology. Coloured segments helped to locate terms and quickly find connections between elements. The visualisations showed that certain categories are hierarchically related. For example, process parameters are controlled by measurements, and the results of the measurements determine the actions in the next steps. Categories of materials and substances are linked to parts through assembly activities (e.g. soldering or gluing). Operator and role categories indicate process participants, while protective equipment is linked to occupational safety.
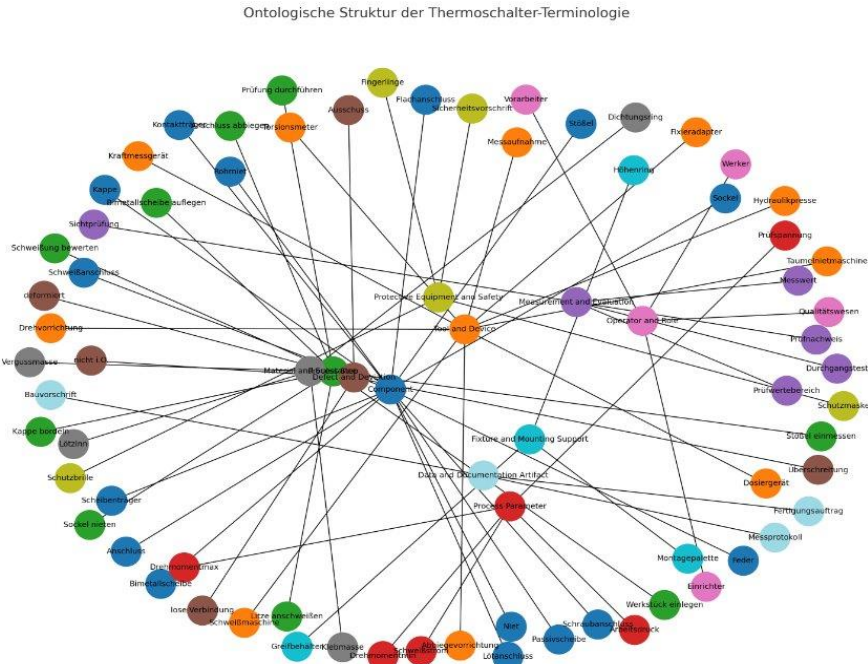


**Figure 6:** Ontological structure of terminology from technical documentation for thermal switches

## 11. Usability assessment for translators and linguists

A translator undertaking the translation of a technical text is obliged to adhere to certain rules, which are determined by many factors, both non-linguistic and specific to the translation process itself. The inherent feature and primary goal of specialised translation is its practical usefulness and functionality, as in the case of all technical documentation, equipment manuals, etc. This necessitates an understanding of the source text and a good knowledge of the field and its specific properties, which determines the quality of the final product. A very important principle of specialised translation is the aforementioned rule of using uniform and consistent terminology. Specialised translation should also be characterised by syntax that is as simple and clear as possible and a certain minimalism in stylistic devices. Just a few years ago, translators carried out the translation process completely on their own, with only a computer, Internet access and fledgling translation applications at their disposal. Traditional dictionaries remained the main tool, later replaced by translation tools. Translators therefore had to demonstrate a high degree of self-discipline and consistency in order to comply with the above rules.

AI tools are therefore a particularly important solution that not only significantly facilitates and speeds up the translator's work, but also provides them with mechanisms and procedures that ensure a an acceptable quality of translation.

The terminological ontology generated by the AI assistant is a tool that facilitates the translation and analysis of technical language. Firstly, it allows for the unambiguous assignment of a term to a category and understanding of its meaning in the context of production. This allows the translator to avoid ambiguity and maintain terminological consistency. Secondly, classification facilitates the creation of glossaries and terminology databases in CAT systems; terms can be linked to their corresponding translations and synonyms. Thirdly, visualising relationships allows you to see the connections between processes, tools and parts, which aids understanding of the source text and identifies gaps in vocabulary. Finally, ontology can be integrated with semantic tools and knowledge management systems, which facilitates the automation of technical documentation analysis and translation.

## 12. Conclusions

The capacity of GPT-4 in the area of restricted domains, and, in this case, electrotechnology, is significant, especially compared to earlier models, but it still has some weaknesses when handling highly specialized tasks in the three examined languages. When prompted to adopt a specific persona, in our case a linguist, the output presents a good quality for this particular specialized domain.

Some weaknesses refer to the accuracy with the precise equivalence of highly domain-specific terminology that was underrepresented in its training data. The model sometimes defaults to incorrect informal translations and showed lexical ontological inconsistency unless prompted with precise professional, technical information. For fields like the researched electrotechnological domain

where, similarly to other domains, accuracy is critical, its translations should always be verified by a human expert.

GPT-4's primary output is text. Therefore, it did not inherently *create* or *render* visual assets like an ontology visualisation. However, when properly prompted the LLM can then be executed by an external environment to create the visualization in a structured format and answer questions about a visualized ontology model, effectively serving as an intelligent query interface.

## Acknowledgement:

## References

Dornseiff, Franz (2003). *Der deutsche Wortschatz nach Sachgruppen.* Berlin & New York: W. de Gruyter (8., völlig neu bearb. ... Aufl.).

Hallig, Rudolf & Wartburg, Walther von (1963). *Begriffssystem als Grundlage fur die Lexikographie. Versuch eines Ordnungsschemas.* Berlin: Akademie-Verlag (2nd ed.).

Li, Ning, Ren, Liang, Liu, Zon-Ssang., Ren, Shu-hang, Xiang, Chong, Wu, Bo-yu., and Cai, Xuan . (2024). Ontology Construction Technology of Knowledge Graph in Oil and Gas Exploration and Development. In: Lin, J. (eds) *Proceedings of the International Field Exploration and Development Conference 2023. IFEDC* 2023. Springer Series in Geomechanics and Geoengineering. Springer, Singapore. https://doi.org/10.1007/978-981-97-0272-5_38

Lewandowska-Tomaszczyk, Barbara (2017). Cluster Equivalence, General Language, and Language for Specific Purposes. In: M. Grygiel (ed.) Cognitive Approaches to Specialist Language. Newcastle upon Tyne: Cambridge Scholars Publishing. 384 – 418.

Lewandowska-Tomaszczyk, Barbara (2022/2023). Emerging AI technologies: ChatGPT
challenges in contemporary university foreign language education. *Konin Language Studies*. KSJ. 11 (2). 2023. 167-186.

Pawłowski, Grzegorz (2023). The Implementation of professional language terminology in Polish production companies. In: Lewandowska-Tomaszczyk, B., Trojszczak, M. (Eds.). *Language Use, Education, and Professional Contexts*. Springer. 55-70.

Pęzik, Piotr (2014). Graph-based analysis of collocational profiles. In V. Jesenšek and P. Grzybek (Eds.). *Phraseologie Im Worterbuch Und Korpus (Phraseology in Dictionaries and Corpora). ZORA* 97 (pp. 227–243). Maribor, Bielsko-Biała, Budapest, Kansas, Praha: Filozofska Fakuteta.

Pęzik, Piotr (2016). Exploring phraseological equivalence with paralela. In: Gruszczyńska E., Leńko-Szymańska, A. (eds) *Polish-language parallel Corpora.* Warsaw: Instytut Lingwistyki Stosowanej UW, 67–81.

Radszuweit, Siegrid & Spalier, Martha (1982). *Knaurs Lexikon der sinnverwandten Worter. 20 000 Stichworter mit ihren Synonymen.* Munchen & Zurich: Droemer Knaur.

Rees, Geraint Paul. & Robert Lew (2023). The Effectiveness of OpenAI GPT generated definitions versus definitions from an English Learners' Dictionary in a lexically orientated reading task. International Journal of Lexicography XX 1-25.

Roget, Peter Mark (1852), *Thesaurus of English words and phrases. classified and arranged so as to facilitate the expression of ideas and assist in literary composition.* London etc.: Longmans, Green & Co.

Roget, Peter Mark (1984), *Roget's II. The new thesaurus.* New York: Berkley (Condensed version of the homonymous publication by Houghton Mifflin, 1980).

Schryver, de Gilles-Maurice 2023. Generative AI and Lexicography: The Current State of the Art Using ChatGPT. International Journal of Lexicography, 2023, 36, 355–387 https://doi.org/10.1093/ijl/ecad021

Sierra, Gerardo, 2000. The onomasiological dictionary: a gap in lexicography. In: Heid, U., Evert, S., Lehmann, E., and C. Rohrer (Eds.), Proceedings of EURALEX 2000 . Stuttgart University. 223-235. https://euralex.org/wp

content/themes/euralex/proceedings/

Veseli, Blerta, Singhania, Sneha, Razniewski, Simon and Gerhard Weikum. 2023. Evaluating language models for knowledge base completion, In: Proceedings of ESWC, vol. 13870 of LNCS, Springer. 227–243. doi:10.1007/ 978-3-031-33455-9\_14.

Vrolijk, Jarno, Reklos, Ioannis, Vafaie, Mahsa, Massari, Arcangelo, Mohammadi, Maryan and Sebastian Rudolph. 2022.Toward a comparison framework for interactive ontology enrichment methodologies, in: Proceedings of VOILA@ISWC, vol. 3253 of CEUR Workshop Proceedings, CEUR-WS. 41–50

## Authors' Bios:

**Barbara Lewandowska-Tomaszczyk** is professor dr habil. in Linguistics and English Language at the Department of Language and Communication at the University of Applied Sciences in Konin (Poland). Her research focuses on cognitive semantics and pragmatics of language contrasts, corpus linguistics, and information retrieval with their applications in translation studies and lexicography.

**Grzegorz Pawłowski,** of the Department of Language and Communication at the University of Applied Sciences in Konin (Poland), is a researcher and a practicing translator. He obtained his PhD in German, and his recent research focuses on the applications of ChatGPT in foreign language teaching and translation studies.