

A New Approach for Semi-automatic Building and Extending a Multilingual Terminology Thesaurus*

Aleš Horák

*Natural Language Processing Centre
Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
hales@fi.muni.cz*

Vít Baisa

*Natural Language Processing Centre
Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
vbaisa@fi.muni.cz*

Adam Rambousek

*Natural Language Processing Centre
Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
rambousek@fi.muni.cz*

Vít Suchomel

*Natural Language Processing Centre
Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
vsuchom2@fi.muni.cz*

This paper describes a new system for semi-automatically building, extending and managing a terminological thesaurus—a multilingual terminology dictionary enriched with relationships between the terms themselves to form a thesaurus. The system allows to radically enhance the workflow of current terminology expert groups, where most of the editing decisions still come from introspection. The presented system supplements the lexicographic process with natural language processing techniques, which are seamlessly integrated to the thesaurus editing environment. The system's methodology and the resulting thesaurus are closely connected to new domain corpora in the six languages involved. They are used for term usage examples as well as for the automatic extraction of new candidate terms. The terminological thesaurus is now accessible via a web-based application, which a) presents rich detailed information on each term, b) visualizes term relations, and c) displays real-life usage examples of the term in the domain-related documents and in the context-based similar terms. Furthermore, the specialized corpora

*Preprint of an article submitted for consideration in International Journal on Artificial Intelligence Tools © 2019 copyright World Scientific Publishing Company
<https://www.worldscientific.com/worldscinet/ijait>

2 Aleš Horák, Vít Baisa, Adam Rambousek, Vít Suchomel

are used to detect candidate translations of terms from the central language (Czech) to the other languages (English, French, German, Russian and Slovak) as well as to detect broader Czech terms, which help to place new terms in the actual thesaurus hierarchy.

This project has been realized as a terminological thesaurus of land surveying, but the presented tools and methodology are reusable for other terminology domains.

Keywords: Thesaurus building; terminology dictionary; domain-corpus exploitation; knowledge extraction; term extraction; DEB platform; knowledge-rich contexts.

1. Introduction

Specialists in any branch inevitably rely on domain-specific vocabulary as a basis for sharing exact terminology among professionals. Such detailed domain terminology cannot be included in general language dictionaries, which is why specialized terminology dictionaries are being built and managed. With the need to share information unambiguously in different languages, terminology dictionaries often link original terms to their translations. The taxonomic ordering of the terminology is described by means of term relations such as synonymy or hypernymy/hyponymy.^a In the system presented in this paper, information about the terms is described and visualized in a way that helps the readers (both specialists and the general public) to understand the meaning of the term and its usage in contexts.

Any human language is continuously evolving—new words and terms appear while usage and meanings evolve. This evolution is even more noticeable in specialized vocabularies.¹ A system for terminology thesauri thus needs to cope with frequent updates of the data. In this article, we present the details of a thesaurus development system which prepares the underlying information for updates automatically, and displays them during the entry editing process for terminologists' authorization and completion.

The Natural Language Processing Centre (NLP Centre) at the Faculty of Informatics, Masaryk University, in cooperation with the Czech Office for Surveying, Mapping and Cadastre (CUZK) has developed a new system for building and extending a specialized terminological thesaurus for the domain of land surveying and land cadastre, which we refer to as TeZK.^b The TeZK project consists of several tightly interconnected parts—a web-based application to create, edit, browse and visualize the terminological thesaurus, and a set of tools to build large corpora of domain oriented documents which allows for the detection of newly emerging terms, or terms missing from the thesaurus. General tools already developed by the NLP Centre for corpus building and term extraction and a platform for dictionary applications are utilized here alongside newly developed tools for extracting

^a*Synonymy* is usually used in a “weakened” form as *near synonymy* or the *see also* relation. *Hypernymy/hyponymy* refers to *broader/narrower terms*.

^bIn Czech “Tezaurus pro obor zeměměřictví a katastru nemovitost” (Thesaurus for the field of land surveying and land cadastre).

candidate translations and for identifying candidate broader terms.^c These have been developed on top of the general corpus tools. During the TeZK project, we enhanced the existing corpus tools so as to support comparable multilingual corpora. We also developed a new thesaurus web application (not limited to single domains) with new methods that interconnect the domain corpus with the terminological thesaurus. Each term in the thesaurus is supplemented with knowledge-rich context information—a term explanation, its relationships to other terms, usage examples, term translations, or specific links to related resources, such as e.g. the corresponding legislative.² In this sense, the system offers a radical improvement of the work of terminology expert groups who are in charge of organizing, constructing and managing the official terminology in their respective area.

After evaluating the TeZK project, the Czech e-Government Office decided to use it as a basis for a new official Czech registry of terminological thesauri, currently in the early development phase. In the follow-up project, the terminological thesaurus system is being updated to support easy and user friendly deployment at any organization, with the possibility to customize work processes based on specific organization requirements. Furthermore, each instance of the terminological thesaurus system will share selected data with the central registry and all other terminological thesauri.

This article is structured as follows. Section 2 contains an overview of related work regarding systems for the advanced building of terminology thesauri. In Section 3, we present the process of creating a large specialized domain corpus with the functionality of extracting new candidate terms in the selected domain. In Section 4, we describe the multiplatform web-based editor and browser application which is based on DEB, an existing dictionary writing platform now enhanced with new functions for the thesauri building system. Although the TeZK project aimed at building and managing a terminological thesaurus of the land surveying domain, our newly developed tools may be re-used for any other domain dictionary, which would facilitate the sharing of information and stimulate a general awareness of a selected domain. We also introduce the linking capabilities of the terminological thesaurus data in the system as related to Linked Open Data methodology. The final section provides an overview of its usage and outlines some possible areas for future work.

2. Related work

The idea of developing a specialized system for creating and sharing terminological thesaurus content is not new, as several tools for dictionary management are available, both free and commercial, such as Lexique Pro^d or t!Term^e. However, these

^cIn lexicography a “broader term” denotes a “hypernym,” we use these two denotations interchangeably in the text.

^d<http://www.lexiquepro.com/>

^e<http://tshwanedje.com/terminology/>

4 *Aleš Horák, Vít Baisa, Adam Rambousek, Vít Suchomel*

applications concentrate on features for dictionary compilation and presentation. We, on the other hand, are aiming at enhancing the whole methodological process by adding knowledge-rich context information for each term and automatic candidate information about new terms, their taxonomical relations and translations. For the same reason, we do not reuse any complete existing tool for building the TeZK terminological thesaurus, only include software tools developed for single tasks and enhance them to form the complete system.

Our thesaurus management system exploits our previous experience in designing and developing several applications for terminological thesaurus creation, such as the Multilingual Glossary of Fine Art Terminology with the Faculty of Fine Arts, Brno University of Technology, or the Czech-English Dictionary of Ethnological Terminology with the National Institute of Folk Culture.^f The current terminology thesaurus project considerably extends the features developed in the other projects.

2.1. Data visualization

Previous research showed that visualization and rich information significantly help dictionary users in understanding the terms and context. Graphic representations of the relationships between terms proved indispensable to users, as is the case in the DiCoInfo Visuel project³ and in EcoLexicon^{4,5,6}. We thus decided to add options that would visualize the term relationships, and a tree representation of hyponymy/hypernymy relations. Several studies also found that rich information is very helpful for students and translators in the field. Marshman's case study concluded that "Studying terms in context can also help clarify differences between concepts expressed by polysemous terms."⁷

2.2. Term extraction

Automatic term extraction provides candidate term lists as support materials for subject field experts. Corpora-based term extraction uses lexical collocability measures combined with other evidence and approaches. For example, the TExSIS system⁸ proposes an algorithm that works with bilingual parallel corpora: it finds aligned chunks of texts (using statistical word alignment) and extracts the translated terms using statistical filters. Even though the results look promising, they are limited to language pairs for which precisely aligned corpora are available.

Another approach was proposed by García-Silva et al.⁹ to automatically construct domain ontologies based on social tagging systems. Such systems (also called folksonomies) allow users to publicly share webpage links and add labels or tags used to categorize the content. However, this approach is not useful for the TeZK project because of the lack of data for the specialized domain.

We decided to integrate a technique (described in Section 3) that has been successfully included in the Sketch Engine corpus manager and verified in a ter-

^f<http://www.nul.k.cz/>

minology extraction project from patent data of the World Intellectual Property Organisation.¹⁰

2.3. *Term translation*

The TeZK system offers a list of candidate term translations based on specifically prepared domain corpora^g. The idea of extracting translation candidates from comparable corpora^h has been studied by Morin et al.¹¹, who have shown that the quality of comparable corpora might alleviate the data sparsity problem. This is also the case of the TeZK system where the selected domain is rather limited. Sadat et al.¹² proposed a system (for the Japanese-English language pair) which first extracts possible translation pairs and then filters out non-promising candidates using linguistic rules. Gu et al.¹³ used a similar approach to discover semantically similar sentences, however only within a single language.

Daille and Morin¹⁴ used contexts for aligning possible translation candidates of previously extracted monolingual multiword expressions from French-English technical documents. Lee et al.¹⁵ used an EM-basedⁱ algorithm for the extraction which required an alignment of comparable documents prior to the actual candidate extraction. They demonstrated a language independent approach on English-Chinese and English-Malay language pairs. In one part of the extraction procedure, they used co-occurrence statistics. Sorg and Cimiano¹⁶ proposed multi-lingual concept linking with the help of explicit semantic analysis, using Wikipedia categorization and cross-language links.

2.4. *Semantic relations*

Another feature of the TeZK system is the extraction of semantic lexical relations, particularly hypernyms and hyponyms, i.e. broader and narrower terms. This technique is generally used for augmenting or verifying existing lexicons and for identifying semantically related terms as proposed by Hearst.¹⁷

In the TeZK project, hypernym candidate identification is used when adding a new term to the ontology or taxonomy built within the system. Hearst identifies the lexico-syntactic patterns by bootstrapping from manually discovered patterns or existing lexicons, and deriving new rules from common syntactic environments. Hearst also argues that this technique does not work well for English meronymy/holonymy.

Snow et al.¹⁸ propose learning the patterns automatically via a logistic regression classifier trained over texts containing hypernym/hyponym word pairs from the WordNet semantic network. Banko et al.¹⁹ presented an open information extraction method, which is based on extracting occurrences of different relations

^gText corpora with documents devoted to a selected field or problem domain, see Section 3 for further details.

^hComparable corpora (as opposed to “parallel corpora”) are text corpora in different languages, whose documents talk about the same topics but are not direct translations of each other.

ⁱExpectation Maximization

using a small set of general relation patterns common to all kinds of relations and then deciding the relations by a CRF-based[‡] unsupervised extraction. The best recall and precision was achieved by a combination of supervised and unsupervised approaches. Arnold and Rahm²⁰ proposed an algorithm to extract semantic relations from Wikipedia corpora which may also prove useful for lexicon enhancement; however, the algorithm was tested only for English data.

A case study by Lefever et al.²¹ describes the HypoTerm system for hypernym detection in Dutch and English. The paper evaluated multiple approaches for relationship detection (pattern-based, morpho-syntactic analysis, statistical, WordNet-based) and discovered that the pattern-based approach provides the highest precision scores.

The TeZK system follows the pattern-based approach for the sake of higher precision and ease of maintenance within the thesaurus system. With the growth of semantic web technologies, machine learning from semantic web data may prove better results, as discussed by Rettinger et al.²²

3. Specialized Corpora and Term Extraction

The advanced automatic functionality of the TeZK system relies on a set of large text corpora in six predetermined languages containing domain-specific texts, i.e. mostly technical or popular text devoted (in this case) to topics related to land surveying, geodesy and land cadastre in general. Our previously developed web crawler, SpiderLing²³, and tools for web text cleaning²⁴ were used to build new domain corpora from publicly available online resources in Czech, English, German, French, Russian and Slovak. The process of constructing the specialized domain corpora started with the “pivot” language, Czech. Firstly, a set of main websites related to the land surveying, the cadastre of real estates, and related topics was listed. See Table 1 for details of the primary sources for the Czech corpus. We use standard corpus size measures, such as the number of documents and the number of tokens, i.e. text units like words, punctuation characters and structure tags.

Secondly, a broader set of documents from a large set of websites was obtained by our WebBootCat tool²⁵ based on the “term content” (see below) of the primary websites. Table 2 presents the number of acquired documents and tokens for all languages covered in the project. This method needs a set of seed words to search the web for relevant documents. For the seed word sets, we have used the main domain terms obtained from the publicly available terminology dictionary.²⁶ The representativeness of the created corpora and their thematic coverage of the selected areas can be seen in Table 3 with details about the document and token distribution among different sub-topics (as divided in the authoritative terminology dictionary). Non-textual and low quality content was automatically identified

[‡]Conditional Random Fields

Table 1. Primary website sources for the Czech domain corpus. Tokens are particular word, punctuation or structure tag occurrences. Unique documents refer to the number of documents after removing (near) duplicate documents. Unique tokens describe the vocabulary richness of the respective source.

Website	Docs	Tokens	U docs	U tokens
www.cuzk.cz	16,405	3,137,795	15,289	340,943
www.vugtk.cz	4,659	6,419,950	3,212	4,386,238
csgk.fce.vutbr.cz	241	77,255	198	58,561
www.kgk.cz	417	44,814	414	29,890
www.sfdp.cz	192	35,287	106	11,279
www.czechmaps.cz	94	108,506	90	98,914
www.zememeric.cz	8,634	6,100,751	6,200	2,638,308

and removed from the downloaded documents utilizing our JusText tool.²⁴ Finally, duplicate documents or paragraphs were purged with our Onion tool.

Table 2. Statistics detailing the web-crawled domain corpora for the six languages.

Language	Documents	Tokens	Web domains
English	8,149	40,225,064	4,946
French	5,326	15,789,761	3,291
German	3,373	9,744,313	2,220
Russian	2,914	19,015,734	1,770
Slovak	2,943	10,252,449	1,528
Czech	27,389	12,689,548	1,061

3.1. Automatic Extraction of New Candidate Terms

These domain corpora offer a sufficient basis for the intelligent functionality of the TeZK system. All the corpora can be continuously updated (extended) by adding new documents (or websites), which will undergo the same processing pipeline as described above, i.e. text extraction, deduplication and tokenization. The first function allows us to identify “candidate terms”, i.e. proposals to be checked by experts in the subject field and easily added to the thesaurus. The candidate terms were extracted from the domain corpora using a hybrid approach: a) the first step is a linguistically motivated rule-based extraction of noun phrases and other term patterns²⁷, and b) the second step lies in sorting all these noun phrases by relevance which is computed by comparing the relative frequency of each noun phrase in the given domain corpus with its relative frequency in a (general) reference corpus^(28,29).

The first step, i.e. identification of phrases in the corpus text, which could form a (complex) term, is based on a predefined (language-dependent) set of patterns denoted as a “term grammar.” These patterns are expressed in the form of Corpus Query Language (CQL) and describe phrases such as a (complex) noun phrase (e.g.

Table 3. Distribution of sub-topics of the resulting corpora (in percentages), per language, by the number of documents and by the number of tokens. Cadastre, cartography and geodesy are the most frequently represented topics.

language sub-topic	English		French		German	
	docs	tokens	docs	tokens	docs	tokens
cadastre	9.7	10.4	6.2	6.8	13.9	10.1
cartography	17.6	20.2	16.5	15.9	26.0	22.3
engineering surveying	3.6	11.2	7.8	6.5	12.1	9.5
theory of errors	5.3	2.6	4.4	2.9	3.6	9.0
geodesy	20.1	19.4	25.4	24.2	1.6	3.7
geoinformation	13.0	6.1	5.7	12.7	7.8	7.9
GPS system	7.7	4.1	5.0	5.2	0.8	0.6
instrumental technology	7.2	5.1	4.3	2.2	2.2	9.1
mapping	6.8	10.0	9.8	11.5	21.4	14.6
metrology	6.2	8.4	4.2	5.5	6.5	7.1
photogrammetry	3.0	2.4	10.8	6.5	4.1	6.0
	Russian		Slovak		Czech	
cadastre	6.2	3.4	15.2	10.7	15.4	14.5
cartography	11.4	9.1	16.0	21.0	21.7	21.0
engineering surveying	6.3	7.8	7.1	5.0	6.7	4.5
theory of errors	7.4	10.5	1.5	2.0	4.5	4.1
geodesy	18.2	27.1	21.0	15.4	11.2	9.0
geoinformation	16.9	14.2	6.0	5.9	3.3	10.3
GPS system	2.8	1.0	3.3	2.5	6.9	4.1
instrumental technology	4.8	6.1	5.5	1.9	6.7	3.2
mapping	7.3	4.9	12.1	18.6	12.6	13.5
metrology	12.6	9.0	2.5	2.8	8.5	11.5
photogrammetry	6.2	6.7	9.8	14.2	2.5	4.3

“*digital photogrammetric workstation*”) or a combination of a noun phrase and a prepositional phrase (e.g. “*parallactic figure with an auxiliary base*”).

In the second step, the resulting “term rank” of each identified phrase is determined by the formula

$$\text{rank}(\text{term_candidate}) = \frac{f + n}{f_{\text{ref}} + n},$$

where f is the domain corpus relative frequency of a given candidate term, and f_{ref} its relative frequency in a reference corpus. The parameter n (called simple math) can be used to fine-tune the results based on the size of the analyzed corpora and on user’s preferences. High values of n cause the algorithm to prefer more frequent phrases and vice versa. In the default setup the value of $n = 1$ is chosen. This approach allows to adapt the term extraction technique to the specific language data in cases where standard statistical methods (e.g. mutual information (MI) score, Log-Likelihood, or Fisher’s Exact Test) fail due to their assumption that the language phenomena are independent of each other – see ²⁹ for a detailed explanation.

For each language, the respective corpus of the TenTen corpus family was used as a reference corpus.³⁰ The TenTen corpus family contains very large general language corpora^k built from web.

3.2. Automatic Term Relations Identification

The methodology of the TeZK systems aims at continuous amendments of the thesaurus taxonomy using new terms. The term inclusion process is supported by two other (semi)automatic techniques: the identification of candidate hypernyms/broader terms and the candidate term translations. The technique of broader terms identification relies on two methods of automatic hypernym extraction: a pattern extraction from a domain corpus and a term similarity based approach.

Within the *pattern extraction method*, the specialized domain corpus (of the pivot language) is filtered^l to obtain a list of possible hypernym candidates, which are then ordered using the logDice similarity score:

$$\logDice(t_1, t_2) = \log_2 \left(\frac{2f_{t_1, t_2}}{f_{t_1} + f_{t_2}} \right),$$

where f_{t_1, t_2} is the number of co-occurrences of terms t_1 and t_2 .

The number of possible patterns can be generally extended without limitations. The TeZK system uses three of the most productive patterns:

- Pattern 1: *The hyponym + is/are + the hypernym,*
- Pattern 2: *The hyponym + and/or another/other/similar + the hypernym,*
- Pattern 3: *The hyponym + is/are a kind/type/part/example/way of + the hypernym.*

Although the accuracy of Pattern 1 and Pattern 2 queries is above 50% , not all successfully extracted hypernym pairs are suitable for the particular term database. For instance, some identified hypernym terms are too general to be included in the thesaurus or, vice versa, explicit hyponyms are particular instances, which are not to be included in the definitions according to the editor's decision. Another approach to finding hypernyms of a term involves searching the current term database and identifying *lexically similar terms*, e.g. “Cartesian coordinate system” and “coordinate system”. The most similar terms are expected to be good generalizations of the term, and thus either good hypernym candidates or synonym terms, which help to identify a common hypernym. The lexical similarity measure between two terms is based on the Jaccard distance of bigrams of characters with a threshold of 0.5:

^kThe sizes of the TenTen corpora range from billions of words to tens of billions of words, ergo 10^{10} words. The TenTen corpus family currently covers 31 languages.

^lThe queries are evaluated via the concordance API of Sketch Engine³¹ with the patterns specified in the formal Corpus Query Language (CQL).

$$lexsim(t_1, t_2) = \frac{|t_1 \text{ bigrams} \cap t_2 \text{ bigrams}|}{|t_1 \text{ bigrams} \cup t_2 \text{ bigrams}|}$$

Both these methods are combined in the system and the best candidates for hypernyms are available for the terminologists in the user interface in the form of a shortcut select box (see Section 4.3). The final decision as to which candidate to select or whether to input the hypernym manually is still left to a terminologist. During the development of the relation extraction module, we have evaluated 7 different patterns. We have evaluated 50 candidate relation pairs per each of the extraction patterns, where every instance was annotated as either a correct hypernym/hyponym pair or not. The most productive patterns (patterns 1 and 2) reached the accuracy of 56% and 60% respectively. Pattern 3 was much less reliable (<10% accuracy) and has been allocated a lower weight. Other patterns (e.g. “to be known/denoted as”) reached less than 5% accuracy and have thus been excluded from the final system. The term similarity method correctly identified the correct hypernym among the top three hypernym candidates in 56% of cases when measured on random terms from the database having at least one hypernym.

3.3. *Translation Candidates Extraction*

All entries in the TeZK terminological thesaurus are organized around the pivot language (Czech), which is also used for definitions of particular term meanings. Each term entry may contain translations to equivalent terms in five foreign languages: English, German, French, Slovak and Russian. When editing a term, the TeZK system suggests translation candidates based on the domain corpora of all these languages. Since these corpora are not parallel, i.e. they do not contain translations of documents but rather documents from the same domain, it is not possible to use standard word-alignment tools such as GIZA++³² to identify direct phrase translations: the corpora are not aligned at the sentence level so standard co-occurrence statistics cannot be used.

We exploit the fact that equivalent terms share similar contexts. Consider two terms: “European Parliament” in English and “Parlement européen” in French. In English documents, the term co-occurs with verbs like “revokes”, “decided”, “informed”, “opposed” and in French, “Parlement européen” co-occurs (collocates) with “révoque”, “informé”, “notifié”, “opposé”. In many cases, the collocations are translations of each other. It is possible to extract the collocates for each monolingual candidate term and measure how much these collocates overlap for all possible pairs of English and French terms using an English-French translation dictionary. It may seem paradoxical to require an existing translation dictionary to extract translation candidates but it is important to note that the dictionary does not need to contain special terminology: the collocations are usually frequent, core language items as is the case with the verbs above, thus they appear in even modest-size dictionaries. And since we limit the selection on both sides to the identified “term

candidates”, see 3.1, the computations are not overloaded by very frequent collocating words (e.g. frequent verbs) which have a lot of translation possibilities.

For each pair of terms, i.e. a source language (Czech) term and a candidate term in a target language, the number of collocates which are translation equivalents according to a given dictionary, is computed. These numbers are then used for ordering all pairs of terms and the top 10 candidates for each source language term are selected. The accuracy of this method is definitely not sufficient for automatic translation, thus the system uses the resulting lists as a reservoir of related terms (and thus translation candidates) in the target language. To evaluate the method on the existing data, the resulting candidate translations were compared to the existing term translations.^m In 34% of Czech terms, the existing English translation appeared in the top 20 automatically identified translation candidates. For the other languages this ratio is 40% for German, 21% for French, 24% for Russian and 47% for Slovak. These numbers do not seem adequate for helping the editor, but please note that in this evaluation possible lexical variations and synonyms were not taken into account as the annotations for them were not available. The results thus express only exact matches, but (different) possible translation candidates were often available in the first 10 items.

Examples of correctly identified candidate translations are:

- *katastr* (cadastre): land cover, cadastre, land information, land, land survey, land registry, land management, information system, georeference, control.
- *měření* (measurement): measurement, point positioning, dual frequency, position, fitting, distance measurement, land cover, ellipsoidal height, vertical control, vertical angle.
- *poloha* (position, location): position, location, measurement, height, ellipsoidal height, positioning, orthometric height, navigation system, ground control, accuracy.

Examples of where the method did not yield correct candidates are:

- *katastrální úřad* (cadastral office): registration, indexing, geospatial information, geodetic datum, analysis, factor, bench mark, system, scan line, process control.
- *atlas* (atlas): input, format, control, system, conformity assessment, cadastre, raster image, quality control, monitoring, data, alteration.

For the purpose of the evaluation, which was measured on the existing translations of 2,972 to 8,439 terms (based on the respective language pair), the statistical bilingual dictionaries were built from the parallel corpus OPUS2³³ and the DGT-TM translation memory³⁴.

^mThe current term dictionary contained from 3,070 to 4,575 translations from Czech terms to terms in the other five languages.

4. Thesaurus Management Application

The main entry point of the TeZK system for the user is a web-based application accessible from all major web browsers without the need to install any new components. The application offers different modes of operation depending on the type of user. This includes

- searching and browsing term information including term usage examples, term relations or the term hierarchy,
- term entry editing, and
- full terminological thesaurus management with the processing of both new terms added by terminologists as well as automatically extracted new candidate terms.

The whole application is based on our general dictionary browsing and editing development platform, DEB, which is briefly presented in the following section.

4.1. *Dictionary Editor and Browser Platform*

Exploiting our experience of several lexicographic projects, we have designed and implemented a universal dictionary writing system that can be exploited in lexicographic applications to build and exploit both small and large lexical databases. The system is called Dictionary Editor and Browser, or DEB.³⁵ Since 2005, DEB has been employed in more than 20 international research projects. Examples of applications based on the DEB platform include the Czech Lexical Database³⁶ with detailed information on more than 213,000 Czech words, or the complex lexical database, Cornetto, combining the Dutch WordNet, an ontology, and an elaborate lexicon.³⁷ Current ongoing projects include the Pattern Dictionary of English Verbs tightly interlinked with corpus evidence,³⁸ the Dictionary of Family names in Britain and Ireland³⁹ providing a detailed investigation into over 45,000 surnames to be published by Oxford University Press, and a compilation of the Dictionary of the Czech Sign Language with extensive use of multimedia recordings to present the signs visually. The DEB platform is based on a client-server architecture, which provides a raft of benefits. All the data are stored on a server and a considerable part of its functionality is also implemented on the server, which permits the client application to be very lightweight. The server part is built from small, reusable parts, called servlets, which allow a modular composition of all services. Each servlet provides different functionalities such as database access, dictionary search, morphological analysis or connections to corpora.

The overall design of the DEB platform focuses on modularity. The data stored in a DEB server can be saved in any kind of structural database (or several different databases) and the results are combined in the answers to user queries without the need to use specific query languages for each data source. The main data storage is currently provided by the Sedna XML database,⁴⁰ which is an open-source native XML database providing XPath and XQuery access to a set of document containers.

The user interface, that forms the most important part of a client application, usually consists of a set of flexible complex forms which dynamically cooperate with the server parts. A DEB client application can be implemented in any programming language which allows interaction with the DEB server using the available server interfaces.ⁿ Details regarding the TeZK client application are presented in Section 4.3.

The main assets of the DEB development platform are:

- All the data are stored on a server and a considerable part of the functionalities is also implemented on a server, allowing the client application to be very lightweight.
- It provides very good tools for (remote) team cooperation so that data modifications are immediately seen by all users. The server also provides authentication and authorization tools.
- A DEB server may offer different interfaces using the same data structure. These interfaces can be reused by many client applications.
- Homogeneity of the data structure and presentation. If an administrator commits a change in the data presentation, this change will automatically appear in every instance of the client software.
- Easy integration with external applications via API (Application Programming Interface).

4.2. *Initial Thesaurus Data*

Although the main aim of the TeZK terminological thesaurus development lies in managing and publishing the authoritative specialized terminology and its updates both to experts in the subject field and the general public, the terminological thesaurus also contains a broad vocabulary of related terms. Users may even search for unofficial terms, and thanks to the term relations and detailed information on the source of a given term, users can easily explore all related terms and navigate to the preferred “official” term variant.

To build the initial TeZK terminological thesaurus data covering a broad domain vocabulary, we have combined several resources. In the first stage, the current Czech authoritative terminology dictionary^o (which contained 3,937 term definitions and translations, but did not offer a taxonomy network) was combined with a hyper/hyponymic tree of 6,800 entries^p (containing hyponymic relations, but with-

ⁿClient applications communicate with servlets using HTTP requests in a manner similar to a popular concept in web development called AJAX (Asynchronous JavaScript and XML) or using the W3C standard SOAP protocol. The data are transported over HTTP in a variety of formats: RDF, XML documents, JSON-encoded data, plain-text formats, or marshalled using SOAP.

^o*Terminologický slovník zeměměřičství a katastru nemovitostí* (The Dictionary of geodesy, cartography and cadastre) is published electronically at <http://www.vugtk.cz/slovník> and processed by the Terminology commission of the Czech Office for Surveying, Mapping and Cadastre.

^pAlso provided by the Terminology commission of the Czech Office for Surveying, Mapping and Cadastre.

out any detailed information about terms) and by 450 candidate terms extracted from the Czech domain corpus.

The first two resources were available in HTML with a mostly fixed entry structure, but still leaving portions of text in an unstructured format. We have thus implemented a flexible import module for the TeZK system, which is able to import both structured and unstructured data to the terminological thesaurus. The system allows the administrators to configure the formatting rules of imported data with HTML, CSV and TXT formats currently supported. After the initial parsing of the document using the formatting rules, the import module detects duplicate, misspelled or close terms, normalizes abbreviated forms, and cleans the data (e.g. correcting the punctuation). In the next step, the imported data are converted to a unified XML format for database storage. The import module was employed to combine the two resources (terminology dictionary and the hypernymic tree) resulting in combined term entries containing both detailed term information and term relations. Generally, the module enables easy expansions of the terminological thesaurus employed, e.g. in importing terms from the regularly updated Registry of Territorial Identification (RUIAN)^q.

The resulting terminological thesaurus also contains suggestions of candidate terms automatically extracted from the Czech domain corpus^r. These suggestions are inserted into a separate taxonomic category, each entry including the information regarding the source and the reliability of the term. These terms are then subject to approval or disapproval by subject field specialists using the automatic candidate functionality for emplacing the term into the correct position in the taxonomy and enriching it with the foreign language terminology translations. Table 4 shows details regarding the current size of the terminological thesaurus.

Table 4. The TeZK terminological thesaurus size statistics.

	Number of entries
Total entries	8,427
Hyper/hyponymic relations	8,827
Explanations provided	4,117
Entries categorized to domain	3,905
Total number of translations	24,973
English translation	9,073
German translation	4,513
Slovak translation	3,751
Russian translation	3,068
French translation	4,568

^q<http://www.cuzk.cz/ruian/>

^rAs Czech plays the role of the "pivot" language in the TeZK terminological thesaurus.

4.3. Entry Editing

The TeZK terminological thesaurus editing module is designed and implemented as a client application, with the DEB server providing the database and management backend. The editing interface is a multi-platform web application accessible in any modern browser utilizing open-source technologies^s. The standardized application interface allows for an easy integration of third-party applications that can be built upon the terminological thesaurus data. The interface provides all the functions needed to work with the data (e.g. search queries, browsing the terminological thesaurus structure and detailed entry information, entry creation and updates...). Two standard remote access techniques are available supporting modern web-service standards: REST/JSON⁴¹ and WSDL^t. One of the intended use cases is the integration into the official public Geoportal website^u, where the terminology is to be used for the document metadata and categorization.

The screenshot shows the TeZK web interface. At the top, there is a navigation bar with 'TeZK', 'Thesaurus', 'Information', and 'Contact' tabs, along with a search bar and a dropdown menu for 'all terms'. Below the navigation bar, there are buttons for '+New term', 'Export', and 'Import'. The main content area is titled 'souřadnicový systém (coordinate system)' and includes a 'Terminologický' section with a definition and three numbered points. Below the definition, there are sections for 'Translations' (listing 'coordinate system', 'système de coordonnées (m)', 'Koordinatensystem (s)', 'система координат', and 'súradnicový systém'), 'Domains' (listing 'geodézie'), and 'Relations' (showing a hierarchy: 'souřadnice (sig.)' -> 'parametry' -> 'globální navigační družicový systém' and 'vyjádření prostorových referencí souřadnicemi' -> 'kartometrie' -> 'dokumentace').

Fig. 1. Browsing the terminological thesaurus, with detailed information for one term.

^sJQuery (<http://jquery.com>) is used for communication and SAPUI5 (<https://sapui5.netweaver.ondemand.com/>) libraries for the graphic interface. The client and the server communicate using a standardized interface over HTTP with the data encoded in the JSON format.

^t<http://www.w3.org/TR/wsd1/>

^u<http://geoportal.cuzk.cz/>

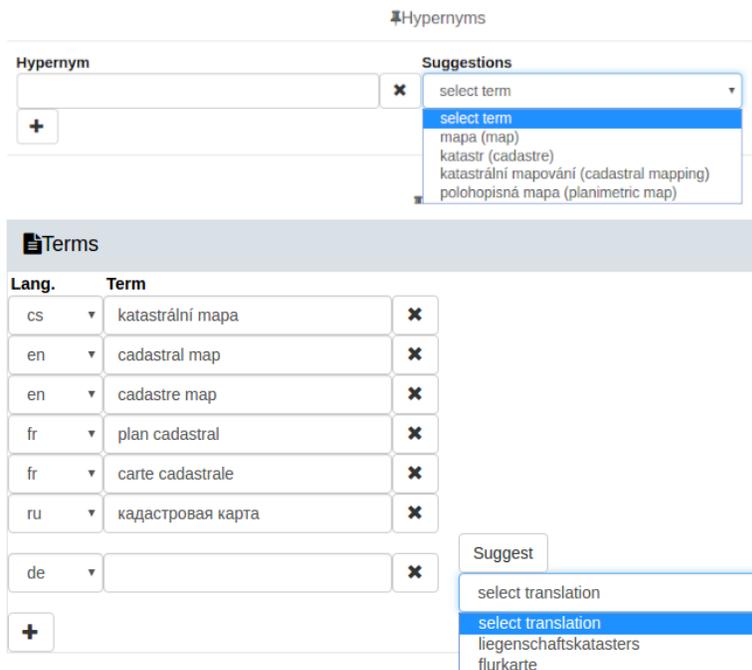


Fig. 2. Automatic broader candidate terms and translations offered within the editor, for the term *katastrální mapa* (cadastre map).

The TeZK terminological thesaurus application offers a graphical interface for browsing the hypernymic tree (see Figure 1). There are several possible visualization techniques for the taxonomy and the TeZK browser works with an expanding multi-level tree. Although it may not display all the relations in a proper graph form (as a term can have more than one hypernym), the expanding tree is the most intuitive representation for users. If a term has more than one hypernym, it is presented multiple times in the tree structure (with the same unique identification number). The automatic hypernym candidates, obtained via the technique described in Section 3.2, are offered to the editor directly within the editing form (see Figure 2).

Each term presentation includes all meaning explanations, translations, and accepted variants. In case of new terms, the translation candidates are selected from the TeZK domain corpora (see Section 3.3 for details) and offered as a list of proposals in the editing form (see also Figure 2). When more sources are incorporated into the terminological thesaurus, the reliability of each source and revision history is presented to the users. Source reliability follows the rating scale of the authoritative body—the most reliable being the terms authorized by the terminology committee, followed by terms used in scientific journals, with terms made up by the general

Usage examples		
exception of Giza, until 2003 there were no general scheme is overlaid on the IKONOS	satellite images	available of the greater part of the pyramid
European Union (EU), which mandates aerial and cartographic and geodetic data, data from using GPS data.	satellite images	of subsidies linked to agricultural land etc. On the other hand the cadastral data
portion of spectrum.	satellite images	- really at the heart of GIS revolution
Image data.	satellite images	often in pixels 100" x 100" on a side and aerial photographs to scanned maps
spatial units, derived from low accurate , or can be identified from low accurate	satellite images	.
outing (multi-date and multi-resolution)	Satellite images	, GIS techniques andground data.
Bibipur and Adi Badris seen clearly on trace of theSarasvati Nadi overlaid on the ofSarasvati channels which are self-evident on	Satellite image	is shown in Fig. 2a. Sarasvati Nadi is
STUDIESPalaeochannels generally appear on the	Satellite image	as serpentine drainage course with high
flowed thousands ofyears ago.	Satellite images	of February 2004 and the topomaps of 1969 and other scientific data, have contributedto
on Yamuna.	Satellite images	of the palaeochannels, geological and sediment
whichultimately met the thirst of millions.	Satellite images	in possession of the ISRO and ONGC have
elevationduring the Late Quaternary uprising.	Satellite images	reveal that had thislandmass not risen,
topographical maps as well as Remote Sensing	Satellite image	maps so that we can learn to navigate together

Fig. 3. Corpus evidence for context usage examples of the selected term *družicový snímek* (satellite image).

Terms used in similar contexts									
geographer	+maker	+researcher	+scientist	+astronomer	+designer	+artist	+scholar	surveyor	
+engineer	+writer	+author							

Fig. 4. Related words (words used in similar contexts) to the selected term—kartograf (cartographer). Words in blue boxes are already present in the TeZK terminological thesaurus.

public at the bottom of the scale. Users and third-party applications may decide which sources or terms they prefer to work with.

Definitions created by specialists are advantageously supplemented with real-world examples of their usage in contexts based on the domain corpora (see Figure 3) or related terms, which are obtained by comparing the term context words and offering a list of terms used in similar contexts (see Figure 4). The “relatedness” of two terms is computed by comparing their collocational tables of words that appear most often in the surrounding contexts of each term.

4.4. *Linked Open Data*

The term Linked (Open) Data (LOD) refers to a methodology for publishing and interlinking remote structured data via online references. This methodology was proposed by Berners-Lee.⁴² The importance of Linked Open Data is acknowledged for example by the European Union, funding projects like LOD21^v (a large integrated project to develop tools, standards and management methods for Linked Open

^v<http://lod2.eu/>

Data) or Open Data Portal^w (catalogue of data available for reuse). The methodology introduces five requirements for data to conform with the LOD principles, which include availability under open license, encoding in open machine-readable structured format, or interlinking with other LOD resources.

Since the Czech Cadastre Office aims to publish the data compiled in the TeZK terminological thesaurus for public use, the system needs to support the Linked Open Data methodology. The DEB platform provides the appropriate tools, but the decisions on how to release the data lie with the author. The DEB platform functionality generally enables publishing all documents as genuine Linked Open Data. All the terminological thesaurus data is accessible via standard web service protocols (WSDL/SOAP), encoded in a standardized XML structure (RDF/SKOS) interlinked to their respective sources. The only requirement outside the system is thus the decision about the (open) license for the terminological thesaurus data reuse. With data in XML format, it is possible to use various techniques for complex data querying, using for example the refinement framework implemented by Bao et al.⁴³

5. Conclusions and Future Work

We have presented the details of a new lexicographic system, which builds extensively on automatic language processing techniques to enable creating and regularly updating a domain terminology. The TeZK system offers a unique combination of a well-prepared lexicographic database system with experimental techniques exploiting information derived from large domain corpora in several languages. The system uses the corpora to provide an automated functionality for extracting terminology candidates based on the available documents as well as offering proposals to incorporate new terms into the terminological thesaurus taxonomy and supplement them with correct term translations into five foreign languages.

The terminology of the land surveying domain is maintained by the Terminology commission of the Czech Office for Surveying, Mapping and Cadastre that approves new terms and their official translations. The Commission will use the TeZK system to review automatically extracted terms and user-submitted terms. With the help of the automatic candidate selection functions, updated versions of the official terminology will be published periodically. Since each term will be rated, public users may quickly check if a term is officially recommended and find the right term for the task. The thesaurus contains knowledge-rich contextual information for each term (definitions, corpus examples, word relations), making it an inestimable resource for all kinds of works related to the specific domain. Thanks to the TeZK system support for Linked Open Data methodology and standardized public API application interface, the published resources and the terminological thesaurus functionality can be seamlessly integrated into third-party applications.

^w<http://data.europa.eu/euodp/en/data>

In the future work, we will further investigate the techniques for candidate translations identification. We plan to employ distributional semantics models as another measure for ordering and classifying the candidate terms in the target language.

The TeZK system will also serve as a basis for the Czech e-Government registry of terminological thesauri, currently in the early development phase. In a follow-up project, the terminological thesaurus system is being updated to support easy and user friendly deployment at any organization (both government organizations, and unofficial interest associations), with the possibility to customize work processes based on specific organization requirements. Furthermore, each instance of the terminological thesaurus system will share data with the central registry and all other terminological thesauri. During 2019, the whole system will be tested with two terminological thesauri – a thesaurus of geospatial information terminology, and the Ministry of Interior law terminology thesaurus.

References

1. R. Fischer, *Lexical change in present-day English: A corpus-based study of the motivation, institutionalization, and productivity of creative neologisms* (Gunter Narr Verlag, Tübingen, 1998).
2. I. Meyer, Extracting knowledge-rich contexts for terminography, *Recent advances in computational terminology* **2** (2001) p. 279.
3. B. Robichaud, Logic Based Methods for Terminological Assessment, in *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, eds. N. Calzolari, K. Choukri, T. Declerck, M. U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis (European Language Resources Association (ELRA), Istanbul, Turkey, may 2012).
4. P. Faber, P. León-Araúz and A. Reimerink, *Representing Environmental Knowledge in EcoLexicon, in Languages for Specific Purposes in the Digital Era*, eds. E. Bárcena, T. Read and J. Arús. (Springer International Publishing, Cham, 2014), Cham, pp. 267–301.
5. P. León-Araúz, A. San Martín and P. Faber, Pattern-based word sketches for the extraction of semantic relations, in *Proceedings of the 5th International Workshop on Computational Terminology (Computerm2016)*2016, pp. 73–82.
6. P. León-Araúz and A. S. Martín, The EcoLexicon Semantic Sketch Grammar: from Knowledge Patterns to Word Sketches, in *Proceedings of the LREC 2018 Workshop "Globalex 2018 - Lexicography & WordNets"*2018.
7. E. Marshman, Enriching terminology resources with knowledge-rich contexts: A case study, *Terminology* **20**(2) (2014) 225–249.
8. L. Macken, E. Lefever and V. Hoste, Taxis: Bilingual terminology extraction from parallel corpora using chunk-based alignment, *Terminology* **19**(1) (2013) 1–30.
9. A. García-Silva, L. J. García-Castro, A. García and O. Corcho, Building domain ontologies out of folksonomies and linked data, *International Journal on Artificial Intelligence Tools* **24**(02) (2015) p. 1540014.
10. A. Kilgarriff, M. Jakubíček, V. Kovář, P. Rychlý and V. Suchomel, Finding Terms in Corpora for Many Languages with the Sketch Engine, in *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics* (The Association for Computational Linguistics, Gothenburg, Sweden, 2014), pp. 53–56.

20 Aleš Horák, Vít Baisa, Adam Rambousek, Vít Suchomel

11. E. Morin, B. Daille, K. Takeuchi and K. Kageura, Brains, not brawn: The use of "smart" comparable corpora in bilingual terminology mining, *ACM Transactions on Speech and Language Processing* **7** (October 2008) 1:1–1:23.
12. F. Sadat, M. Yoshikawa and S. Uemura, Bilingual Terminology Acquisition from Comparable Corpora and Phrasal Translation to Cross-language Information Retrieval, in *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics - Volume 2 ACL '03*, (Association for Computational Linguistics, Stroudsburg, PA, USA, 2003), pp. 141–144.
13. Y. Gu, Z. Yang, G. Xu, M. Nakano, M. Toyoda and M. Kitsuregawa, Exploration on efficient similar sentences extraction, *World Wide Web* **17**(4) (2014) 595–626.
14. B. Daille and E. Morin, French-english Terminology Extraction from Comparable Corpora, in *Proceedings of the Second International Joint Conference on Natural Language Processing IJCNLP'05*, (Springer-Verlag, Berlin, Heidelberg, 2005), pp. 707–718.
15. L. Lee, A. Aw, M. Zhang and H. Li, EM-based Hybrid Model for Bilingual Terminology Extraction from Comparable Corpora, in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters COLING '10*, (Association for Computational Linguistics, Stroudsburg, PA, USA, 2010), pp. 639–646.
16. P. Sorg and P. Cimiano, Exploiting wikipedia for cross-lingual and multilingual information retrieval, *Data & Knowledge Engineering* **74** (2012) 26 – 45, Applications of Natural Language to Information Systems.
17. M. A. Hearst, Automatic Acquisition of Hyponyms from Large Text Corpora, in *Proceedings of the 14th Conference on Computational Linguistics - Volume 2 COLING '92*, (Association for Computational Linguistics, Stroudsburg, PA, USA, 1992), pp. 539–545.
18. R. Snow, D. Jurafsky and A. Y. Ng, Learning syntactic patterns for automatic hypernym discovery, *Advances in Neural Information Processing Systems* **17** (2004).
19. M. Banko, O. Etzioni and T. Center, The tradeoffs between open and traditional relation extraction., *ACL* **8** (2008) 28–36.
20. P. Arnold and E. Rahm, Automatic extraction of semantic relations from wikipedia, *International Journal on Artificial Intelligence Tools* **24**(02) (2015) p. 1540010.
21. E. Lefever, M. Van de Kauter and V. Hoste, HypoTerm: detection of hypernym relations between domain-specific terms in Dutch and English, *Terminology* **20**(2) (2014) 250–278.
22. A. Rettinger, U. Lösch, V. Tresp, C. d'Amato and N. Fanizzi, Mining the semantic web, *Data Mining and Knowledge Discovery* **24**(3) (2012) 613–662.
23. V. Suchomel and J. Pomikálek, Efficient Web Crawling for Large Text Corpora, in *Proceedings of the Seventh Web as Corpus Workshop (WAC7)*, eds. A. Kilgarriff and S. Sharoff2012, pp. 39–43.
24. J. Pomikálek, *Removing Boilerplate and Duplicate Content from Web Corpora*, PhD thesis, Masaryk University, Faculty of Informatics, 2011.
25. M. Baroni, A. Kilgarriff, J. Pomikálek and P. Rychlý, WebBootCaT: instant domain-specific corpora to support human translators, in *Proceedings of EAMT 2006 – 11th Annual Conference of the European Association for Machine Translation* (The Norwegian National LOGON Consortium and The Departments of Computer Science and Linguistics and Nordic Studies at Oslo University (Norway), Oslo, 2006), pp. 247–252.
26. P. Hánek, *Terminologický slovník zeměměřictví a katastru nemovitostí (in Czech, The Terminology Dictionary of Geodesy, Cartography and Cadastre)* (Výzkumný ústav geodetický, topografický a kartografický, v.v.i., 2012).
27. M. Jakubíček, A. Horák and V. Kovář, Mining phrases from syntactic analysis, in *International Conference on Text, Speech and Dialogue, TSD 2009* (Springer, 2009),

- pp. 124–130.
28. A. Kilgarriff, Comparing corpora, *International journal of corpus linguistics* **6**(1) (2001) 97–133.
 29. A. Kilgarriff, Simple maths for keywords, in *Proceedings of the Corpus Linguistics Conference* (University of Liverpool, Liverpool, 2009).
 30. M. Jakubíček, A. Kilgarriff, V. Kovář, P. Rychlý and V. Suchomel, The TenTen Corpus Family, in *7th International Corpus Linguistics Conference CL 2013* (Lancaster, 2013), pp. 125–127.
 31. A. Kilgarriff, V. Baisa, J. Bušta, M. Jakubíček, V. Kovář, J. Michelfeit, P. Rychlý and V. Suchomel, The Sketch Engine: ten years on, *Lexicography* **1**(1) (2014) 7–36.
 32. F. J. Och and H. Ney, A systematic comparison of various statistical alignment models, *Computational Linguistics* **29** (March 2003) 19–51.
 33. J. Tiedemann, News from OPUS-A collection of multilingual parallel corpora with tools and interfaces, in *Recent Advances in Natural Language Processing* **5**2009, pp. 237–248.
 34. R. Steinberger, A. Eisele, S. Klocek, S. Pilos and P. Schlüter, DGT-TM: A freely available translation memory in 22 languages, in *Proceedings of the 8th international conference on Language Resources and Evaluation (LREC'2012)*2012, pp. 454–459.
 35. A. Horák, K. Pala and A. Rambousek, The Global WordNet Grid Software Design, in *Proceedings of the Fourth Global WordNet Conference* (University of Szegéd, Szegéd, Hungary, 2008), pp. 194–199.
 36. A. Horák and A. Rambousek, PRALED – A New Kind of Lexicographic Workstation, in *Computational Linguistics: Applications*, eds. A. Przepiórkowski, M. Piasecki, K. Jassem and P. Fuglewicz (Springer, 2013) pp. 131–141.
 37. A. Horák, P. Vossen and A. Rambousek, A Distributed Database System for Developing Ontological and Lexical Resources in Harmony, in *Lecture Notes in Computer Science: Computational Linguistics and Intelligent Text Processing* (Springer-Verlag, Haifa, Israel, 2008), pp. 1–15.
 38. I. E. Maarouf, J. Bradbury, V. Baisa and P. Hanks, Disambiguating verbs by collocation: Corpus lexicography meets natural language processing, in *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, eds. N. C. C. Chair, K. Choukri, T. Declerck, H. Loftsson, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis (European Language Resources Association (ELRA), Reykjavik, Iceland, may 2014).
 39. P. Hanks, R. Coates and P. McClure, Methods for Studying the Origins and History of Family Names in Britain, in *Facts and Findings on Personal Names: Some European Examples* (Acta Academiae Regiae Scientiarum Upsaliensis, Uppsala, 2011), pp. 37–58.
 40. A. Fomichev, M. Grinev and S. Kuznetsov, Sedna: A Native XML DBMS, *Lecture Notes in Computer Science* **3831** (2006) p. 272.
 41. R. T. Fielding and R. N. Taylor, Principled Design of the Modern Web Architecture, *ACM Transactions on Internet Technology* **2** (May 2002) 115–150.
 42. T. Berners-Lee, Design Issues: Linked Data (2006).
 43. Z. Bao, Y. Yu, J. Shen and Z. Fu, A query refinement framework for XML keyword search, *World Wide Web* **20** (2017) 1–37.