Anna WRÓBLEWSKA, Grzegorz PROTAZIUK, Robert BEMBENIK
Warsaw University of Technology, Institute of Computer Science

Teresa PODSIADŁY-MARCZYKOWSKA
Polish Academy of Science, Institute of Biocybernetics and Bioengineering

# LEXO: A LEXICAL LAYER FOR ONTOLOGIES – DESIGN AND BUILDING SCENARIOS[1]

**Summary**. The paper describes a lexical layer for ontologies and scenarios of its population. This layer extends lexical descriptions of the given ontology. It defines terms and their lexicalized meanings (given with contexts) associated with elements in the ontology. Additionally, it provides links to commonly used lexical knowledge resources.

**Keywords**: ontology engineering, lexical layer for ontology, ontology population

# WARSTWA LEKSYKALNA ONTOLOGII I SCENARIUSZE JEJ BUDOWY

**Streszczenie**. W artykule jest opisana warstwa leksykalna (słownikowa) ontologii oraz scenariusze jej budowy. Ta warstwa rozszerza opisy słowne danej ontologii, definiuje znaczenia leksykalne za pomocą słów występujących w kontekście danego terminu oraz wiąże je z elementami ontologii. Możliwe jest dołączenie znaczeń leksykalnych zdefiniowanych za pomocą znanych zasobów, jak WordNet czy Wikipedia.

**Słowa kluczowe**: inżynieria ontologii, warstwa leksykalna ontologii, metody budowania ontologii

## 1. Linguistic Enrichment of Ontologies

The knowledge bases which are suitable for intelligent automatic processing as well as using in Natural Language Processing (NLP) methods require linking with the ontologies (which define semantics) and lexical resources. Such knowledge bases are crucial for realizing the vision of the Semantic Web, where knowledge sharing, information integration, interoperability and semantic adequacy are main requirements.

One of the shortest dictionary definitions of semantics is the study of meaning [1]. The more complex explanation of this would lead to a relationship that maps words, terms and written expressions into common sense and understanding of objects and phenomena in the specific domain (given by defined domain ontologies). It is worth mentioning that objects, phenomena and relationships between them are to a large extent language independent. Generally, an ontology concept/semantic layer can be defined as a domain ontology that holds the knowledge model for a specific field. To this end we accept an assumption of independence of the semantic layer from the language, so that the same semantic network of concepts can be mapped to multiple languages. It is useful in automatic translations or cross-lingual searches. Such ontology mapping to a specific language can be named linguistic layer for the given ontology. It can be realized by different dictionaries or other specialized frameworks designed particularly for proper linguistic representation, e.g. another ontology.

A linguistic layer has mappings (relations) between words and phrases (that are lexical entries) in written natural language and concepts from a given ontology (which define semantics). The proposed linguistic layer has relationships between lexical entries, e.g. the linguistic relationships (synonymy, homonymy, antonymy, meronymy etc.). A linguistic layer may include all or some of the objects of the linguistic research [2]: a lexicon (a vocabulary), morphology (internal structure of words), syntax (word formations in sentences), lexical semantics (including lexical relations), phraseology (describing context in which terms are used), discourse analysis, and areas associated with speech: phonetics, phonology and pragmatics. The usefulness of the above linguistic objects in applications varies with different level of extent. For example, lexicons should be operable in all types of search applications, especially multilingual ones. Phonology or phonetic aspects may be usable in applications for speech analysis or reading texts to aid for the blind.

A linguistic layer associated with a given ontology (semantic layer) allows determining easily which concepts included in the ontology occur in texts written in a natural language. Additionally, the phrases extracted from texts may be linked to their semantics. Thus, the following relations between terms and concepts can be expressed: polysemy (a term can refer

to many concepts) and synonymy (many terms refer to the same concept). The idea is illustrated in Fig. 1.

The linguistic layer provides a possibility to understand text (written in a natural language) and produce text (in a natural language as well) using selected concepts. A separation between a semantic layer and a linguistic layer gives an opportunity for plugging a variety of ontologies specific for certain application, express them in several languages or localize them in specific domains. Moreover, it makes the semantic/concept layer independent of its linguistic realization.
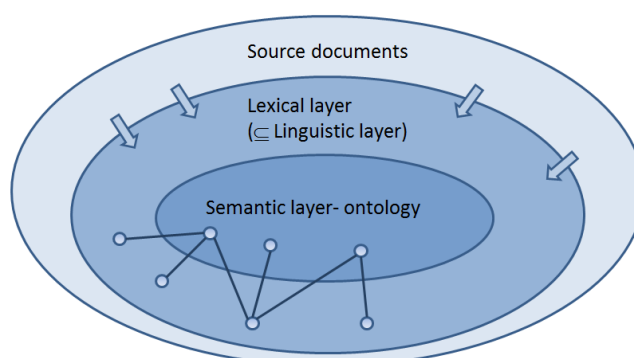


Fig. 1.  Overall idea of associations of semantic and lexical layers: representation of semantics in text and text considering its meaning (circles in the semantic layer), not words (circles in the lexical layer)

Rys. 1.  Ogólna koncepcja powiązania warstw semantycznej i leksykalnej: reprezentacja semantyki w tekście oraz tekstu, biorąc pod uwagę jego znaczenie (kółka w warstwie semantycznej), a nie słowa (kółka w warstwie leksykalnej)

The rest of the paper is organized as follows. In section 2 we provide a the short review of linguistic layer models. In section 3 our lexical layer is presented. Section 4 gives a short comparison of our model with other models. In sections 5 an example of populated lexical layer is presented. Section 6 outlines scenarios of populating the layer. Section 7 summarises the paper.

## 2. Linguistic Layer Models – State of the Art

Lexical knowledge contained in terminological and linguistic resources can be expressed in various ways. Over the last few years, a number of works aiming at interfacing ontologies and lexical resources have been initiated. The following initiatives may be considered the most important: Lexical Markup Framework (LMF, [3]), Linguistic Information Repository (LIR, [4]), and Lexicon Model for Ontologies (LEMON, [5]).

LMF is a meta-model that provides a standardized framework, allowing for the creation and use of computational lexicons. The LMF meta-model is organized into packages. From our per-

spective the most relevant are: Core Package, Morphology Package, NLP Morphological Patterns Package, NLP Syntax Package, Constraint Expression Package, NLP Semantic and NLP Multilingual Notations Package. The core package contains the basic elements of the model and their dependencies. The central entity in the LMF meta-model is the *Lexical Resource*, which has an associated *Global Information* object capturing administrative details and information related to encoding. The *Lexical Resource* consists of several language-specific *Lexicons*. A *Lexicon* then comprises *Lexical Entries* (i.e., words, multi-word entities such as terms and idioms, etc.) which are realized in different *Forms* and can have different meanings (*Senses*). LMF was used as a  basis for the construction of other linguistic frameworks.

LIR is a model inspired from the LMF for associating lexical information with OWL ontologies. The main goal of LIR is to provide a model allowing for the enrichment of ontology with a lexico-cultural layer for capturing the language-specific terminology used to refer to certain concepts in the ontology. The LIR model has focused on multilingual aspects, as well as, on capturing specific variants of terms (such as abbreviations, short forms, acronyms, transliterations, etc.) which are all modeled as subclasses of the property *hasVariant*. To account for multilinguality, the classes *LexicalEntry*, *Lexicalization*, *Sense*, *Definition*, *Source* and *UsageContext* are all associated with a certain *Language* to model variants of expression across languages. It also allows to document the meaning of certain concepts in different cultural settings.

LEMON is a simplified version of the previous frameworks (LMF, LIR, LexInfo [6], SKOS) with a strong focus on usability in information extraction.

## 3. The Lexical Layer Model

In this section we introduce our proposal of a lexical layer model – LEXO. The goal was to design a practical framework for expressing lexical knowledge of a given semantic layer. It should be usable for knowledge engineers and linguists as well as for automatic approaches. In the lexical layer we include such linguistic information as lexicon, lexical semantics and phraseology that may be defined by contexts or links to external resources such as WordNet, Wikipedia or DBPedia. The assumption of a simplicity of the model causes that our lexical layer does not include any representation of morphology and syntax.

The lexical layer in LEXO describes terms and lexicalized semantic meanings associated with elements of a semantic layer (ontology). The lexicalized semantic meanings can be expressed with the use of words or phrases in a natural language. The meanings may be given

by contexts including terms associated with the meanings using various relations, foremost linguistic ones.

### 3.1. Main Modules

Our lexical layer is composed of three main modules (Fig. 2):

- **LexicalKnowledgeSource** – used to link meaning descriptions with lexical knowledge sources, e.g. Wikipedia, DBPedia, WordNet.
- **LexicalLayerElement** – the core lexical layer representing words and phrases, their lexical meanings and contexts.
- **OntologicalLayer** – a gateway between a semantic layer of an ontology and a lexical layer.
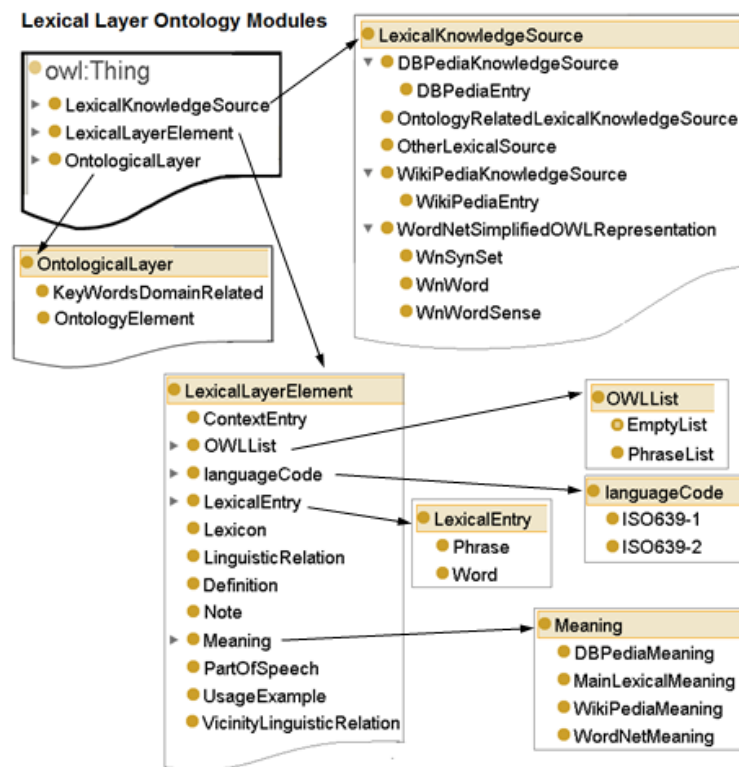


Fig. 2. The main modules of the lexical layer (a panel in the top left corner). Arrows in the figure indicate class hierarchy of the lexical layer

Rys. 2. Główne moduły warstwy leksykalnej (panel w lewym górnym rogu). Strzałki wskazują hierarchię klas w warstwie słownikowej

### 3.2. Linking to a Semantic Layer

Class *KeyWordsDomainRelated* groups terms describing generally a domain of the semantic layer. It is an auxiliary class which is helpful in searching elements of the ontology related to a given word. For example in the scientific community domain, described in section 5, we can define terms: "scientific", "academic", "science", "research".

Class *OntologyElement* realizes a reified relation from the lexical layer to a semantic layer. It allows a user to characterize all types of elements of a semantic layer, e.g. properties, instances, classes, and annotations. Each instance of *OntologyElement* has to have association with at least one instance of *LexicalEntry* or, in more advanced solutions, with an instance of *Meaning*. It may be associated with objects from both classes. Additionally, the *OntologyElement* class may indicate the proper name of an element from a semantic layer using object property *hasOntologyElementName*. Necessary conditions for the *OntologyElement* class are shown in Fig. 3.
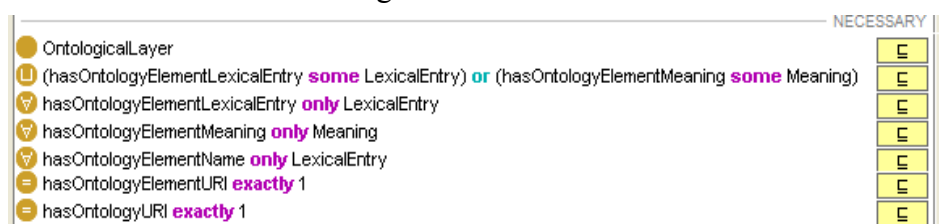


Fig. 3.   Necessary conditions for the *OntologyElement* class (visualized in Protégé editor v. 3.4.7)
Rys. 3.   Warunki konieczne dla klasy *OntologyElement* (wizualizacja w edytorze Protégé v. 3.4.7)

### 3.3.  Lexical Entries

*LexicalEntry* class represents particular words or compound phrases (e.g. names). Our assumption was to represent only base form of lexical entries. However, different variations of the forms can be attached as simple data properties (subhierarchy of *lexicalForm* property). The difference between words and compound terms is modeled by means of the classes *Word* and *Phrase*. A phrase also has its base form saved as a data property. Additionally, it can be decomposed into a list of words. Some not important words, e. g. particles or prepositions can be omitted in that list.

### 3.4.  Main Lexical Meanings

*MainLexicalMeaning* is a the most important class of the whole model (Fig. 4). It defines a lexicalized sense of one instance of *LexicalEntry* class. The main meaning is defined by means of context (described below in this section). The class *MainLexicalMeaning* is a reified relation between a lexical entry and a context. It defines also the quantity parameters: the probability and the priority number. The probability parameter indicates probability of occurrence of a given lexical entry in the meaning defined by a given instance of *MainLexicalMeaning* class. The priority number denotes the position of a given meaning in the sequence of meanings associated with a given lexical entry, starting from the most probable one.

A context is defined as a set of phrases or words (instances of *ContextEntry* class) co-occurring in any kind of linguistic relation with the defined instance of *LexicalEntry*. They are

associated with the main meaning within the lexical relations, e.g. antonymy, synonymy, equivalence, broader or narrower terms. The proximity between a term represented by an instance of *LexicalEntry* and terms from a context associated with this instance is defined as the special reified relation – *VinicinityLinguisticRelation* class. The strength of proximity may be expressed by the data property *hasCoefficient*.
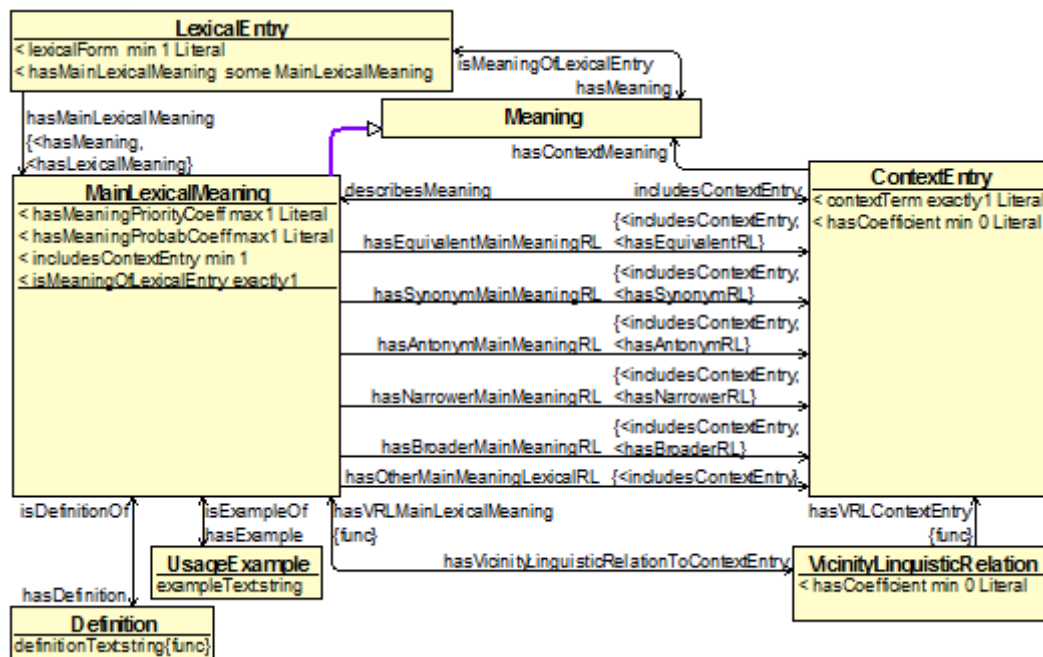


Fig. 4.   UML diagram illustrating definition of *MainLexicalMeaning* class and a context
Rys. 4.   Diagram UML, ilustrujący definicję klasy  *MainLexicalMeaning* i kontekstu

### 3.5.  Links to Lexical Knowledge Resources

Meanings may be also defined as links to commonly used lexical knowledge sources, e.g. Wikipedia or WordNet. Wikipedia-based Meaning (*WikiPediaMeaning* class) is expressed in the model as a connection to a Wikipedia entry. This entry has its own name (Wikipedia term), URI and Wikipedia disambiguation pages (listing various-sensed articles associated with the same term), related links (hyperlinks or links in See-also sections indicating Wikipedia articles about related or often confused terms), equivalent links related to redirects in Wikipedia and category Wikipedia pages. These links can additionally help in methods of defining lexical meaning by analyzing Wikipedia pages.

*WordNetMeaning* places our lexical meaning in the WordNet structure. We defined simplified WordNet structure and its lexical relations (Fig. 5). The meanings in our lexical layer can be extended to other knowledge resources available on the Internet, e.g. BabelNet, Wiktionary.

Additionally, we model hierarchy of lexical relations  to gather all such relations used in our framework (used in defining *WordNetMeaning*, *MainLexicalMeaning, WikiPediaMeaning* and the simplified WordNet structure).

## 4. Comparison with Other Linguistic Layer Models

LMF, as a first standard, lacks any references to a semantic layer of ontology. It is an extended and highly complicated model. The next models (LIR and LEMON) overcome these drawbacks. They connect semantic layer in a different way: LIR links *LexicalEntries* and LEMON links *Senses* with the ontology. All the models provide extended representation of language morphology and syntax.

In contrast to the above linguistic frameworks, our approach does not provide any representation of morpho-syntactic properties, but offers expanded structure for defining lexical meanings. For example , we define the base lexical meanings as a set of contexts, which is a convenient representation for  NLP applications Our framework defines also representations of links to various commonly used lexical knowledge resources: WordNet, Wikipedia, DBPedia and others. Similarly to LEMON, our framework represents compound terms as a list of words and it associates lexicalized meanings with ontology elements. Our model gives the opportunity to link *LexicalEntries* with the semantic layer for simpler models (and for concepts name representation) and to associate *Meanings* with elements from a semantic layer for more advanced models.
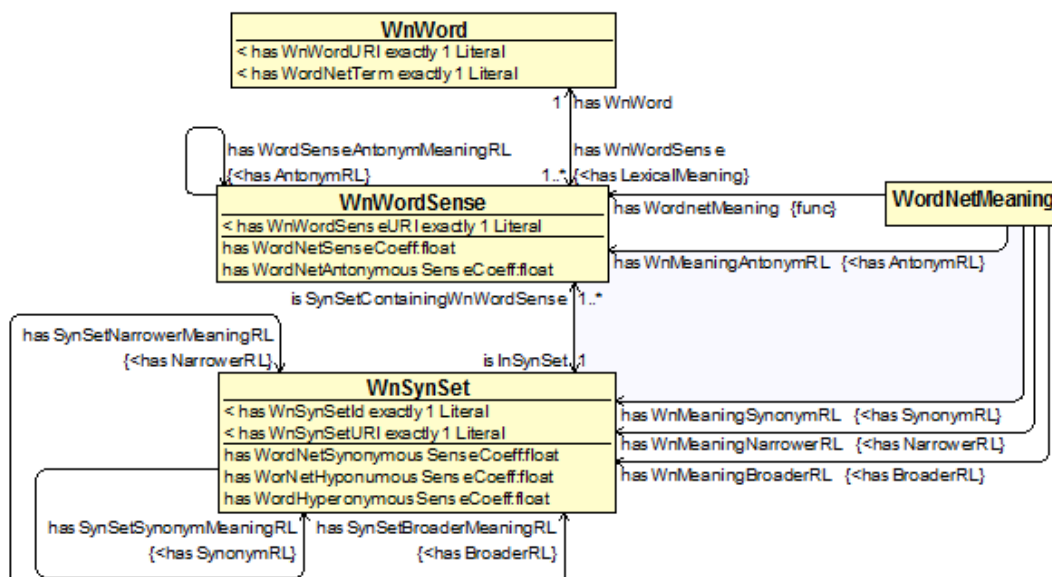


Fig. 5.   Representation of a meaning based on simplified structure of WordNet
Rys. 5.   Reprezentacja znaczenia oparta na uproszczonej strukturze WordNet-u

## 5. Use Case of the Lexical Layer

Given a domain ontology of scientific community with a class defining persons and their place in science [7] we can define lexical layer for the ontology containing different expressions about a person, e.g. a *scientist,* an *entrepreneur* (Fig. 6) as an example. A given term (an instance of the class *LexicalEntry*) is related to the instance of *MainLexicalMeaning* class associated with a set of instances of the *ContextEntry*. The set expresses a context, in which the term is used within the given meaning (as a person in science).

In the Fig. 6. with a notion Person represented in a semantic layer two lexical meaning are associated, namely: *scientist_LM* indicating meaning of the scientist term and *entrepreneur_LM* indicating meaning of the entrepreneur term. The context of the former consists of phrases: *lecturer*, *is an author*, *delivers a seminar*, whereas the context of the latter consists of the phrases: *delivers a seminar*, *sponsor*, manage.
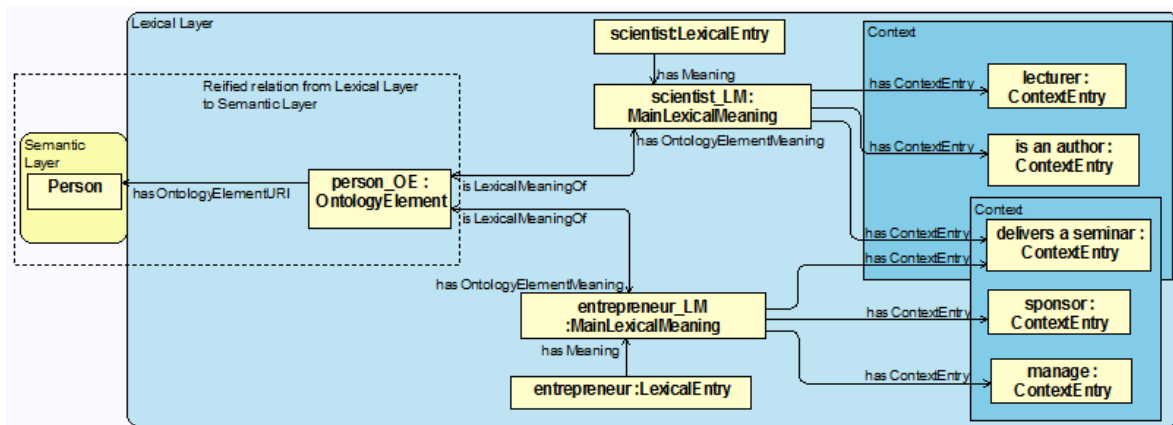


Fig. 6.   Example of a lexicalized description using lexical layer for scientific community ontology
Rys. 6.   Przykład opisu w warstwie leksykalnej dla ontologii społeczności naukowej

## 6. Scenarios for Building a Lexical Layer

In this section we present exemplary scenarios of building lexical layer for a given semantic layer of a domain ontology. Although, in scenarios the candidate entries for a lexical layer are created automatically, the final evaluation and acceptance of these entries should be done by an expert.

### 6.1.  Building a Lexical Layer for a Given Semantic/Concept Layer

In this scenario a lexical layer is built based on an input ontology and a text repository from the domain covered by that ontology. In the method the techniques used for discovering frequent sets and association rules are adapted. Actions in the scenario:

For each element of the ontology semantic layer *SE*

1.  Create a lexical entry *LE* in the lexical layer (LEXO).  The name of that entry is the name of the *SE* element (fragment of element URI). Additionally, extract *rdf:labels* of SE and add them as synonymous entries in the new created context of the *LE*. (Compound names should be divided into separate words and incorporated in LEXO as phrases.)

2.  Create collection CN of names of elements being in any kind of relation with the SE object. For example: when the analysed element is a concept, the associated elements are among others: its object properties and names of data properties.

3.  Discover maximal frequent sets MFSets. Each set consists of a name of the considered LE and not empty collection of names from CN. The sets are discovered in paragraphs extracted from documents stored in the input text repository.

4.  Calculate the support of the name of LE and generate association rules in the following form:

    $$\{LE.\text{name}\} \rightarrow mfSet \setminus \{LE.\text{name}\}$$
    where $mfSet \in$ MFSets

5.  For each association rule generated at step 4 which has the lift parameter greater than one create a lexical meaning of *LE* which context includes all names from consequence of the rule. The probability of occurrence of the name of LE in this context is equal to the confidence parameter of the rule.

### 6.2.  Adding Synonyms and Extending Contexts Based on Analysis of Text Corpus

In this scenario we utilize an assumption that synonyms in texts do not appear together in a sentence, but they appear quite frequently in similar contexts. It means synonyms are often used with the same words. Having the pairs of terms that do not co-occur, we can define a similarity measure for the contexts of the terms, such as in [8]. In this experiment we use text documents concerning the same domain and LEXO built in scenario 6.1. Actions in the scenario for generating pairs of terms that are likely to be synonyms and for adding to contexts frequently co-occurring terms:

1.  Text corpus is preprocessed (tokenization, sentence splitting, lemmatization) and tagged with parts of speech and some cleaning is done (e.g. removing stop-words).

2. The *Apriori* based algorithm for finding frequent itemsets is executed.

3. For every frequent term that is in field of our interest (the lexical entries in their meanings from LEXO), frequent itemsets containing that word are found. The terms frequently co-occurring are considered as a *temporary context* of a word.

4. Finally synonymy measure is computed for every pair of words with use of their temporary contexts. Based on this, a decision is taken whether the pair can candidate for synonymy or not.

5. The frequent itemsets contain words or phrases frequently co-occurring with the analyzed terms from LEXO. The terms frequently co-occurring are considered as candidate contexts can be attached to *MainLexicalMeanings* through *VincityLexicalRelation* in LEXO. (Firstly, we have to check their similarity and associations with other terms existing as a context in particular meanings in LEXO.)

### 6.3. Adding Information from WordNet

In this scenario we utilize a structure of WordNet (shown in Fig. 5) and collect WordNet data associated with our lexicalized meanings prebuilt in LEXO (in scenario 6.1). Actions in this scenario are as follows:

1. We search through WordNet for lexical entries defined in LEXO.

2. For each *MainLexicalMeaning* we have a lexical entry and an associated context AC (a set of instances of *ContextEntry* class). We identify a *WnWord* as the analyzed lexical entry and *WnSynSet* as the pointer to a synset containing most terms from the AC. We incorporate into LEXO the chosen *WnWord*, *WnSynSet* and *WnWordSense* with their properties (URI and coefficients). Additionally, we can incorporate also all the entities that are in linguistic relations (narrower, broader, antonym, synonym) with the chosen word, synset and word sense or all senses of the chosen word.

### 6.4. Building from Wikipedia, Disambiguation

In this scenario we use Wikipedia as a source of vast valuable descriptive knowledge – an extended thesaurus with semantic relations between concepts represented as articles. A structure of Wikipedia, besides definitive articles, contains [9]: article redirects (generally linking equivalent terms to the main article), article links (representing relationships between articles, having not identified types and strength), categories (a hierarchy providing broader and narrower terms), disambiguation pages (containing a set of links to articles about different senses of the term) and "see also sections" or "listing sections" defining related (in the broader sense) articles and potential instances of the article main title, respectively.

In this scenario we have to identify articles from Wikipedia regarding particular meanings of elements (concepts and instances) from the semantic layer. In this scenario we search through Wikipedia articles, use LEXO built in scenario 6.1 and a list of defined words (LRW) related to the whole domain of the ontology. Actions in the scenario for finding representation of meanings in Wikipedia are as follows:

1. We search in Wikipedia structure for lexical entries connected with meanings (*MainLexicalMeanings*) associated with concepts (or instances) from the semantic layer.

2. When we find only one Wikipedia article without a disambiguation, we can check if the LRW and the context terms (associated with related *MainLexicalMeaning*) is used within the article text.

3. When we are redirected to other article title we can incorporate it in LEXO as equivalent name in Wikipedia and add *WikipediaEntry* in LEXO for the article we were redirected to.

4. When we find several articles or are redirect to a disambiguation page, we have to resolve the suitable article, checking frequencies of terms associated with the analysed meaning and with the whole ontology (LRW), similarly as in point 2).

5. Additionally, we can place the main article in a hierarchy of Wikipedia articles throughout Wikipedia categories (Wikipedia broader and narrower senses). We can add associated articles to the main one with mutual cross-links between articles (Wikipedia related terms).

## 7. Summary

In the paper we presented LEXO – a lexical layer model for ontologies. In this model we emphasize linguistic meanings of words or phrases and the linguistic relations between them rather than  other linguistic property of words such as morphology or syntax. This approach causes that the proposed model is especially useful in algorithms dedicated to discovering knowledge from text repositories. We outlined several such methods for population of the lexical layer in which the knowledge from different resources: domain text repositories, a given ontology (semantic layer), Wordnet and Wikipedia is explored.

**BIBLIOGRAPHY**

1.    Nielson H. R., Nielson F.: Semantics with Applications, a Formal Introduction. John Wiley & Sons, Chicester, England 1995.

2.  Meyer C. F.: Introducing English Linguistics. Cambridge University Press, 2009.

3.  The LMF Working Group. Language resource management, Lexical Markup Framework (LMF). Technical Report ISO/TC 37/SC 4 N453 (N330 Rev.16), ISO, 2008. http://www.lexicalmarkupframework.org/.

4.  Linguistic Information Repository, http://mayor2.dia.fi.upm.es/oeg-upm/index.php/en/downloads/63-lir.

5.  LEMON ontology, http://www.monnet-roject.eu/Monnet/Monnet/English/Navigation/DemosAndDownloads?init=true.

6.  Buitelaar P., Cimiano P., Haase P., Sintek M.: Towards Linguistically Grounded Ontologies. The Semantic Web: Research and Applications. LNCS, Vol. 5554, Springer, Heidelberg 2009, p. 111÷125.

7.  Wróblewska A., Podsiadły-Marczykowska T., Bembenik R., Protaziuk G., Rybiński H.: Methods and Tools for Ontology Building, Learning and Integration – Application in the SYNAT Project. Bembenik R., Skonieczny L., Rybiński H., Niezgódka M. (eds.): Intelligent Tools for Building a Scientific Information Platform. SCI, Vol. 390. Springer, Heidelberg, 2012, p. 121÷152.

8.  Rybinski H., Kryszkiewicz M., Protaziuk G., Jakubowski A., Delteil A.: Discovering Synonyms based on Frequent Termsets. Rough Sets and Intelligent Systems Paradigms, Springer, LNAI, Vol. 4595, 2007, p. 516÷525.

9.  Stoutenburg S., Kalita J.: Extracting Semantic Relationships between Wikipedia Articles. Proc. 35th International Conference on Current Trends in Theory and Practice of Computer Science, Czech Republic, 2009.

**Omówienie**

W artykule został opisany proponowany model warstwy lingwistycznej – LEXO – łączącej teksty napisane w języku naturalnym i warstwę semantyczną – daną ontologię dziedzinową (1). Podano krótką charakterystykę istniejących modeli takiej pomocniczej warstwy: LMF, LIR, LEMON. Przedstawiono moduły modelu (2) oraz sprecyzowano jego najważniejsze aspekty: połączenie warstwy leksykalnej z semantyczną (3), klasę *LexicalEntry*, opisującą terminy z języka naturalnego, oraz klasę definiującą podstawowe znaczenie leksykalne (4) i inne reprezentacje znaczeń przy użyciu popularnych zasobów leksykalnych, tj. WordNet (5),

Wikipedia. Podano przykład zastosowania warstwy LEXO (6) oraz naszkicowano scenariusze wypełnienia danymi takiej warstwy leksykalnej.

**Addresses**

Anna WRÓBLEWSKA: Institute of Computer Science, Warsaw University of Technology, Nowowiejska 15/19, 00-665 Warszawa, Poland, A.Wroblewska@ii.pw.edu.pl.

Grzegorz PROTAZIUK: Institute of Computer Science, Warsaw University of Technology, Nowowiejska 15/19, 00-665 Warszawa, Poland, G.Protaziuk@ii.pw.edu.pl.

Robert BEMBENIK: Institute of Computer Science, Warsaw University of Technology, Nowowiejska 15/19, 00-665 Warszawa, Poland, R.Bembenik@ii.pw.edu.pl.

Teresa PODSIADŁY-MARCZYKOWSKA: Institute of Biocybernetics and Bioengineering, Trojdena 4, 02-109 Warszawa, Poland, tpodsiadly@ibib.waw.pl.