

Human Associations Help to Detect Conventionalized Multiword Expressions

Natalia Loukachevitch

Lomonosov Moscow State University
Leniskie Gory,1
Moscow, Russia
louk_nat@mail.ru

Anastasia Gerasimova

Lomonosov Moscow State University
Leniskie Gory,1
Moscow, Russia
anastasiagerasimova432@gmail.com

Abstract

In this paper we show that if we want to obtain human evidence about conventionalization of some phrases, we should ask native speakers about associations they have to a given phrase and its component words. We have shown that if component words of a phrase have each other as frequent associations, then this phrase can be considered as conventionalized. Another type of conventionalized phrases can be revealed using two factors: low entropy of phrase associations and low intersection of component word and phrase associations. The association experiments were performed for the Russian language.

1 Introduction

A lot of approaches have been proposed for automatic extraction of idioms, collocations, or multiword terms from texts as potential candidates for inclusion in lexical or terminological resources (Bonial et al., 2014; Gelbukh and Kolesnikova, 2014; Pecina, 2010; Piasecki et al., 2015).

However, developers of computational resources need clear guidelines for the introduction of phrases into their resources. Special instructions on introducing multiword terms exist for constructing information-retrieval thesauri (ANSI/NISO, 2005). Developers of WordNet-like thesauri, a very popular type of resources, discuss the problem of introducing multiword expressions in their resources in several works (Maziarz et al., 2015; Piasecki et al., 2015; Vincze and Almasi, 2014). For example, it is supposed that wordnets have to include only lexicalized concepts as synsets (Miller, 1998). However, Agirre et al., (2006) stress that boundaries of lexicalization are very difficult to draw. Bentivogli and Pianta

(2004) argue that there is a necessity to include non-lexicalized phrases into wordnets.

Multiword expressions comprise a broad scope of phrases including idiomatic expressions, noun compounds, technical terms, proper names, verb-particle and light verb constructions, conventionalized phrases, and others (Calzolari et al., 2002; Sag et al., 2002; Baldwin and Kim, 2010). For some of these constructions, such as idioms, it is evident that they should be included in computational lexicons. But for many of other expressions, for example, conventionalized phrases, it is not easy to make a decision about the necessity of their inclusion. To distinguish a multiword expression, it is important to analyze if it has any "idiosyncrasies", which can be lexical, syntactical, semantical or statistical.

Conventionalized phrases have statistical idiosyncrasy and usually only one approach is proposed in literature to distinguish such phrases from other compositional phrases. This is so-called substitutionability test, which shows if the phrase components can be easily substituted with their synonyms (Sag et al., 2002; Farahmand et al., 2015; Farahmand and Henderson, 2016; Pearce, 2001; Senaldi et al., 2016).

In this paper, we show that there are at least two more types of statistical idiosyncrasy (and related tests) to distinguish conventionalized expressions:

- association idiosyncrasy when components of a phrase are highly associated with each other, and
- relational idiosyncrasy when a phrase has lexical associations that significantly differ from the associations of its component words; usually it means that the phrase denotes a specific entity or process with a set of its own properties and relations.

We provide evidence for these types of phrase idiosyncrasy in association experiments in Russian, in which we asked Russian native speakers what associations they had for phrases and their component words. We have found that the human association experiment is a very efficient tool to detect conventionalized phrases with high accuracy. To the best of our knowledge, this is the first attempt to use human associations for distinguishing conventionalized phrases.

The structure of the paper is as follows. In Section 2 we consider types of phrase idiosyncrasy. Section 3 describes the specificity of RuThes thesaurus, from which we take phrases for the experiments. Section 4 presents the association experiment and its results. In Section 5 we test embedding models on their capability to distinguish conventionalized phrases. Section 6 reviews related work concerning approaches of annotating compositionality/non-compositionality/conventionalization of noun phrases.

2 Types of Idiosyncrasy of Multiword Expressions

Multiword expressions are phrases that have some specificity (idiosyncrasy). Because of this, it is useful to collect them and store in lexicons and thesauri (Calzolari et al., 2002; Sag et al., 2002; Baldwin and Kim, 2010).

The idiosyncrasy can be lexical when a component of a phrase appears only within this phrase (Baldwin and Kim, 2010). It can be syntactical when the syntactic behavior of a phrase differs from usual (for example, fixed word order). Semantical idiosyncrasy can be revealed when the meaning of a phrase cannot be inferred from the meanings of its components. If a phrase has one of the above-mentioned types of idiosyncrasy it can be called a lexicalized expression (Sag et al., 2002; Baldwin and Kim, 2010).

Statistical idiosyncrasy presupposes that the components of a phrase co-occur more often than expected by chance. Besides, the frequency of phrases with statistical idiosyncrasy is much higher than the frequency of the phrase with one component changed to its near-synonym (*weather forecast* vs. *weather prediction*), as the result of the substitutionability test (Sag et al., 2002; Farahmand and Henderson, 2016). Phrases with statistical idiosyncrasy (often called *convention-*

alized phrases) can be syntactically and semantically compositional.

In many cases conventionalized phrases are difficult to distinguish. For example, one of the often mentioned conventionalized phrase *traffic lights* looks fully compositional. However, if we examine the meaning of this phrase, we can see that the denoted entity can be categorized as a road facility; it has signals; it is usually constructed on road intersections; it is needed for regulating road traffic, etc. This means that the phrase *traffic lights* has thesaurus relations with the corresponding words (facilities, road, signals, regulation) that cannot be inferred from the meanings of its component words *traffic* and *lights*.

A lot of similar examples can be found. Compositional *seat belt* has relation to the safety concept. Food courts are usually located in shopping centers, and therefore compositional phrase *food court* has relation with the *shopping center* concept, etc. These relations can be very useful in such NLP applications as textual entailment.

Thus, we can suppose that conventionalized phrases have not only statistical idiosyncrasy, but also *relational idiosyncrasy*, which can be revealed easier than using the substitutionability test. The same idiosyncrasy can be found in unclear cases of possible lexicalized expressions.

In (Mel'čuk, 2012) so-called quasi-idioms are discussed. According to Mel'čuk, a phrase *AB* is a quasi-idiom or weak idiom iff its meaning: 1) includes the meaning of both of its lexical components, neither as the semantic pivot, and 2) includes an additional meaning *C* as its semantic pivot. Mel'čuk (2012) gives an example of *barbed wire*, which is an obstacle, but neither *barbed* nor *wire* are obstacles. Thus, it seems that *semantic pivot* in this case is the hypernym relation, that cannot be inferred from the phrase component words. It means that the quasi-idiom is a subtype of relational idiosyncrasy.

In this paper we show that this relational idiosyncrasy can be found in association experiments with native speakers. Besides, we can also reveal the association idiosyncrasy of conventionalized phrases in these experiments.

3 RuThes Thesaurus as a Source of Conventionalized Expressions

For the present work, we utilized multiword expressions included in the Russian-language

thesaurus RuThes¹ (Loukachevitch and Dobrov, 2014). The RuThes thesaurus is a linguistic ontology for natural language processing, i.e. an ontology, where the majority of concepts are introduced on the basis of actual language expressions.

RuThes has considerable similarities with WordNet: the inclusion of concepts based on senses of real text units, representation of lexical senses, detailed coverage of word senses. At the same time, the differences include attachment of different parts of speech to the same concepts, formulating names of concepts, attention to multiword expressions, the set of conceptual relations, etc.

In particular, the developers of the RuThes thesaurus have special rules for including phrases that appear compositional into the thesaurus. Such phrases are introduced if they have specificity in relations with other single words and/or expressions (Loukachevitch and Lashevich, 2016). The following subtypes of these expressions can be considered:

- A phrase is a synonym to a single word; for example, земельный участок (landing lot) is a synonym to word земля (land), or a phrase has a frequent abbreviation: заработная плата – зарплата (employee wages);
- A phrase has a synonymous phrase and this fact cannot be simply inferred from the components of the phrase: мобильный телефон (mobile phone) – сотовый телефон (cell phone);
- A phrase generalizes several single words. Such phrases as транспортное происшествие (transport accident) or учебное заведение (educational institution) often look compositional but they have a very important function of knowledge representation: they gather together similar concepts;
- A phrase has relations that do not follow from its component words. For example, the compositional phrase дорожное движение (road traffic) has numerous relations with other phrases that cannot be inferred from its components, for example, hyponyms (left-hand traffic, one-way traffic), related concepts (car accident, traffic jam), etc.

Thus, phrases from RuThes without evident non-compositionality were selected for the association experiment in order to understand correlations between choice of phrases made by experts and associations of native speakers.

4 Association Experiment

For the experiment, we took two-word noun phrases (*Adjective + Noun* and *Noun + Noun-in-Genitive*) that have high frequency in Russian newswire text collections.

The multiword expressions were of two main groups. The first group (Thesaurus group) included multiword expressions from the RuThes thesaurus. We chose phrases that either look fully compositional (*increase of prices*) or that have one of components is used in a known (=described in dictionaries) metaphoric sense. This group contained 15 phrases. Another group of phrases comprised fully compositional noun phrases not included in the thesaurus, for example, *end of January, mighty earthquake, result of work*, etc. The non-thesaurus group contained 36 phrases.

We asked respondents (mainly university students) to think of single-word associations to noun phrases. In a separate experiment, we collected associations to the component words of the same phrases. We wanted to understand if the collected associations can serve as a base for distinguishing thesaurus phrases from non-thesaurus phrases (and as a consequence, conventionalized phrases from non-conventionalized). Twenty six native speakers gave their associations for the thesaurus phrases and twenty nine respondents participated in the experiment with non-thesaurus phrases. Forty seven people gave associations for single words.

The study was conducted via Google Forms. The respondents were asked to provide single-word associations. However, some participants could think only of multiword expressions. Such associations were also taken into account. Table 1 contains examples of obtained associations and their frequencies for some thesaurus phrases.

From the associations obtained, we calculated the following characteristics (Tables 2, 3):

- entropy of answers for single words and phrases (currently, only entropy of phrase associations was found useful and included in the tables);

¹http://www.labinform.ru/pub/ruthes/index_eng.htm

- intersection between associations of component words and phrase associations (columns Ph1 and Ph2 in Tables 2, 3); and
- number of times when one component word served as an association of another component word (columns A12 and A21 in Tables 2, 3).

Table 2 contains the results for the thesaurus phrases, and Table 3 shows partial results for the non-thesaurus phrases.

We can see that for thesaurus phrases, the components are associated with each other more often than for non-thesaurus phrases. The average value of such associations for thesaurus phrases is 10 times greater than for non-thesaurus phrases. For some thesaurus phrases, both components are highly connected with another component. Withing non-thesaurus phrases, such frequent mutual associations were not found.

source	associations	freq
w1: земельный (landing)	участок (lot) вопрос (issue)	38 2
w2: участок (lot)	земля (land) дача (dacha) полицейский (police) дорога (road) дом (house)	11 11 4 3 2
phrase: земельный участок (landing lot)	дача (dacha) дом (house) надел (allotment)	12 2 2
w1: повышение (increase)	должность (post) зарплата (wages) работа (job)	8 7 6
w2: цена (price)	ценник (price-tag) стоимость (cost) высокая (high) качество (quality)	5 4 3 2
phrase: повышение цен (increase of prices)	инфляция (inflation) кризис (crisis) нефть (oil)	12 5 2

Table 1: Examples of the most frequent associations for thesaurus phrases and its components

Therefore, we think that mutual associations between phrase components are an important sign of phrase **conventionalization**. It seems that such phrases are stored as single units in the human memory. In our case such conventionalized

Phrase	A12	A21	Ph1	Ph2	Entr
транспортное происшествие (transport accident)	0	7	0	6	2.16
учебное заведение (education institute)	1	8	1	1	2.52
программное обеспечение (software program)	13	14	6	1	2.77
повышение цен (increase in prices)	0	0	0	0	2.85
земельный участок (landing lot)	38	13	0	14	2.89
квадратный метр (square meter)	10	20	0	1	3.22
электронная почта (electronic mail)	6	12	3	4	3.27
дорожное движение (road traffic)	0	2	0	0	3.33
заработная плата (employee wage)	18	10	8	2	3.42
главный герой (main hero)	5	1	0	3	3.56
медицинская помощь (medical aid)	0	5	0	4	3.58
торговый центр (shopping center)	26	0	1	0	3.62
лента новостей (news feed)	18	0	1	1	3.79
мобильный телефон (mobile phone)	26	12	3	7	3.81
температура воздуха (air temperature)	4	0	6	1	3.81
Average	11	6.27	1.93	2.93	3.24

Table 2: Results of association experiments for the thesaurus phrases

phrases included: программное обеспечение (*software program*), земельный участок (*landing lot*), квадратный метр (*square meter*), электронная почта (*electronic mail*), заработная плата (*employee wages*), мобильный телефон (*mobile phone*), лента новостей (*news feed*), and торговый центр (*shopping center*).

Besides, we found that the average level of entropy (4.07) of phrase associations is much higher for non-thesaurus phrases than for thesaurus phrases (3.24). This means that associations of thesaurus phrases are more concentrated, more motivated by the phrase. But at the same time some clearly compositional non-thesaurus phrases also have fairly low entropy of associations, for example, пресс-служба администрации (*press-service of the administration*).

We can also see that the phrases differ in the number of intersections between the associations obtained for a phrase and for its components. It seems natural that the already found conventionalized phrases have numerous intersections of this kind (Table 2) because the phrase and its components are closely related to each other.

On the contrary, other thesaurus phrases have a relatively small number of such intersections.

Phrase	A12	A21	Ph1	Ph2	Entr
финал лиги (league final)	0	0	2	18	2.16
начало года (beginning of the year)	0	0	1	0	2.66
ежедневный обзор (daily review)	1	0	3	14	2.90
пресс-служба администрации (press-service of administration)	0	0	14	6	2.97
еженедельный обзор (weekly review)	0	0	6	10	3.19
необходимый документ (necessary document)	1	0	0	14	3.64
конец января (end of January)	0	0	2	7	3.81
должность главы (post of the head)	0	0	5	2	3.85
новое поколение (new generation)	0	4	1	8	3.90
член совета (member of council)	5	0	2	5	3.96
повышение эффективности (increase in efficiency)	0	1	6	6	3.98
увеличение объема (growth in volume)	3	0	7	4	4.02
крупный размер (large size)	4	0	6	4	4.07
миллион евро (million of euros)	0	0	5	4	4.11
...					
особое внимание (special attention)	0	2	0	5	4.63
председатель комитета (chairman of committee)	5	0	5	3	4.63
экономический форум (economic forum)	0	2	2	3	4.65
интересный комментарий (interesting comment)	0	0	2	6	4.69
Average	1.22	0.36	2.80	6.36	4.07

Table 3: Results of association experiments for non-thesaurus phrases

It means that the thesaurus phrases evoke their own associations more often. For example, the phrase *повышение цен* (*increase of prices*) has frequent associations with the words *инфляция* (*inflation*) (16 of 25) and *кризис* (*crisis*), which were not mentioned as associations for its component words. On average, intersection between associations of the phrase and its component associations for non-thesaurus phrases is four times less than for thesaurus phrases.

It can also be seen that non-thesaurus phrases with low entropy of associations can have large numbers of intersections between the component associations and the phrase associations. In such cases, low entropy of the phrase associations is mainly determined by its components, for example, their probable syntactic dependencies. Only one of the non-thesaurus phrases has both low entropy of phrase associations and a few number of intersections of the phrase and component associations at the same time: *начало года* (begin-

ning of the year). It is highly associated with calendar months: *January* and *September*. For thesaurus phrases, a relatively high number of intersections between the phrase and component associations was revealed for most arguable thesaurus phrases: *транспортное происшествие* (*transport accident*) and *температура воздуха* (*air temperature*).

Thus, we can suppose that if a phrase has a low level of entropy of associations together with a small number of the same associations for the phrase and its components then it is also conventionalized.

We can introduce the threshold as $0.8 * \text{MaxEntropy}$ of answers. MaxEntropy is the maximal entropy we can obtain if respondents give equiprobable answers. In the current experiment, the threshold is equal to 3.76 for thesaurus phrases and 3.89 for non-thesaurus phrases. In our experiment, such conventionalized phrases include учебное заведение (*educational institute*), повышение цен (*increase in prices*), дорожное движение (*road traffic*), главный герой (*main hero*), медицинская помощь (*medical aid*).

As a result, we can say that we have found two signs of phrase conventionalization in the association experiment described:

- component words are frequently associated with each other, and
- associations of a phrase have both low entropy (less than $0.8 * \text{MaxEntropy}$) and a low level of intersection between component and phrase associations (less than 20%).

Using all three factors (association of component words to each other, entropy of phrase associations, and intersection of component word associations and phrase associations), it is possible to differentiate thesaurus phrases and non-thesaurus phrases with greater than 94% accuracy.

It is interesting to compare current results with the smaller amounts of associations. With this aim, we took the first 15 associations obtained for single words and phrases. The same above-mentioned thesaurus phrases have frequent mutual associations between components (that is, have association idiosyncrasy).

Phrases *медицинская помощь* (*medical aid*) and *температура воздуха* (*air temperature*) had entropy of associations more than $0.8 * \text{MaxEntropy}$. Only two non-thesaurus phrases had both

low entropy (less than $0.8 * \text{MaxEntropy}$) and the low level of intersection between associations of the phrase and its components: финал лиги (*league final*) and начало года (*beginning of the year*). As a result, in this smaller experiment, the obtained associations can distinguish thesaurus phrases with accuracy more than 92%.

5 Detecting the Conventionalized Expressions with Distributional Models

We compared the results of the association experiment with the results of distributional models. In previous works, it was supposed that non-compositional phrases can be distinguished with comparison of the phrase distributional vector and distributional vectors of their components: it was supposed that the similarity is less for non-compositional phrases (Cordeiro et al., 2016a; Gharbieh et al., 2016).

We used a Russian news collection (0.45 B tokens) and generated phrase and word embeddings with word2vec tool. For the phrases under consideration, we calculated cosine similarity between the phrase vector $v(w_1 w_2)$ and the sum of normalized vectors of phrase components $v(w_1 + w_2)$ according to formula from (Cordeiro et al., 2016a).

$$v(w_1 + w_2) = \left(\frac{v(w_1)}{|v(w_1)|} + \frac{v(w_2)}{|v(w_2)|} \right)$$

To evaluate different parameter sets, we located all phrases in the ascending order of similarity scores. We wanted to check if the thesaurus phrases with idiosynrasy obtain lesser values of word2vec similarity than non-thesaurus phrases without any specificity. We utilized MAP (mean average precision measure) to evaluate the quality of ordering.

We experimented with different parameters of word2vec and evaluated them with MAP on our data. We found that the best word2vec model (200 dimensions, 3 word window size) achieved quite low value of MAP (**0.391**), which means that it is very difficult for current embedding models to differentiate thesaurus and non-thesaurus phrases in our experiment.

We can also calculate MAP for the same phrase list ordered accoring to the increased entropy of phrase associations. And here we obtain MAP equal to **0.642**. Thus, entropy of human associations without accounting additional factors pre-

dicts thesaurus phrases significantly better than the embedding models.

6 Related Work

The annotation of multiword expressions on compositionality/non-compositionality of noun compounds has been studied in several works (Cordeiro et al., 2016b; Reddy et al., 2011; Ramisch et al., 2016).

Reddy et al. (2011) created the set of 90 noun compounds. The phrases were taken from WordNet. For each compound, the following types of tasks have been given: a judgement on how literal the phrase is and a judgement on how literal each noun is within the compound. They used 30 turkers to obtain judgements on the compound compositionality in each task.

Ramisch et al. (2016) asked respondents about the degree to which the meaning of an expression follows from its components: separately from each component and from both components in total. The authors of the paper stress that such indirect annotation provides reliable and stable data. However, this approach was confronted with difficulties concerning the inconsistency of the answers in some cases. For example, English speakers agreed on the level of head and head + modifier compositionality for phrase *dirty word*, but disagreed when judging the modifier: it was fully idiomatic for some, but others thought that the phrase just contained an uncommon sense of *dirty*.

Maziarz et al. (2015) try to formulate the procedural definition of multiword lexical units that should be included in the Polish wordnet so that lexicographers could apply these principles consistently. Then they asked linguists to classify phrases using this definition into three categories: *multiword lexical unit*, *not multiword lexical unit*, and *don't know*. They concluded that a group of 5-7 linguists is able to decide whether multiword lexical units should be introduced in a wordnet with the appropriate agreement. However, this approach was considered too expensive.

In another experiment, Maziarz et al. (2015) directed linguists to answer questions based on non-compositionality criteria of phrases including metaphoric character, hyponymy toward the syntactic head, ability to be paraphrased, non-separability, fixed word order, terminological register, etc. Then the answers were used to train the decision tree algorithm to predict inclusion or non-

inclusion of an expression into the Polish wordnet. However, the obtained decision trees were different for the various phrase sets under analysis.

Farahmand et al. (2015) describe the annotation of non-compositionality and conventionalization of noun compounds. They asked the annotators to make binary decisions about compositionality of phrases. Compositional compounds were further annotated as conventionalized or non-conventionalized. A compound was considered as conventionalized in neither of its constituents can be substituted for their near-synonyms. Sometimes the decision was difficult because such phrases could really exist (*floor space* vs. *floor area*).

To annotate the compounds, five experts were hired. In such a way, the authors (Farahmand et al., 2015) tried to avoid problems with crowdsourcing, which can lead to flaws in the results (Reddy et al., 2011). The authors stress that identifying conventionalization is not a trivial task and that human agreement on this property can be quite low. The examples of found compositional, but conventionalized phrases included: *cable car*, *food court*, *speed limit*, etc. The task of this study to distinguish conventionalized or non-conventionalized phrases among compositional compounds is the closest to our work.

For Russian there are two large resources of human associations. The well-known Russian Association dictionary (Karaulov et al., 1994) is currently obsolete. Another association-oriented project Sociation.org² collected a lot of current Russian associations but it does not have associations for the phrases under analysis.

Practical conclusions from the above-described experiments and related work are as follows:

- In annotating compositionality/non-compositionality of multiword expressions by crowdsourcing as in (Cordeiro et al., 2016b; Reddy et al., 2011; Ramisch et al., 2016), it is also useful to ask respondents about their associations for the phrase and its components to detect relational idiosyncrasy,
- In expert analysis of multiword expressions for inclusion into computational resources as in (Maziarz et al., 2015; Farahmand et al., 2015), it is useful to ask experts about additional lexical or conceptual relations that the

phrase have and that do not follow from the phrase components,

- In computational approaches of extracting non-compositional multiword expressions, it is useful to compare contexts of phrase occurrence and contexts of its component word occurrences trying to detect *weirdness* in the phrase context.

7 Conclusion

In this paper, we have shown that if we want to obtain human evidence about conventionalization of some phrases, we can ask native speakers about associations they have for a phrase and its component words.

We have found that there are two forms of manifesting conventionalized phrases. First, we can consider that a phrase is conventionalized if its component words have frequent associations to each other. The second type of conventionalized phrases can be revealed on the basis of two factors: low entropy of phrase associations and a low number of intersections between component word and phrase associations. These three factors allow predicting conventionalized phrases with high accuracy. We have also shown that the existing embedding models distinguish conventionalized phrases from non-conventionalized significantly worse.

In our opinion, developers of thesauri should consider the relational specificity (idiosyncrasy) of multiword expressions, which can help them to decide on inclusion of specific phrases into their resources. Weird word co-occurrences with the phrase in comparison with its component contexts can be considered as an additional factor to detect conventionalized expressions in computational approaches.

Acknowledgments.

This study is supported by Russian Scientific Foundation (project N16-18-02074).

References

- Eneko Agirre, Izaskun Aldezabal, and Eli Pociello. 2006. Lexicalization and multiword expressions in the basque wordnet. In *Proceedings of Third International WordNet Conference*. pages 131–138.
- ANSI/NISO. 2005. *Z39.19. Guidelines for the Construction, Format and Management of Monolingual Thesauri*. ANSI/NISO.

²<http://sociation.org/>

- Timothy Baldwin and Su Nam Kim. 2010. Multiword expressions. In *Handbook of Natural Language Processing, Second Edition*, Chapman and Hall/CRC, pages 267–292.
- Luisa Bentivogli and Emanuele Pianta. 2004. Extending wordnet with syntagmatic information. In *Proceedings of second global WordNet conference*, pages 47–53.
- Claire Bonial, Meredith Green, Jenette Preciado, and Martha Palmer. 2014. An approach to take multiword expressions. In *Proc. of the 10th Workshop on Multiword Expressions*, pages 94–98.
- Nicoletta Calzolari, Charles J Fillmore, Ralph Grishman, Nancy Ide, Alessandro Lenci, Catherine MacLeod, and Antonio Zampolli. 2002. Towards best practice for multiword expressions in computational lexicons. In *Proceedings of LREC-2002*.
- Silvio Cordeiro, Carlos Ramisch, Marco Idiart, and Aline Villavicencio. 2016a. Predicting the compositionality of nominal compounds: Giving word embeddings a hard time. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, pages 1986–1997.
- Silvio Cordeiro, Carlos Ramisch, and Aline Villavicencio. 2016b. Filtering and measuring the intrinsic quality of human compositionality judgments. In *ACL 2016*, pages 32–37.
- Meghdad Farahmand and James Henderson. 2016. Modeling the non-substitutability of multiword expressions with distributional semantics and a log-linear model. In *Proceedings of the 12th Workshop on Multiword Expressions, ACL 2016*, pages 61–66.
- Meghdad Farahmand, Aaron Smith, and Joakim Nivre. 2015. A multiword expression data set: Annotating non-compositionality and conventionalization for english noun compounds. In *Proceedings of NAACL-HLT*. Association for Computational Linguistics, pages 29–33.
- Alexander Gelbukh and Olga Kolesnikova. 2014. Multiword expressions in nlp: General survey and a special case of verb-noun constructions. *Computational Linguistics: Concepts, Methodologies, Tools, and Applications*, pages 178–197.
- Waseem Gharbieh, Virendra C Bhavsar, and Paul Cook. 2016. A word embedding approach to identifying verb–noun idiomatic combinations pages 112–118.
- Yuri Karaulov, Yu. Sorokin, E. Tarasov, N. Ufimtseva, and G. Cherkasova. 1994. *Russian Association Dictionary*.
- Natalia Loukachevitch and Boris Dobrov. 2014. Ruthes linguistic ontology vs. russian wordnets. In *Proceedings of Global WordNet Conference GWC-2014*, pages 154–162.
- Natalia Loukachevitch and German Lashevich. 2016. Multiword expressions in russian thesauri ruthes and ruwordnet. In *Proceedings of the AINL FRUCT 2016*. FRUCT, pages 66–71.
- Marek Maziarz, Stan Szpakowicz, and Maciej Piasecki. 2015. A procedural definition of multi-word lexical units. In *Proceedings of Recent Advances in NLP Conference RANLP-2015*, pages 427–435.
- Igor Mel’čuk. 2012. Phraseology in the language, in the dictionary, and in the computer. *Yearbook of Phraseology* 3(1):31–56.
- George A Miller. 1998. Nouns in wordnet. *WordNet: An electronic lexical database* pages 24–45.
- Darren Pearce. 2001. Synonymy in collocation extraction. In *Proceedings of the workshop on WordNet and other lexical resources, second meeting of the north american chapter of the association for computational linguistics*, pages 41–46.
- Pavel Pecina. 2010. Lexical association measures and collocation extraction. *Language resources and evaluation* 44(1-2):137–158.
- Maciej Piasecki, Michal Wendelberger, and Marek Maziarz. 2015. Extraction of the multi-word lexical units in the perspective of the wordnet expansion. In *RANLP-2015*, pages 512–520.
- Carlos Ramisch, Silvio Cordeiro, Leonardo Zilio, Marco Idiart, Aline Villavicencio, and Rodrigo Wilkens. 2016. How naked is the naked truth? a multilingual lexicon of nominal compound compositionality. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, pages 114–133.
- Siva Reddy, Diana McCarthy, and Suresh Manandhar. 2011. An empirical study on compositionality in compound nouns. In *IJCNLP*, pages 210–218.
- Ivan A Sag, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002. Multiword expressions: A pain in the neck for nlp. In *Proceedings of International Conference on Intelligent Text Processing and Computational Linguistics, CICLING-2002*. Springer Berlin Heidelberg, pages 1–15.
- Marco SG Senaldi, Gianluca E Leboni, and Alessandro Lenci. 2016. Lexical variability and compositionality: Investigating idiomaticity with distributional semantic models pages 21–31.
- Veronika Vincze and Attila Almasi. 2014. Non-lexicalized concepts in wordnets: A case study of english and hungarian. In *Proceedings of Global WordNet, Conference GWC-2014*. Global WordNet Association.