

A Multilingual Evaluation Dataset for Monolingual Word Sense Alignment

Sina Ahmadi*, John P. McCrae*,

Sanni Nimb¹, Fahad Khan³, Monica Monachini³, Bolette S. Pedersen⁸, Thierry Declercq^{2,12}, Tanja Wissik²,
 Andrea Bellandi³, Irene Pisani⁴, Thomas Troelsgård¹, Sussi Olsen⁸, Simon Krek⁵, Veronika Lipp⁶,
 Tamás Váradi⁶, László Simon⁶, András Gyórfy⁶, Carole Tiberius⁹, Tanneke Schoonheim⁹, Yifat Ben Moshe¹⁰,
 Maya Rudich¹⁰, Raya Abu Ahmad¹⁰, Dorielle Lonke¹⁰, Kira Kovalenko¹¹, Margit Langemets¹³, Jelena Kallas¹³,
 Oksana Dereza⁷, Theodorus Fransen⁷, David Cillessen⁷, David Lindemann¹⁴, Mikel Alonso¹⁴, Ana Salgado¹⁵
 José Luis Sancho¹⁶, Rafael-J. Ureña-Ruiz¹⁶, Jordi Porta Zamorano¹⁶, Kiril Simov¹⁷, Petya Osenova¹⁷,
 Zara Kancheva¹⁷, Ivaylo Radev¹⁷, Ranka Stanković¹⁸, Andrej Perdih¹⁹, Dejan Gabrovšek¹⁹

*Insight Centre for Data Analytics, National University of Ireland, Galway

{sina.ahmadi,john.mccrae}@insight-centre.org

(other affiliations in Appendix A)

Abstract

Aligning senses across resources and languages is a challenging task with beneficial applications in the field of natural language processing and electronic lexicography. In this paper, we describe our efforts in manually aligning monolingual dictionaries. The alignment is carried out at sense-level for various resources in 15 languages. Moreover, senses are annotated with possible semantic relationships such as broadness, narrowness, relatedness, and equivalence. In comparison to previous datasets for this task, this dataset covers a wide range of languages and resources and focuses on the more challenging task of linking general-purpose language. We believe that our data will pave the way for further advances in alignment and evaluation of word senses by creating new solutions, particularly those notoriously requiring data such as neural networks. Our resources are publicly available at <https://github.com/elexis-eu/MWSA>.

Keywords: lexical semantic resources, sense alignment, lexicography, language resource

1. Introduction

Lexical semantic resources (LSRs) are knowledge repositories that provide the vocabulary of a language in a descriptive and structured way. One of the famous examples of LSRs are dictionaries. Dictionaries form an important foundation of numerous natural language processing (NLP) tasks, including word sense disambiguation, machine translation, question answering and automatic summarization. However, the task of combining dictionaries from different sources is difficult, especially for the case of mapping the senses of entries, which often differ significantly in granularity and coverage. Approaches so far have mostly only been evaluated on named entities and quite specific domain language. In order to support a shared task at the GLOB-ALEX workshop¹, we have developed a new baseline that covers 15 languages and will provide a new baseline for the task of monolingual word sense alignment.

Different dictionaries and related resources such as word-nets and encyclopedia have significant differences in structure and heterogeneity in content, which makes aligning information across resources and languages a challenging task. Word sense alignment (WSA) is a more specific task of linking dictionary content at sense level which has been proved to be beneficial in various NLP tasks, such as word-sense disambiguation (Navigli and Ponzetto, 2012), semantic role labeling (Palmer, 2009) and information extraction (Moro et al., 2013). Moreover, combining LSRs can enhance domain coverage in terms of the number of lexical items and types of lexical-semantic information (Shi and

Mihalcea, 2005; Ponzetto and Navigli, 2010; Gurevych et al., 2012).

Given the current progress of artificial intelligence and the usage of data to train neural networks, annotated data with specific features play a crucial role to tackle data-driven challenges, particularly in NLP. In recent years, a few efforts have been made to create *gold-standard* dataset, i.e., a dataset of instances used for learning and fitting parameters, for aligning senses across monolingual resources including collaboratively-curated ones such as Wikipedia², and expert-made ones such as WordNet. However, the previous work is limited to a handful of languages and much of it is not on the core vocabulary of the language, but instead on named entities and specialist terminology. Moreover, despite the huge endeavour of lexicographers to compile dictionaries, proper lexicographic data are rarely openly accessible to researchers. In addition many of the resources are quite small and the extent to which the mapping is reliable is unclear.

In this paper, we present a set of datasets for the task of WSA containing manually-annotated monolingual resources in 15 languages. The annotation is carried out at sense level where four semantic relationships, namely, relatedness, equivalence, broadness, and narrowness, are selected for each pair of senses in the two resources by native lexicographers. Given the lexicographic context of this study, we have tried to provide lexicographic data from expert-made dictionaries. We believe that our datasets will pave the way for further developments in exploring statistical and neural methods, as well as for evaluation purposes. The rest of the paper is organized as follows: we first describe the previous work in Section 2. After having de-

* Contact Authors

¹<https://globalex2020.globalex.link/>

²<https://www.wikipedia.org>

Headword (POS)	R1-IDs	R1 senses	Semantic relation	Sense match	R2 senses	R2-IDs
clog (verb)						
	clog.v.02	dance a clog dance			to become clogged; to become loaded or encumbered, as with extraneous matter.	
	clog.v.03	impede the motion of, as with a chain or a burden			to encumber or load, especially with something that impedes motion; to hamper.	
	clog.v.01	become or cause to become obstructed			to coalesce or adhere; to unite in a mass.	
	clog.v.06	fill to excess so that function is impaired			to obstruct so as to hinder motion in or through; to choke up; .	
	clog.v.04	impede with a clog or as if with a clog			to burden; to trammel; to embarrass; to perplex.	
	clog.v.05	coalesce or unite in a mass				

Figure 1: Sense provided for *clog* (verb) in the English WordNet (R1) and the Webster Dictionary (R2). Drop-down lists are created dynamically for semantic relationship annotation.

scribed our methodology in Section 3, we further elaborate on the challenges of sense annotation in Section 4. We evaluate the datasets in Section 5 and finally, conclude the paper in Section 6.

2. Related work

Aligning senses across lexical resources has been attempted in several lexicographical milieus over the recent years. Such resources mainly include open-source dictionaries, WordNet and collaboratively-curated resources, such as Wikipedia. The latter has been shown to be reliable resources to construct accurate sense classifiers (Dandala et al., 2013).

There has been a significant body of research in aligning English resources, particularly, Princeton WordNet with Wikipedia (including (Ruiz-Casado et al., 2005; Ponzetto and Navigli, 2010; Niemann and Gurevych, 2011; McCrae, 2018)), with the Longman Dictionary of Contemporary English and with Roget’s thesaurus (Kwong, 1998), with Wiktionary³ (Meyer and Gurevych, 2011) or with the Oxford Dictionary of English (Navigli, 2006). Meyer and Gurevych (2011) also present a manually-annotated dataset for WSA between the English WordNet and Wiktionary.

On the other hand, there are a fewer number of manually aligned monolingual resources in other languages. For instance, there have been considerable efforts in aligning lexical semantic resources (LSRs) in German, particularly, the GermaNet–the German Wordnet (Hamp and Feldweg, 1997) with the German Wiktionary (Henrich et al., 2011), with the German Wikipedia (Henrich et al., 2012) and with the Digital Dictionary of the German Language (*Digitales Wörterbuch der Deutschen Sprache* (Klein and Geyken, 2010)) (Henrich et al., 2014). Gurevych et al. (2012) present UKB—a large-scale lexical-semantic resource containing pairwise sense alignments between a subset of nine resources in English and German which are mapped to a uniform representation. For Danish, aligning senses across modern lexical resources has been carried out in several projects in recent years (Pedersen et al., 2018), and a next natural step is to link these to historical Danish dictionaries.

Pedersen et al. (2009) describe the semi-automatic compilation of a WordNet for Danish, *DanNet*, based on a monolingual dictionary, the Danish Dictionary (*Den Danske Ordbog* (DDO)). Later, the semantic links between these two resources facilitated the compilation of a comprehensive thesaurus (*Den Danske Begrebsordbog*) (Nimb et al., 2014). The semantic links between thesaurus and dictionary made it possible to combine verb groups and dictionary valency information, used as input for the compilation of the Danish FrameNet Lexicon (Nimb, 2018). Furthermore, they constitute the basis for the automatically integrated information on related words in DDO, on the fly for each dictionary sense (Nimb et al., 2018). Similarly, Simov et al. (2019) report the manual mapping of the Bulgarian Word-Net BTB-WN with the Bulgarian Wikipedia.

Given the amount of the effort required to construct and maintain expert-made resources, various solutions have been proposed to automatically link and merge existing LSRs at different levels. LSRs being very diverse in domain coverage (Meyer, 2010; Burgun and Bodenreider, 2001), previous works have focused on methods to increase domain coverage, enrich sense representations and decrease sense granularity (Miller, 2016). Miller and Gurevych (2014) describe a technique for constructing an *n*-way alignment of LSRs and applied it to the production of a three-way alignment of the English WordNet, Wikipedia and Wiktionary. Niemann and Gurevych (2011) propose a threshold-based Personalized PageRank method for extracting a set of Wikipedia articles as alignment candidates and automatically aligning them with WordNet synsets. This method yields a sense inventory of higher coverage in comparison to taxonomy mapping techniques where Wikipedia categories are aligned to WordNet synsets (Ponzetto and Navigli, 2009). Matuschek and Gurevych present the Dijkstra-WSA algorithm as a graph-based approach (Matuschek and Gurevych, 2013) and a machine learning approach where features such as sense distances and gloss similarities are used for the task of WSA (Matuschek and Gurevych, 2014). It should be noted that all of these approaches produce results that are of lower reliability than gold standard datasets such as the ones presented in this paper.

³<https://www.wiktionary.org/>

3. Methodology

The main goal of the current study is to provide semantic relationships between two sets of senses for the same lemmas in two monolingual dictionaries. As an example, Figure 1 illustrates the senses for the entry “clog” (verb) in the English WordNet (Miller, 1995) (left) and the Webster’s Dictionary 1913 (Webster and Slater, 1828) (right). For further clarification, we provide two case studies of Danish and Italian in Section 4

The actual annotation was implemented by means of dynamic spreadsheets that provide a simple but effective manner to complete the annotation. This also had the added advantage that the annotation task could be easily completed from any device. In order to collect the data that was required for the annotation, each of the participating institutes provided their data in some form. We asked them, where possible, to organize their two dictionaries either in OntoLex-Lemon (Cimiano et al., 2016), TEI-Lex0 (Romary and Tasovac, 2018) or by following a simple TSV (tab-separated values) or Excel format providing the following data:

- An entry identifier, that locates the entry in the resource
- A sense identifier marking the sense in the resource, for example the sense number
- The lemma of the entry
- The part-of-speech of the entry
- The sense text, including the definition

In order to facilitate the task of annotation, we convert the initial data into spreadsheets. These spreadsheets provided an easy mapping and had the following columns:

- The headword and part of speech (given in parentheses after the headword);
- The sense text (definition) in the first resource;
- An interactive drop-down to specify one of the 5 semantic relations (see below) from the sense in the first resource;
- The sense text (abbreviated) in a drop-down list from the second resource, which the first resource is matched to;
- The full sense text of the second resource.

The fifth column played no technical role in the annotation, but was provided for reference, however as it was formatted with text wrapping on, it allowed the annotators to see the full definition of the second resource. In general we arranged the spreadsheets such that there were more senses for the first resource. In cases where the number of senses between the two resources were roughly equal, we created two spreadsheets based on which of the two datasets had more senses for those entries. In other cases, such as the English WordNet-Webster mapping where one resource (in this case WordNet) has many more senses, we used this as the first resource. Even still, there were some cases where the resource with more senses may contain a sense that corresponds to multiple senses in the second resource and in this case the annotators were instructed to simply use the “Insert Row Below” feature of the spreadsheet, which also duplicated the drop-down lists.

3.1. Semantic relationships

One of the challenges is that sense granularity between two dictionaries is rarely such that we would expect one-to-one mapping between the senses of an entry. In this respect, we followed a simple approach such as that in SKOS (Miles and Bechhofer, 2009) providing different kinds of linking predicates, which are described in Table 1. While it is certainly not easy to decide which relationship is to be used (we discuss this below), we found that this methodology was broadly effective and we believe will simplify the development of machine-learning-based classifiers for sense alignment prediction.

3.2. Data selection

The selection of the initial set of lemmas and senses to be aligned is guided by the following criteria:

- The lemmas should represent all open class words, namely nouns, verbs, adjectives and adverbs.
- Another criterion was that the lemmas should represent different degrees of polysemy, i.e. both highly polysemous lemmas as well as monosemous ones should be included.
- The lemmas in the two resources have the same part-of-speech tags. Spelling variations are normalized to a unique variation.

3.3. Dictionaries used in the creation of the dataset

For alignment we used the following dictionaries:

Basque The Basque Wordnet (MCR 3.0) and the Basque Monolingual Dictionary “*Euskal Hiztegia*” (copyright by the author, Ibon Sarasola) were linked.

Bulgarian The BulTreeBank Wordnet (BTB-WN) (Osenova and Simov, 2017) and the Bulgarian Wiktionary⁴ were used.

Danish We used the *Ordbog over det danske Sprog* (ODS)⁵ (Dahlerup, 1918), a historical dictionary covering 188,000 lemmas in Danish from 1700-1950, and *Den Danske Ordbog* (DDO) (Farø et al., 2003) a dictionary of modern Danish covering Danish from 1950 till today. One additional criterion in data selection was that at least one of the senses in DDO should be linked to a base or core concept in the Princeton WordNet via the Danish WordNet (Pedersen et al., 2019). This resulted in 4,500 DDO lemmas (of 97,500 in the dictionary). The lemma intersection (86%) with ODS was selected for our task.

Dutch We used the *Woordenboek der Nederlandsche Taal* (Dictionary of the Dutch Language, WNT)⁶ and the *Algemeen Nederlands Woordenboek* (Dictionary of Contemporary Dutch, ANW)⁷. The Dutch lemmas

⁴<https://bg.wiktionary.org>

⁵https://ordnet.dk/ods_en

⁶<http://gtb.ivdnt.org/search>

⁷<http://anw.ivdnt.org/search>

exact	The sense are the same, for example the definitions are simply paraphrases
broader	The sense in the first dictionary completely covers the meaning of the sense in the second dictionary and is applicable to further meanings
narrower	The sense in the first dictionary is entirely covered by the sense of the second dictionary, which is applicable to further meanings
related	There are cases when the senses may be equal but the definitions in both dictionaries differ in key aspects
none	There is no match for this sense

Table 1: Semantic relationships according to SKOS used for WSA task

were selected based on the Danish lemma list due to the close relationship between the two languages, facilitated by the information on the English equivalents from the Princeton WordNet.

English (KD) We used the Password and Global dictionary series provided by K Dictionaries through Lexicala⁸.

English (NUIG) As such, we developed a second English dataset using Princeton WordNet (Miller, 1995) (Fellbaum, 2010) and the public domain version of Webster’s dictionary from 1913⁹.

Estonian We used the EKS Dictionary of Estonian and the PSV Basic Estonian Dictionary (Kallas et al., 2014).

German We used the German versions of OmegaWiki¹⁰ and Wiktionary¹¹.

Hungarian We linked the Explanatory Dictionary of Hungarian (1959-1962)¹² containing 60,000 entries and, the Comprehensive Dictionary of Hungarian (2006-)¹³ containing 110,000 entries. Both are typical academic dictionaries.

Irish We used the Wiktionary data¹⁴ and *An Foclóir Beag* (Dónaill and Maoileoin, 1991, ‘The Little Dictionary’), the only two monolingual dictionaries available for this language.

Italian We used ItalWordNet (Roventini et al., 2000) and SIMPLE (Lenci et al., 2000).

Serbian We used the Serbian WordNet (Krstev et al., 2004; Stanković et al., 2018) and the *Rečnik Matice srpske I-VI: Rečnik srpskohrvatskog književnog jezika* (Dictionary of the Serbo-Croatian Literary Language).

Slovene (JSI) Slovene WordNet (Erjavec and Fiser, 2006) and Slovene Lexical Database (Gantar and Krek, 2011) were used.

Slovene (ISJFR) eSSKJ–Dictionary of the Slovenian Standard Language (3rd edition) (Gliha Komac et al., 2016) and the *Kostelski slovar* (Gregorič, 2014) were aligned.

Spanish The *Diccionario de la lengua española* (2011 edition) (RAE, 2001) was linked with the entries in the Spanish Wiktionary¹⁵ (backup dump of late August 2019) sharing the same lemmas.

Portuguese *Dicionário da Língua Portuguesa Contemporânea* (DLPC, (Casteleiro, 2001)) and *Dicionário Aberto* (DA)¹⁶ were used.

Russian Ozhegov and Shvedova’s ”The Dictionary of the Russian Language” (Ozhegov and Shvedova, 1992) and the Dictionary of the Russian Language edited by A.P. Evgenyeva, or *Maliy Akademicheskii Slovar* (Short Academic Dictionary) (Evgenyeva, 1999, MAS) were used.

3.4. Dataset structure

Listing 1 presents the structure of the datasets in JSON format. External keys such as `meta_ID` and `external_ID` will enable future lexicographers to integrate the annotations in external resources. Given that some of the semantic relationships, such as `narrower` and `broader`, are not symmetric, `sense_source` and `sense_target` are important classes in determining the semantic relationship correctly.

```
{
  "lemma": "splenetic",
  "POS_tag": "adjective",
  "gender": "",
  "meta_ID": "",
  "resource_1_senses": [
    {
      "#text": "of or relating to the spleen",
      "external_ID": "splenic.a.01"},
    {
      "#text": "very irritable",
      "external_ID": "bristly.s.01"}
  ],
  "resource_2_senses": [
    {
      "#text": "affected with spleen; malicious;
      ↪ spiteful; peevish; fretful.",
      "external_ID": ""}
  ],
  "alignment": [
    {
      "sense_source": "very irritable",
      "sense_target": "affected with spleen;
      ↪ malicious; spiteful; peevish;
      ↪ fretful.",
      "semantic_relationship": "exact"}
  ]
}
```

Listing 1: An example of the structure of senses and their alignments in the datasets

⁸<https://www.lexicala.com/>
⁹<https://www.websters1913.com/>
¹⁰<http://www.omegawiki.org/>
¹¹<https://de.wiktionary.org/>
¹²<http://mek.oszk.hu/adatbazis/magyar-nyelv-ertelmezo-szotara>
¹³<http://nagyszotar.nytud.hu>
¹⁴<https://ga.wiktionary.org>

¹⁵<https://es.wiktionary.org/>
¹⁶<https://dicionario-aberto.net>

4. Case Studies

We explain some of the challenges in the task based on the qualitative experience of two of the annotation teams. They report challenges in the preparation of data and in the annotation process.

4.1. *Ordbog over det Danske Sprog and Den Danske Ordbog*

The datasets for this task are created using the following steps:

- Extracting senses in ODS and DDO. This was a challenging process as different reference keys which are used for senses, were dealt with differently. For the same reason, we did not take multi-word expressions into account in the extraction process.
- Normalizing orthographies. As a historical dictionary, ODS employs an old Danish orthography. We automatically converted that orthography to the modern one using a mapping between characters.
- Dataset creation. Entries are linked using a common ID, called `metaID`, in ODS and DDO. Using this ID, senses of the same headwords in the two dictionaries are brought together for the annotation task.

When it came to the linking process between the senses of the two dictionaries, all senses and sub-senses within the sense hierarchy are brought together at the same level. This facilitated the annotation task as all possibilities could be visually taken into account easily. However, we believe that such a relaxation over the hierarchy may result in semantically less-representative senses.

Senses were considered to be ‘exact’ matches also in cases where definitions differed slightly due to new techniques and modernisation in society. E.g. the historical sense of the noun *passager* (‘passenger’) (‘person travelling with mail coach etc.’) was considered an exact match to the modern sense ‘person travelling with private or public means of transportation’.

The more vague ‘related’ relation was used when there were differences in ontological type between the two definitions, e.g. the property of ‘being able to sleep’, a sense of the noun *søvn* (‘sleep’) in the historical dictionary, is ‘related’ to ‘the state of sleeping’ sense in the modern dictionary. Often such differences in ontological type across the two dictionaries were due to regular polysemy (act/result, semiotic artifact/content, animal/food, organisation/building etc., see for example (Pustejovsky, 1995)). Two dictionaries will often differ in their descriptions in cases of regular polysemy, focusing on either one or the other sense leaving one of them under-specified, or describing both of them. For instance, while DDO for the noun *afsked* ‘farewell’ describes the act of saying farewell, ODS focuses on the result, namely the phrase ‘farewell’, therefore the senses are only ‘related’ and not exact matches. Likewise, ODS has only one sense for the noun *ambassade* ‘embassy’, namely the ‘organisation’ sense, while DDO has two: the organisation sense, but also the building sense. Moreover, ‘related’ has also been used when the ontological type is in fact the same for the two senses, but where

other parts of the definitions differ slightly, e.g. in the case of the noun *bamse* (‘bear, teddy bear’). The sense in the historical dictionary, i.e. ‘fat, clumsy person, especially a child’, is considered to be ‘related’ to the modern sense of the same lemma, i.e. ‘fat, good-natured person’.

Regarding the ‘broader’ and ‘narrower’ relations, the historical ODS sense was for example considered to be ‘broader’ in the case of the noun *værge* (guardian): ‘a guardian of anything or anybody’ which in the modern dictionary is restricted to only being ‘a guardian in legal context’ (i.e. a guardian for a child not yet legally competent or for an incapacitated adult). An opposite case where the historical sense is ‘narrower’ than the modern one can be illustrated by the adjective *spids* (‘sharp’) where ODS describes two specific senses, one about sound and another one about smell, while DDO merges the two senses into one: ‘pungent in an unpleasant way (about smell, taste or sound)’.

4.2. *ItalWordNet and SIMPLE*

Regarding Italian, the team at ILC-CNR chose ItalWordNet (IWN) and SIMPLE, two Italian language lexical resources which had been previously developed in the institute. The former, IWN, is a lexical semantic network for Italian (Roventini et al., 2002) which is part of the WordNet family (Miller, 1995). As such it is organised around the notion of a synset of word senses and the network structure based on lexical-semantic relations which hold between senses across synsets. The 50,000 Italian synsets contained in IWN are linked to the Princeton Wordnet. The latter resource, SIMPLE, constitutes the semantic level of a quadripartite Italian lexicon. Its structure is inspired by Generative Lexicon theory (Pustejovsky, 1995) and in particular the notion of qualia structure which is used to organise the Semantic Units (SemUs) which constitute the basic structures representing word-sense. SIMPLE contains 20,000 SemUs and we used the definitions of these SemUs for the task. Both lexicons share a set of common “base concepts” that provided the basis of a previous (semi-)automatic mapping of the two lexicons on the basis of their respective ontological organisations (Roventini et al., 2007; Roventini and Ruimy, 2008). Although this mapping did not make the five-fold distinction, i.e., exact, narrower, broader, related, and none, it did constitute a useful starting point and a basis for comparison for the task.

The teams that had originally compiled IWN and SIMPLE shared many members in common and so, the definitions for corresponding senses across the two lexicons are sometimes very similar or differ solely on the basis of an extra clause. This made it easy to determine, in many cases, if two senses were ‘exact’ matches or if one was ‘broader’ or ‘narrower’ than the other by just comparing strings. The applicability of the ‘related’ category was less clear than the others but the annotator made use of it in cases where two senses referred to different concepts which did not match but were semantically related, as well as in cases of metaphoric senses in which one sense refers to the concrete and the other to the metaphorical meaning.

The annotator found the most challenging aspect of the task to lie in the necessity of having to choose the type of match-

Language	Resource	Nouns	Verbs	Adjectives	Adverbs	Other	All
Basque	Basque Wordnet	929 (6836)	0 (0)	0 (0)	0 (0)	0 (0)	929 (6836)
	<i>Euskal Hiztegia</i>	971 (7754)	0 (0)	0 (0)	0 (0)	0 (0)	971 (7754)
Bulgarian	BTB-WN	1394 (15649)	175 (1698)	305 (3187)	50 (338)	0 (0)	1924 (20872)
	Bulgarian Wiktionary	1273 (12883)	164 (1107)	194 (1418)	39 (306)	0 (0)	1670 (15714)
Danish	<i>Ordbog over det danske Sprog</i>	2176 (282040)	983 (119163)	436 (60599)	0 (0)	0 (0)	3595 (461802)
	<i>Den Danske Ordbog</i>	1036 (12326)	383 (4045)	248 (2228)	0 (0)	0 (0)	1667 (18599)
Dutch	<i>Woordenboek der Nederlandsche Taal</i>	1459 (28979)	405 (5185)	527 (7878)	106 (2662)	0 (0)	2497 (44704)
	<i>Algemeen Nederlands Woordenboek</i>	497 (8443)	140 (1542)	109 (1393)	13 (172)	0 (0)	759 (11550)
English (KD)	Global	92 (532)	107 (617)	80 (457)	57 (257)	61 (283)	397 (2146)
	Password	66 (536)	72 (417)	62 (324)	33 (177)	46 (188)	279 (1642)
English (NUIG)	<i>Webster</i>	1131 (11606)	741 (4622)	373 (2585)	45 (269)	0 (0)	2290 (19082)
	<i>Princeton WordNet</i>	730 (12166)	496 (6980)	249 (2892)	24 (207)	0 (0)	1499 (22245)
Estonian	Dictionary of Estonian (EKS)	543 (4012)	273 (1598)	151 (747)	98 (451)	78 (370)	1143 (7178)
	Estonian Basic Dictionary (PSV)	543 (4492)	273 (1983)	151 (1097)	98 (596)	79 (468)	1144 (8636)
German	German Wiktionary	2026 (15160)	0 (0)	0 (0)	0 (0)	0 (0)	2026 (15160)
	German OmegaWiki	1266 (14354)	0 (0)	0 (0)	0 (0)	0 (0)	1266 (14354)
Hungarian	Comprehensive						1355 (14654)
	Explanatory						1038 (10934)
Irish	<i>An Foclóir Beag</i>	891 (8053)	11 (95)	55 (267)	10 (56)	36 (171)	1003 (8642)
	Irish Wiktionary	1209 (6696)	8 (45)	61 (181)	10 (41)	36 (109)	1324 (7072)
Italian	ItalWordNet	408 (3128)	352 (2411)	0 (0)	0 (0)	0 (0)	760 (5539)
	SIMPLE	290 (1990)	218 (1240)	0 (0)	0 (0)	0 (0)	508 (3230)
Serbian	Serbian WordNet	691 (5864)	985 (6522)	92 (713)	0 (0)	0 (0)	1768 (13099)
	Dictionary of Serbo-Croatian Literary Language	289 (2360)	281 (1527)	29 (215)	0 (0)	0 (0)	599 (4102)
Slovenian (JSI)	Slovene WordNet	409 (1106)	303 (901)	237 (733)	44 (133)	0 (0)	993 (2873)
	Slovene Lexical Database	284 (2237)	191 (1047)	220 (1486)	29 (102)	0 (0)	724 (4872)
Slovenian (ISJFR)	Standard Slovenian Dictionary (eSSKJ)	229 (2060)	109 (911)	76 (620)	0 (0)	60 (588)	474 (4179)
	<i>Kostelski slovar</i>	151 (1050)	61 (308)	45 (257)	0 (0)	38 (263)	295 (1878)
Spanish	<i>Diccionario de la lengua española</i>	617 (7986)	225 (2426)	305 (3269)	26 (161)	24 (250)	1197 (14092)
	Spanish Wiktionary	602 (6421)	227 (2045)	294 (2825)	25 (129)	22 (123)	1170 (11543)
Portuguese	<i>Dicionário da Língua Portuguesa Contemporânea</i>	285 (4060)	58 (686)	110 (1287)	9 (143)	1 (9)	463 (6185)
	<i>Dicionário Aberto</i>	199 (1521)	53 (203)	67 (372)	3 (15)	1 (5)	323 (2116)
Russian	<i>Ozhegov-Shvedova</i>	258 (2038)	109 (615)	101 (533)	15 (77)	44 (368)	527 (3631)
	Dictionary of the Russian Language (MAS)	310 (2811)	173 (1338)	190 (1219)	20 (114)	71 (1010)	764 (6492)

Table 2: Statistics of the datasets. This table shows the number of senses in the resources (number of the words in the definitions are provided in parentheses).

ing relationship from out of the five options available. This choice was not always an intuitive one and the procedure often called for a careful analysis in order to achieve as objective an assessment of the case under consideration as possible. The annotator also found it useful to consult other lexical resources, in particular two online versions of the well known *Treccani*¹⁷ and *Garzanti*¹⁸ reference dictionaries.

5. Evaluation

We performed an intrinsic evaluation on our datasets by computing a number of resource statistics on the senses. Table 2 provides resource statistics based on part-of-speech tags and languages. As most of the lemmas available in the resources belong to open classes, namely nouns, verbs, adjectives and adverbs, we carried out our experiments with respect to those part-of-speech tags. Moreover, there are few languages, such as German, Italian and Serbian, for which only a certain number of the part-of-speech tags

¹⁷<http://www.treccani.it>

¹⁸<https://www.garzantilinguistica.it>

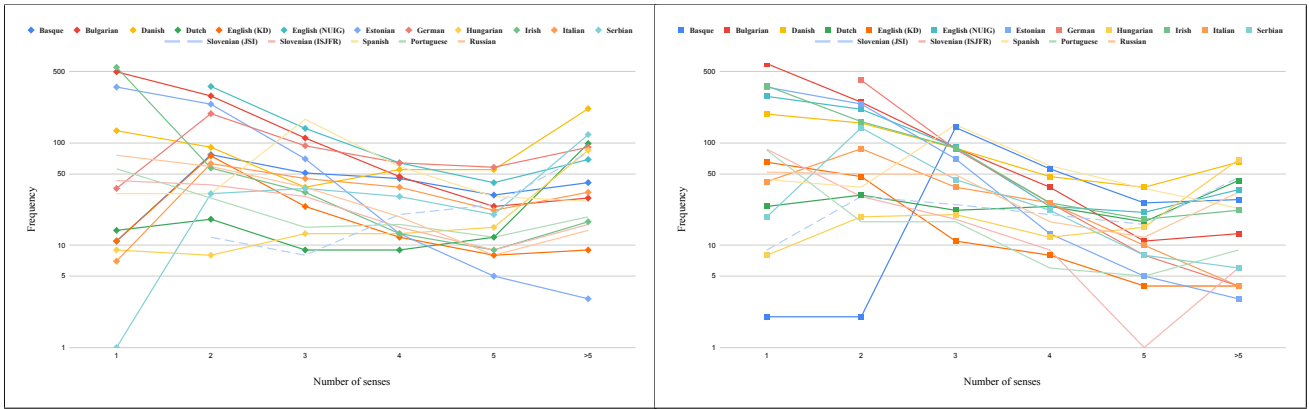


Figure 2: Frequency of the number of senses in the datasets per language and resource (left resources at left and right resources at right)

are available. As a unique case, the Hungarian entries are aligned at lemma-level without taking the POS tags into account. The POS tags are provided within the senses, upon the lexicographers' request.

Moreover, the distribution of the frequency of number of senses is presented in Figure 2, where we show for each resource how many entries had 1, 2, 3, 4, 5 or more senses.

5.1. Sense granularity

The granularity of senses is a determining factor in applying automatic approaches for semantic similarity evaluation. Sense granularity does not follow an identical pattern across resources and languages. The type of the resource, the preference of the lexicographer and the historical period of the resource edition are some of the factors on how senses are shaped.

Figure 3 illustrates the correlation between number of tokens in the first and second resource of the languages provided in our datasets. To calculate the correlation, we divide the number of space-separated tokens in one of the resources by the other resource. Although most of the resources have a correlation of [1, 2] which indicates a relatively similar granularity of senses in the two resources, Danish and English (NUIG) represent higher correlations. In the case of Danish, a correlation of 24.8 demonstrates a huge difference in how senses are expanded in the resources. This can be justified by the fact that ODS as a historical resource provides many senses which are no longer used in the language. In addition, the structure of the resource is in such a way that citations and further details are provided at sense-level rather than separately.

5.2. Sense alignments

One of the main challenges in aligning senses are due to the structure of the senses. A resource which provides senses in a hierarchy based on main senses and their sub-senses represents semantically context-dependent senses in comparison to one in which senses are *semantically independent*, which are stand-alone senses not influenced by the hierarchy. On the other hand, senses may contain descriptions beyond the definition, such as usage examples and idioms. To evaluate the distribution of the alignments with respect to the senses, we assume that each entry is a *lexicographic*

network (Ahmadi et al., 2018), i.e., a graph where the nodes and edges are the senses and the alignments, respectively. Given a set of aligned senses, we denote the number of senses in resource 1 and resource 2 by n_1 and n_2 , respectively. We also denote the number of alignments in each entry by m . Therefore, the average degree of senses in each resource is defined as $k_1 = \frac{m}{n_1}$ and $k_2 = \frac{m}{n_2}$. Similarly, the average degree of the whole dataset can be calculated as $k = \frac{2 \times m}{n_1 + n_2} = \frac{n_1 \times k_1 + n_2 \times k_2}{n_1 + n_2}$. Finally, we define the number of existing alignments divided by the number of possible alignments as the density $\delta = \frac{m}{n_1 \times n_2}$.

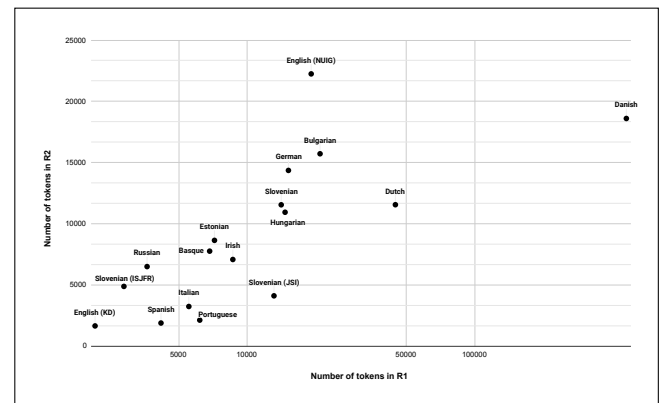


Figure 3: Correlation between number of tokens

Table 3 represents the results of our evaluations on the aligned senses. The degree indicates the distribution of the alignments with respect to the senses. For instance, a degree of 1.182 (k_1) in the case of Russian shows that every sense is at least aligned with another one. On the other hand, a low degree of 0.250 (k_1) in the case of Dutch indicates the sparsity of alignments over the senses. Moreover, density δ provides an insight into how alignments are distributed over the combination of all senses. In other words, a higher density represents a higher probability that two senses are aligned in the two resources. Estonian and German resources, for example, have the highest density among the resources.

Language	Semantic relationship					k_1	k_2	k	δ
	exact	narrower	broader	related	all				
Basque	399	138	94	184	815	0.877	0.839	0.858	9.03E-04
Bulgarian	958	274	254	492	1978	1.028	1.184	1.101	6.16E-04
Danish	1103	316	189	36	1644	0.457	0.986	0.625	1.04E-07
Dutch	489	30	64	42	625	0.250	0.823	0.384	3.30E-04
English (KD)	107	78	28	88	301	0.758	1.079	0.891	2.72E-03
English (NUIG)	885	339	42	67	1333	0.582	0.889	0.704	3.88E-04
Estonian	1025	61	54	4	1144	1.001	1.000	1.000	5.00E-01
German	354	311	426	126	1217	0.601	0.961	0.739	3.70E-01
Hungarian	465	214	227	43	949	0.700	0.914	0.793	6.75E-04
Irish	731	45	67	132	975	0.972	0.736	0.838	7.34E-04
Italian	327	132	44	89	592	0.779	1.165	0.934	1.53E-03
Serbian	325	47	73	146	591	0.334	0.987	0.499	5.58E-04
Slovenian (JSI)	306	183	169	54	712	0.717	0.983	0.829	9.90E-04
Slovenian (ISJFR)	110	88	10	39	247	0.521	0.837	0.642	1.77E-03
Spanish	867	185	114	93	1259	1.052	1.076	1.064	8.99E-04
Portuguese	207	38	2	28	275	0.594	0.851	0.700	1.84E-03
Russian	363	15	159	86	623	1.182	0.815	0.965	1.55E-03

Table 3: A description of the semantic relationship alignments using basic graph measures

5.3. Inter-annotator agreement

While the linking for most of the languages was only developed by a single annotator, we collected multiple annotations for four languages which enabled us to evaluate the alignment agreement over the same senses. Given the invariable number of annotators depending on the language and, the categorical nature of the problem, we used the Krippendorff’s alpha-reliability (Krippendorff, 2011) for calculating the inter-annotator agreement (IAA) where we considered each possible sense pair as an item for the agreement. Thus, if a pair of senses was not chosen by any of the annotators, they are considered to agree that the link between this is *none*. Table 4 presents the IAA in a 5-class model, that is the five semantic relationships. Moreover, we provide a 2-class model where all types of semantic relationships, namely exact, broader, narrower and related, are merged and compared against ‘none’ as the other class. Regarding the number of senses, 561, 4979, 185 and 270 senses were annotated by more than one annotator for English, German, Irish and Danish, respectively, which made it possible to calculate IAA.

Regarding the English (KD) resources, an internal evaluation of the annotated data with two annotators show an agreement for 76% of the annotators.

	Agreement (5-class)	Agreement (2-class)
Irish (3)	0.83	0.99
English (NUIG) (3)	0.43	0.73
Danish (2)	0.95	0.92
German (2)	0.71	0.58

Table 4: Inter-annotator agreement using Krippendorff’s alpha. Number of annotators provided in parentheses.

6. Conclusion

In this paper, we presented a set of 17 datasets for the task of monolingual word sense alignment covering 15 lan-

guages. This dataset innovates on previous datasets by focusing on general vocabulary, which is much harder to link than the focus of previous works. In addition to the collaboratively-curated resources such as Wiktionary, many expert-made resources are used in our datasets for the task. We developed the alignment using 5 categories of links, namely exact, broader, narrower, related and not related, i.e. *none*, and found that our annotators were able to perform this task with high agreement. Given the significant size of the datasets, we believe that they can be beneficial not only for evaluation purposes, but also for training new statistical and neural models for various tasks such as word sense alignment, semantic relationship detection, paraphrasing and semantic entailment, to mention but a few.

As future work, we are planning to evaluate the performance of various methods for the tasks of sense alignment and semantic relationship detection using these datasets. Moreover, we would like to explore language-independent techniques to facilitate monolingual lexical data linking and increase the interoperability of both monolingual and multilingual dictionaries.

7. Acknowledgements

The authors would like to thank the three anonymous reviewers for their insightful suggestions and careful reading of the manuscript. This work has received funding from the EU’s Horizon 2020 Research and Innovation programme through the ELEXIS project under grant agreement No. 731015. The contributions in Bulgarian were partially funded by the Bulgarian National Interdisciplinary Research e-Infrastructure for Resources and Technologies in favor of the Bulgarian Language and Cultural Heritage, part of the EU infrastructures CLARIN and DARIAH – CLaDA-BG, Grant number DO1-272/16.12.2019. This work is also supported by Sci-



ence Foundation Ireland (SFI) under the Insight Center for Data Analytics (Grant Number SFI/12/RC/2289) and the Irish Research Council under the “Cardamom” Consolidator Laureate Grant (IRCLA/2017/129).

8. Bibliographical References

- Ahmadi, S., Arcan, M., and McCrae, J. (2018). On lexicographical networks. In *Workshop on eLexicography: Between Digital Humanities and Artificial Intelligence*.
- Burgun, A. and Bodenreider, O. (2001). Comparing terms, concepts and semantic classes in WordNet and the Unified Medical Language System. In *in WordNet and the Unified Medical Language System. Proc NAACL Workshop, WordNet and Other Lexical Resources: Applications, Extensions and Customizations*, pages 77–82.
- Casteleiro, J. M. (2001). Dicionário da língua portuguesa contemporânea. *Lisboa: Academia das Ciências de Lisboa e Editorial Verbo*, 2.
- Cimiano, P., McCrae, J. P., and Buitelaar, P. (2016). Lexicon Model for Ontologies: Community Report.
- Dahlerup, V. (1918). *Ordbog over det danske sprog*, volume 1. Gyldendal.
- Dandala, B., Mihalcea, R., and Bunescu, R. (2013). Word sense disambiguation using Wikipedia. In *The People’s Web Meets NLP*, pages 241–262. Springer.
- Dónaill, N. O. and Maoileoin, P. U. (1991). *An Foclóir Beag*. An Gum.
- Erjavec, T. and Fiser, D. (2006). Building Slovene WordNet. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation*, pages 1678–1683.
- A. P. Evgenyeva, editor. (1999). *Dictionary of the Russian Language*, volume 1-4. Russkiy yazyk.
- Ken Farø, et al., editors. (2003). *Den danske ordbog, bd. 1-6*, volume 1-6. Gyldendal. https://ordnet.dk/ddo_en.
- Fellbaum, C. (2010). WordNet. In *Theory and applications of ontology: computer applications*, pages 231–243. Springer.
- Gantar, P. and Krek, S. (2011). Slovene lexical database. In *Natural language processing, multilinguality*, pages 72–80.
- Gliha Komac, N., Jakop, N., Ježovnik, J., Kern, B., Klemenčič, S., Krvina, D., Ledinek, N., Meterc, M., Michelizza, M., Pavlič, M., et al. (2016). *eSSKJ: Dictionary of the Slovenian Standard Language*. ZRC SAZU, 3rd edition.
- Gregorič, J. (2014). *Kostelski slovar*. Založba ZRC.
- Gurevych, I., Ecker-Kohler, J., Hartmann, S., Matuschek, M., Meyer, C. M., and Wirth, C. (2012). Uby: A large-scale unified lexical-semantic resource based on LMF. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 580–590. Association for Computational Linguistics.
- Hamp, B. and Feldweg, H. (1997). Germanet-a lexical-semantic net for German. In *Automatic information extraction and building of lexical semantic resources for NLP applications*.
- Henrich, V., Hinrichs, E., and Vodolazova, T. (2011). Semi-automatic extension of GermaNet with sense definitions from Wiktionary. In *Proceedings of the 5th Language and Technology Conference (LTC 2011)*, pages 126–130.
- Henrich, V., Hinrichs, E. W., and Suttner, K. (2012). Automatically linking GermaNet to Wikipedia for harvesting corpus examples for Germanet senses. *JLCL*, 27(1):1–19.
- Henrich, V., Hinrichs, E., and Barkey, R. (2014). Aligning word senses in GermaNet and the DWDS dictionary of the German language. In *Proceedings of the Seventh Global Wordnet Conference*, pages 63–70.
- Kallas, J., Tuulik, M., and Langemets, M. (2014). The basic Estonian dictionary: the first monolingual L2 learner’s dictionary of Estonian. In *Proceedings of the XVI Euralex Congress*.
- Klein, W. and Geyken, A. (2010). Das digitale Wörterbuch der deutschen Sprache (DWDS). In *Lexicographica: International annual for lexicography*, pages 79–96. De Gruyter.
- Krippendorff, K. (2011). Computing Krippendorff’s alpha-reliability. *Annenberg School for Communication Departmental Papers. Philadelphia*.
- Krstev, C., Palović-Lažetić, G., Vitas, D., and Obradović, I. (2004). Using textual and lexical resources in developing Serbian Wordnet. *SCIENCE AND TECHNOLOGY*, 7(1-2):147–161.
- Kwong, O. Y. (1998). Aligning WordNet with additional lexical resources. *Usage of WordNet in Natural Language Processing Systems*.
- Lenci, A., Bel, N., Busa, F., Calzolari, N., Gola, E., Monachini, M., Ogonowski, A., Peters, I., Peters, W., Ruimy, N., et al. (2000). SIMPLE: A general framework for the development of multilingual lexicons. *International Journal of Lexicography*, 13(4):249–263.
- Matuschek, M. and Gurevych, I. (2013). Dijkstra-WSA: A graph-based approach to word sense alignment. *Transactions of the Association for Computational Linguistics*, 1:151–164.
- Matuschek, M. and Gurevych, I. (2014). High performance word sense alignment by joint modeling of sense distance and gloss similarity. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 245–256.
- McCrae, J. P. (2018). Mapping wordnet instances to Wikipedia. In *Proceedings of the 9th Global WordNet Conference (GWC 2018)*, pages 62–69.
- Meyer, C. M. and Gurevych, I. (2011). What psycholinguists know about chemistry: Aligning Wiktionary and WordNet for increased domain coverage. In *Proceedings of 5th International Joint Conference on Natural Language Processing*, pages 883–892.
- Meyer, C. M. (2010). How web communities analyze human language: Word senses in Wiktionary. In *In Second Web Science Conference*.
- Miles, A. and Bechhofer, S. (2009). SKOS simple knowledge organization system reference. *W3C recommendation*, 18:W3C.

- Miller, T. and Gurevych, I. (2014). WordNet—Wikipedia—Wiktionary: Construction of a three-way alignment. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation*, pages 2094–2100.
- Miller, G. A. (1995). WordNet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- Miller, T. (2016). *Adjusting Sense Representations for Word Sense Disambiguation and Automatic Pun Interpretation*. Ph.D. thesis, Technische Universität Darmstadt, Darmstadt, January.
- Moro, A., Li, H., Krause, S., Xu, F., Navigli, R., and Uszkoreit, H. (2013). Semantic rule filtering for web-scale relation extraction. In *International Semantic Web Conference*, pages 347–362. Springer.
- Navigli, R. and Ponzetto, S. P. (2012). BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193:217–250.
- Navigli, R. (2006). Meaningful clustering of senses helps boost word sense disambiguation performance. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 105–112. Association for Computational Linguistics.
- Niemann, E. and Gurevych, I. (2011). The people’s web meets linguistic knowledge: Automatic sense alignment of Wikipedia and WordNet. In *Proceedings of the Ninth International Conference on Computational Semantics (IWCS 2011)*.
- Nimb, S., Trap-Jensen, L., and Lorentzen, H. (2014). The Danish thesaurus: Problems and perspectives. In *Proceedings of the XVI EURALEX International Congress: The User in Focus*, pages 15–19.
- Nimb, S., Sørensen, N. H., and Troelsgård, T. (2018). From standalone thesaurus to integrated related words in the Danish dictionary. In *Proceedings of the XVIII EURALEX International Congress: Lexicography in Global Contexts*, pages 915–923.
- Nimb, S. (2018). The Danish FrameNet lexicon: method and lexical coverage. In *Proceedings of the International FrameNet Workshop at LREC 2018: Multilingual FrameNets and Constructions*, pages 51–55.
- Osenova, P. and Simov, K. (2017). Challenges behind the data-driven bulgarian wordnet (bultreebank bulgarian wordnet). In John P. McCrae, et al., editors, *Proceedings of the LDK 2017 Workshops: 1st Workshop on the OntoLex Model (OntoLex-2017), Shared Task on Translation Inference Across Dictionaries & Challenges for Wordnets co-located with 1st Conference on Language, Data and Knowledge (LDK 2017), Galway, Ireland, June 18, 2017*, volume 1899 of *CEUR Workshop Proceedings*, pages 152–163. CEUR-WS.org.
- Ozhegov, S. I. and Shvedova, N. Y. (1992). *Explanatory Dictionary of the Russian Language*. Az, Moscow.
- Palmer, M. (2009). SemLink: Linking PropBank, VerbNet and FrameNet. In *Proceedings of the generative lexicon conference*, pages 9–15. GenLex-09, Pisa, Italy.
- Pedersen, B. S., Nimb, S., Asmussen, J., Sørensen, N. H., Trap-Jensen, L., and Lorentzen, H. (2009). DanNet: the challenge of compiling a wordnet for Danish by reusing a monolingual dictionary. *Language resources and evaluation*, 43(3):269–299.
- Pedersen, B. S., Nimb, S., Olsen, S., and Sørensen, N. H. (2018). Combining dictionaries, wordnets and other lexical resources—advantages and challenges. In *Globalex Proceedings 2018, Miyasaki, Japan*.
- Pedersen, B. S., Nimb, S., Olsen, I. R., , and Olsen, S. (2019). Linking DanNet with Princeton WordNet. In *Global WordNet 2019 Proceedings*.
- Ponzetto, S. P. and Navigli, R. (2009). Large-scale taxonomy mapping for restructuring and integrating Wikipedia. In *Twenty-First International Joint Conference on Artificial Intelligence*.
- Ponzetto, S. P. and Navigli, R. (2010). Knowledge-rich word sense disambiguation rivaling supervised systems. In *Proceedings of the 48th annual meeting of the association for computational linguistics*, pages 1522–1531. Association for Computational Linguistics.
- Pustejovsky, J. (1995). *The Generative Lexicon*. MIT Press, Cambridge, MA.
- RAE, R. A. E. (2001). *Diccionario de la Lengua Espanola RAE*. Diccionario de la lengua española. Planeta Publishing Corporation.
- Romary, L. and Tasovac, T. (2018). TEI Lex-0: A target format for TEI-encoded dictionaries and lexical resources. In *TEI Conference and Members’ Meeting*.
- Roventini, A. and Ruimy, N. (2008). Mapping events and abstract entities from PAROLE-SIMPLE-CLIPS to ItalWordNet. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC’08)*, Marrakech, Morocco, May. European Language Resources Association (ELRA).
- Roventini, A., Alonge, A., Calzolari, N., Magnini, B., and Bertagna, F. (2000). ItalWordNet: a large semantic database for italian. In *Proceedings of the Second International Conference on Language Resources and Evaluation*.
- Roventini, A., Ulivieri, M., and Calzolari, N. (2002). Integrating two semantic lexicons, SIMPLE and ItalWordNet: What can we gain? In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC’02)*, Las Palmas, Canary Islands - Spain, May. European Language Resources Association (ELRA).
- Roventini, A., Ruimy, N., Marinelli, R., Ulivieri, M., and Mammini, M. (2007). Mapping concrete entities from PAROLE-SIMPLE-CLIPS to ItalWordNet: Methodology and results. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 161–164, Prague, Czech Republic, June. Association for Computational Linguistics.
- Ruiz-Casado, M., Alfonseca, E., and Castells, P. (2005). Automatic assignment of Wikipedia encyclopedic entries to wordnet synsets. In *International Atlantic Web Intelligence Conference*, pages 380–386. Springer.
- Shi, L. and Mihalcea, R. (2005). Putting pieces together:

- Combining FrameNet, VerbNet and WordNet for robust semantic parsing. In *International conference on intelligent text processing and computational linguistics*, pages 100–111. Springer.
- Simov, K., Osenova, P., Laskova, L., Radev, I., and Kancheva, Z. (2019). Aligning the bulgarian btb wordnet with the bulgarian wikipedia. In Christiane Fellbaum, et al., editors, *Proceedings of the Tenth Global Wordnet Conference*, pages 290–297.
- Stanković, R., Mladenović, M., Obradović, I., Vitas, M., and Krstev, C. (2018). Resource-based WordNet augmentation and enrichment. In Svetla Koeva, editor, *Proceedings of the Third International Conference Computational Linguistics in Bulgaria (CLIB 2018)*, pages 104–114, Sofia, Bulgaria, May. Institute for Bulgarian Language “Prof. Lyubomir Andreychin”, Bulgarian Academy of Sciences.
- Webster, N. and Slater, R. J. (1828). *Noah Webster’s first edition of an American dictionary of the English language*. Foundation for American Christian Education San Francisco.
- ¹⁵Academia das Ciências de Lisboa, Lisbon, Portugal
anacastrosalgado@gmail.com
- ¹⁶Centro de estudios de la Real Academia Española, Madrid, Spain
{sancho,rafa,porta}@rae.es
- ¹⁷Bulgarian Academy of Sciences, Bulgaria
{kivs,petya,zara,radev}@bultreebank.org
- ¹⁸University of Belgrade, Belgrade, Serbia ranka@rgf.rs
- ¹⁹Research Centre of the Slovenian Academy of Sciences and Arts, Fran Ramovš Institute of the Slovenian Language, Ljubljana, Slovenia
{andrej.perdih,dejan.gabrovsek}@zrc-sazu.si

A. Authors’ affiliations

- ¹Society for Danish Language and Literature (DSL), Copenhagen, Denmark
{sn,tt}@dsl.dk
- ²Austrian Centre for Digital Humanities and Cultural Heritage, Austrian Academy of Sciences, Vienna, Austria
tanja.wissik@oeaw.ac.at
- ³Istituto di Linguistica Computazionale “A. Zampolli– CNR”, Pisa, Italy
{monica.monachini,fahad.khan, andrea.bellandi}@ilc.cnr.it
- ⁴Università di Pisa, Italy
i.pisanil@studenti.unipi.it
- ⁵Jožef Stefan Institute, Ljubljana, Slovenia
simon.krek@guest.arnes.si
- ⁶Research Institute for Linguistics, Budapest, Hungary
{lipp.veronika,varadi.tamas, simon.laszlo,gyorffy.andras}@nytud.hu
- ⁷Insight Centre for Data Analytics, National University of Ireland, Galway
{oksana.dereza,theodorus.fransen}@insight-centre.org,
d.cillessen1@nuigalway.ie
- ⁸Centre for Language Technology, University of Copenhagen, Denmark
{bspedersen, saolsen}@hum.ku.dk
- ⁹Dutch Language Institute, Leiden, the Netherlands
{carole.tiberius,tanneke.schoonheim}@ivdnt.org
- ¹⁰K Dictionaries, Tel Aviv, Israel
{yifat,maya,raya,dorielle}@Kdictionaries.com
- ¹¹Institute for Linguistic Studies of the Russian Academy of Sciences, St. Petersburg, Russia
kira.kovalenko@gmail.com
- ¹²DFKI GmbH, Multilinguality and Language Technology, Germany
declerck@dfki.de
- ¹³Institute of the Estonian Language, Estonia
{margit.langemets,jelena.kallas}@eki.ee
- ¹⁴Euskal Herriko Unibertsitatea, Universidad del País Vasco, Leioa, Spain
david.lindemann@ehu.eus,
mikelalon@gmail.com