

# Unsupervised, Knowledge-Free, and Interpretable Word Sense Disambiguation

Alexander Panchenko<sup>‡</sup>, Fide Marten<sup>‡</sup>, Eugen Ruppert<sup>‡</sup>, Stefano Faralli<sup>†</sup>,  
Dmitry Ustalov<sup>\*</sup>, Simone Paolo Ponzetto<sup>†</sup>, and Chris Biemann<sup>‡</sup>

<sup>‡</sup>Language Technology Group, Department of Informatics, Universität Hamburg, Germany

<sup>†</sup>Web and Data Science Group, Department of Informatics, Universität Mannheim, Germany

<sup>\*</sup>Institute of Natural Sciences and Mathematics, Ural Federal University, Russia

{panchenko,marten,ruppert,biemann}@informatik.uni-hamburg.de

{simone,stefano}@informatik.uni-mannheim.de

dmitry.ustalov@urfu.ru

## Abstract

Interpretability of a predictive model is a powerful feature that gains the trust of users in the correctness of the predictions. In word sense disambiguation (WSD), *knowledge-based* systems tend to be much more interpretable than *knowledge-free* counterparts as they rely on the wealth of manually-encoded elements representing word senses, such as hypernyms, usage examples, and images. We present a WSD system that bridges the gap between these two so far disconnected groups of methods. Namely, our system, providing access to several state-of-the-art WSD models, aims to be *interpretable* as a knowledge-based system while it remains completely *unsupervised* and *knowledge-free*. The presented tool features a Web interface for all-word disambiguation of texts that makes the sense predictions human readable by providing interpretable word sense inventories, sense representations, and disambiguation results. We provide a public API, enabling seamless integration.

## 1 Introduction

The notion of word sense is central to computational lexical semantics. Word senses can be either *encoded manually* in lexical resources or *induced automatically* from text. The former knowledge-based sense representations, such as those found in the BabelNet lexical semantic network (Navigli and Ponzetto, 2012), are easily interpretable by humans due to the presence of definitions, usage examples, taxonomic relations, related words, and images. The cost of such interpretability is that every element mentioned above is encoded

manually in one of the underlying resources, such as Wikipedia. Unsupervised knowledge-free approaches, e.g. (Di Marco and Navigli, 2013; Bartunov et al., 2016), require no manual labor, but the resulting sense representations lack the above-mentioned features enabling interpretability. For instance, systems based on sense embeddings are based on dense uninterpretable vectors. Therefore, the meaning of a sense can be interpreted only on the basis of a list of related senses.

We present a system that brings interpretability of the knowledge-based sense representations into the world of unsupervised knowledge-free WSD models. The contribution of this paper is the first *system* for word sense induction and disambiguation, which is unsupervised, knowledge-free, and interpretable at the same time. The system is based on the WSD approach of Panchenko et al. (2017) and is designed to reach interpretability level of knowledge-based systems, such as Babelfy (Moro et al., 2014), within an unsupervised knowledge-free framework. Implementation of the system is open source.<sup>1</sup> A live demo featuring several disambiguation models is available online.<sup>2</sup>

## 2 Related Work

In this section, we list prominent WSD systems with openly available implementations.

**Knowledge-Based and/or Supervised Systems**  
IMS (Zhong and Ng, 2010) is a supervised all-words WSD system that allows users to integrate additional features and different classifiers. By default, the system relies on the linear support vector machines with multiple features. The AutoExtend (Rothe and Schütze, 2015) approach can be used to learn embeddings for lexemes and synsets

<sup>1</sup><https://github.com/uhh-lt/wsd>

<sup>2</sup><http://jobimtext.org/wsd>

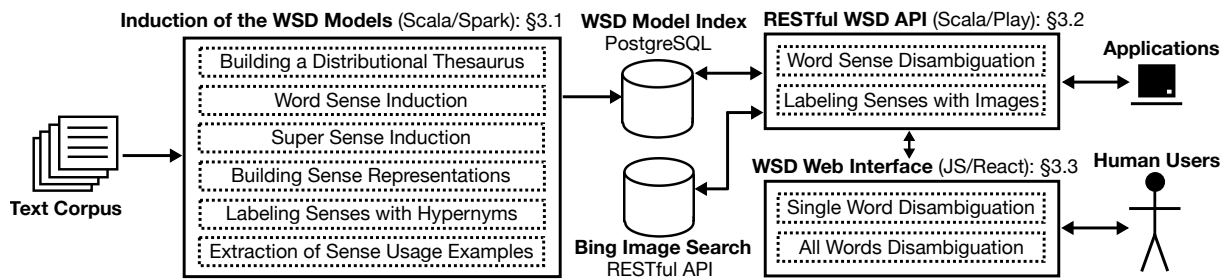


Figure 1: Software and functional architecture of the WSD system.

of a lexical resource. These representations were successfully used to perform WSD using the IMS.

DKPro WSD (Miller et al., 2013) is a modular, extensible Java framework for word sense disambiguation. It implements multiple WSD methods and also provides an interface to evaluation datasets. PyWSD<sup>3</sup> project also provides implementations of popular WSD methods, but these are implemented in the Python language.

Babelfy (Moro et al., 2014) is a system based on the BabelNet that implements a multilingual graph-based approach to entity linking and WSD based on the identification of candidate meanings using the densest subgraph heuristic.

**Knowledge-Free and Unsupervised Systems** Neelakantan et al. (2014) proposed a multi-sense extension of the Skip-gram model that features an open implementation. AdaGram (Bartunov et al., 2016) is a system that learns sense embeddings using a Bayesian extension of the Skip-gram model and provides WSD functionality based on the induced sense inventory. SenseGram (Pelevina et al., 2016) is a system that transforms word embeddings to sense embeddings via graph clustering and uses them for WSD. Other methods to learn sense embeddings were proposed, but these do not feature open implementations for WSD.

Among all listed systems, only Babelfy implements a user interface supporting interpretable visualization of the disambiguation results.

### 3 Unsupervised Knowledge-Free Interpretable WSD

This section describes (1) how WSD models are learned in an unsupervised way from text and (2) how the system uses these models to enable human interpretable disambiguation in context.

#### 3.1 Induction of the WSD Models

Figure 1 presents architecture of the WSD system. As one may observe, no human labor is used to learn interpretable sense representations and the corresponding disambiguation models. Instead, these are induced from the input text corpus using the JoBimText approach (Biemann and Riedl, 2013) implemented using the Apache Spark framework<sup>4</sup>, enabling seamless processing of large text collections. Induction of a WSD model consists of several steps. First, a graph of semantically related words, i.e. a distributional thesaurus, is extracted. Second, word senses are induced by clustering of an ego-network of related words (Biemann, 2006). Each discovered word sense is represented as a cluster of words. Next, the induced sense inventory is used as a pivot to generate sense representations by aggregation of the context clues of cluster words. To improve interpretability of the sense clusters they are labeled with hypernyms, which are in turn extracted from the input corpus using Hearst (1992) patterns. Finally, the obtained WSD model is used to retrieve a list of sentences that characterize each sense. Sentences that mention a given word are disambiguated and then ranked by prediction confidence. Top sentences are used as sense usage examples. For more details about the model induction process refer to (Panchenko et al., 2017). Currently, the following WSD models induced from a text corpus are available:

**Word senses based on cluster word features.** This model uses the cluster words from the induced word sense inventory as sparse features that represent the sense.

**Word senses based on context word features.** This representation is based on a sum of word vectors of all cluster words in the induced sense inventory weighted by distributional similarity scores.

<sup>3</sup><https://github.com/alvations/pywds>

<sup>4</sup><http://spark.apache.org>

### Super senses based on cluster word features.

To build this model, induced word senses are first globally clustered using the Chinese Whispers graph clustering algorithm (Biemann, 2006). The edges in this sense graph are established by disambiguation of the related words (Faralli et al., 2016; Ustalov et al., 2017). The resulting clusters represent semantic classes grouping words sharing a common hypernym, e.g. “animal”. This set of semantic classes is used as an automatically learned inventory of super senses: There is only one global sense inventory shared among all words in contrast to the two previous traditional “per word” models. Each semantic class is labeled with hypernyms. This model uses words belonging to the semantic class as features.

### Super senses based on context word features.

This model relies on the same semantic classes as the previous one but, instead, sense representations are obtained by averaging vectors of words sharing the same class.

## 3.2 WSD API

To enable fast access to the sense inventories and effective parallel predictions, the WSD models obtained at the previous step were indexed in a relational database.<sup>5</sup> In particular, each word sense is represented by its hypernyms, related words, and usage examples. Besides, for each sense, the database stores an aggregated context word representation in the form of a serialized object containing a sparse vector in the Breeze format.<sup>6</sup> During the disambiguation phrase, the input context is represented in the same sparse feature space and the classification is reduced to the computation of the cosine similarity between the context vector and the vectors of the candidate senses retrieved from the database. This back-end is implemented as a RESTful API using the Play framework.<sup>7</sup>

## 3.3 User Interface for Interpretable WSD

The graphical user interface of our system is implemented as a single page Web application using the React framework.<sup>8</sup> The application performs disambiguation of a text entered by a user. In particular, the Web application features two modes:

**Single word disambiguation mode** is illustrated in Figure 2. In this mode, a user specifies

an ambiguous word and its context. The output of the system is a ranked list of all word senses of the ambiguous word ordered by relevance to the input context. By default, only the best matching sense is displayed. The user can quickly understand the meaning of each induced sense by looking at the hypernym and the image representing the sense. Faralli and Navigli (2012) showed that Web search engines can be used to acquire information about word senses. We assign an image to each word in the cluster by querying an image search API<sup>9</sup> using a query composed of the ambiguous word and its hypernym, e.g. “jaguar animal”. The first hit of this query is selected to represent the induced word sense. Interpretability of each sense is further ensured by providing to the user the list of related senses, the list of the most salient context clues, and the sense usage examples (cf. Figure 2). Note that all these elements are obtained without manual intervention.

Finally, the system provides the reasons behind the sense predictions by displaying context words triggered the prediction. Each common feature is clickable, so a user is able to trace back sense cluster words containing this context feature.

**All words disambiguation mode** is illustrated in Figure 3. In this mode, the system performs disambiguation of all nouns and entities in the input text. First, the text is processed with a part-of-speech and a named entity taggers.<sup>10</sup> Next, each detected noun or entity is disambiguated in the same way as in the single word disambiguation mode described above, yet the disambiguation results are represented as annotations of a running text. The best matching sense is represented by a hypernym and an image as depicted in Figure 3. This mode performs “semantification” of a text, which can, for instance, assist language learners with the understanding of a text in a foreign language: Meaning of unknown to the learner words can be deduced from hypernyms and images.

## 4 Evaluation

In our prior work (Panchenko et al., 2017), we performed a thorough evaluation of the method implemented in our system on two datasets showing the state-of-the-art performance of the approach as compared to other unsupervised knowledge-free

<sup>5</sup><https://www.postgresql.org>

<sup>6</sup><https://github.com/scalanlp/breeze>

<sup>7</sup><https://www.playframework.com>

<sup>8</sup><https://facebook.github.io/react>

<sup>9</sup><https://azure.microsoft.com/en-us/services/cognitive-services/search>

<sup>10</sup><http://www.scalanlp.org>

Sentence  
Jaguar is a large spotted predator of tropical America similar to the leopard. **A**

Word  
Jaguar **B**

Model  
Word Senses based on Cluster Word Features **C**

**PREDICT SENSE** **RANDOM SAMPLE**

---

**1. jaguar (animal)**  
 Similarity score: 0.00184 / Confidence: 99.87% / Sense ID: jaguar#0 / BabelNet ID: bn:00033987n

Hypernyms: animal **D**, wildlife, bird, mammal

Sample sentences:  
 The **jaguar**, a compact and well-muscled animal, is the largest cat in the New World.  
**Jaguar** may leap onto the back of the prey and sever the cervical vertebrae, immobilizing the target.

Cluster words: lion, tiger, leopard, wolf, monkey, otter, crocodile, alligator, deer, cat, elephant, fox, eagle, owl, snake

Context words:  
 elephant: 0.012, tiger: 0.012, fox: 0.0099, wolf: 0.0097, cub: 0.0086, monkey: 0.0083, leopard: 0.0074, eagle: 0.0062  
 den: 0.0043, elk: 0.0040, 32078 more not shown

Matching features: leopard: 0.0011, predator: 0.00040, spotted: 0.00038, large: 0.0000041, similar: 0.0000015, tropical: 5.6e-7, america: 2.0e-7

**BABELNET LINK** **F** **SHOW LESS** **E**

Figure 2: Single word disambiguation mode: results of disambiguation of the word “Jaguar” (B) in the sentence “*Jaguar* is a large spotted predator of tropical America similar to the leopard.” (A) using the WSD disambiguation model based on cluster word features (C). The predicted sense is summarized with a hypernym and an image (D) and further represented with usage examples, semantically related words, and typical context clues. Each of these elements is extracted automatically. The reasons of the predictions are provided in terms of common sparse features of the input sentence and a sense representation (E). The induced senses are linked to BabelNet using the method of Faralli et al. (2016) (F).

Sentence  
Jaguar is a large spotted predator of tropical America similar to the leopard. **A**

Model  
Word Senses based on Cluster Word Features **C**

**DISAMBIGUATE SENTENCE** **RANDOM SAMPLE**

---

**Detected Entities**  
 The system has detected these entities in the given sentence.




 animal		 animal		 country
Jaguar <b>D</b>	is a large spotted	predator <b>D</b>	of tropical	America <b>D</b>

Figure 3: All words disambiguation mode: results of disambiguation of all nouns in a sentence.

# Words	# Senses	Avg. Polysemy	# Contexts
863	2,708	3.13	11,712

Table 1: Evaluation dataset based on BabelNet.

methods for WSD, including participants of the SemEval 2013 Task 13 (Jurgens and Klapaftis, 2013) and two unsupervised knowledge-free WSD systems based on word sense embeddings (Bartunov et al., 2016; Pelevina et al., 2016). These evaluations were based on the “lexical sample” setting, where the system is expected to predict a sense identifier of the ambiguous word.

In this section, we perform an extra evaluation that assesses how well hypernyms of ambiguous words are assigned in context by our system. Namely, the task is to assign a correct hypernym of an ambiguous word, e.g. “animal” for the word “Jaguar” in the context “*Jaguar* is a large spotted predator of tropical America”. This task does not depend on a fixed sense inventory and evaluates at the same time WSD performance and the quality of the hypernymy labels of the induced senses.

#### 4.1 Dataset

In this experiment, we gathered a dataset consisting of definitions of BabelNet 3.7 senses of 1,219 frequent nouns.<sup>11</sup> In total, we collected 56,003 sense definitions each labeled with gold hypernyms coming from the IsA relations of BabelNet.

The average polysemy of words in the gathered dataset was 15.50 senses per word as compared to 2.34 in the induced sense inventory. This huge discrepancy in granularities lead to the fact that some test sentences cannot be correctly predicted by definition: some (mostly rare) BabelNet senses simply have no corresponding sense in the induced inventory. To eliminate the influence of this idiosyncrasy, we kept only sentences that contain at least one common hypernym with all hypernyms of all induced senses. The statistics of the resulting dataset are presented in Table 1, it is available in the project repository.

#### 4.2 Evaluation Metrics

WSD performance is measured using the accuracy with respect to the sentences labeled with the direct hypernyms (*Hypers*) or an extended set of hypernym including hypernyms of hypernyms (*Hy-*

<sup>11</sup>Most of the nouns come from the TWSI (Biemann, 2012) dataset, while the remaining nouns were manually selected.

WSD Model		Accuracy	
Inventory	Features	Hypers	HyperHypers
Word Senses	Random	0.257	0.610
Word Senses	MFS	0.292	0.682
Word Senses	Cluster Words	0.291	0.650
Word Senses	Context Words	<b>0.308</b>	<b>0.686</b>
Super Senses	Random	0.001	0.001
Super Senses	MFS	0.001	0.001
Super Senses	Cluster Words	<b>0.174</b>	<b>0.365</b>
Super Senses	Context Words	0.086	0.188

Table 2: Performance of the hypernymy labeling in context on the BabelNet dataset.

*perHypers*). A correct match occurs when the predicted sense has at least one common hypernym with the gold hypernyms of the target word in a test sentence.

#### 4.3 Discussion of Results

**Word Senses.** All evaluated models outperform both random and most frequent sense baselines, see Table 2. The latter picks the sense that corresponds to the largest sense cluster (Panchenko et al., 2017). In the case of the traditional “per word” inventories, the model based on the context features outperform the models based on cluster words. While sense representations based on the clusters of semantically related words contain highly accurate features, such representations are sparse as one sense contains at most 200 features. As the result, often the model based on the cluster words contain no common features with the features extracted from the input context. The sense representations based on the aggregated context clues are much less sparse, which explains their superior performance.

**Super Senses.** In the case of the super sense inventory, the model based solely on the cluster words yielded better results than the context-based model. Note here that (1) the clusters that represent super senses are substantially larger than word sense clusters and thus less sparse, (2) words in the super sense clusters are unweighted in contrast to word sense cluster, thus averaging of word vectors is more noise-prone. Besides, the performance scores of the models based on the super sense inventories are substantially lower compared to their counterparts based on the traditional “per word” inventories. Super sense models are able to perform classification for any unknown word missing in the training corpus, but their disambiguation task is more complex (the models need

to choose one of 712 classes as compared to an average of 2–3 classes for the “per word” inventories). This is illustrated by the near-zero scores of the random and the MFS baselines for this model.

## 5 Conclusion

We present the first openly available word sense disambiguation system that is unsupervised, knowledge-free, and interpretable at the same time. The system performs extraction of word and super sense inventories from a text corpus. The disambiguation models are learned in an unsupervised way for all words in the corpus on the basis on the induced inventories. The user interface of the system provides efficient access to the produced WSD models via a RESTful API or via an interactive Web-based graphical user interface. The system is available online and can be directly used from external applications. The code and the WSD models are open source. Besides, in-house deployments of the system are made easy due to the use of the Docker containers.<sup>12</sup> A prominent direction for future work is supporting more languages and establishing cross-lingual sense links.

## Acknowledgments

We acknowledge the support of the DFG under the “JOIN-T” project, the RFBR under project no. 16-37-00354 mol\_a, Amazon via the “AWS Research Grants” and Microsoft via the “Azure for Research” programs. Finally, we also thank four anonymous reviewers for their helpful comments.

## References

- S. Bartunov, D. Kondrashkin, A. Osokin, and D. Vetrov. 2016. [Breaking Sticks and Ambiguities with Adaptive Skip-gram](#). In *Proc. AISTATS*. Cadiz, Spain, pp. 130–138.
- C. Biemann. 2006. [Chinese Whispers: An Efficient Graph Clustering Algorithm and Its Application to Natural Language Processing Problems](#). In *Proc. TextGraphs*. New York, NY, USA, pp. 73–80.
- C. Biemann. 2012. [Turk Bootstrap Word Sense Inventory 2.0: A Large-Scale Resource for Lexical Substitution](#). In *Proc. LREC*. Istanbul, Turkey, pp. 4038–4042.
- C. Biemann and M. Riedl. 2013. [Text: now in 2D! A framework for lexical expansion with contextual similarity](#). *Journal of Language Modelling* 1(1):55–95.
- A. Di Marco and R. Navigli. 2013. [Clustering and Diversifying Web Search Results with Graph-Based Word Sense Induction](#). *Computational Linguistics* 39(3):709–754.
- S. Faralli and R. Navigli. 2012. [A New Minimally-Supervised Framework for Domain Word Sense Disambiguation](#). In *Proc. EMNLP-CoNLL*. Jeju Island, Korea, pp. 1411–1422.
- S. Faralli, A. Panchenko, C. Biemann, and S. P. Ponzetto. 2016. [Linked Disambiguated Distributional Semantic Networks](#). In *Proc. ISWC, Part II*. Kobe, Japan, pp. 56–64.
- M. A. Hearst. 1992. [Automatic Acquisition of Hyponyms from Large Text Corpora](#). In *Proc. COLING*. Nantes, France, pp. 539–545.
- D. Jurgens and I. Klapaftis. 2013. [SemEval-2013 Task 13: Word Sense Induction for Graded and Non-Graded Senses](#). In *Proc. SemEval*. Atlanta, GA, USA, pp. 290–299.
- T. Miller et al. 2013. [DKPro WSD: A Generalized UIMA-based Framework for Word Sense Disambiguation](#). In *Proc. ACL*. Sofia, Bulgaria, pp. 37–42.
- A. Moro, A. Raganato, and R. Navigli. 2014. [Entity Linking meets Word Sense Disambiguation: A Unified Approach](#). *Transactions of the Association for Computational Linguistics* 2:231–244.
- R. Navigli and S. P. Ponzetto. 2012. [BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network](#). *Artificial Intelligence* 193:217–250.
- A. Neelakantan, J. Shankar, A. Passos, and A. McCallum. 2014. [Efficient Non-parametric Estimation of Multiple Embeddings per Word in Vector Space](#). In *Proc. EMNLP*. Doha, Qatar, pp. 1059–1069.
- A. Panchenko, E. Ruppert, S. Faralli, S. P. Ponzetto, and C. Biemann. 2017. [Unsupervised Does Not Mean Uninterpretable: The Case for Word Sense Induction and Disambiguation](#). In *Proc. EACL*. Valencia, Spain, pp. 86–98.
- M. Pelevina, N. Arefiev, C. Biemann, and A. Panchenko. 2016. [Making Sense of Word Embeddings](#). In *Proc. RepLANLP*. Berlin, Germany, pp. 174–183.
- S. Rothe and H. Schütze. 2015. [AutoExtend: Extending Word Embeddings to Embeddings for Synsets and Lexemes](#). In *Proc. ACL-IJCNLP*. Beijing, China, pp. 1793–1803.
- D. Ustalov, A. Panchenko, and C. Biemann. 2017. [Watset: Automatic Induction of Synsets from a Graph of Synonyms](#). In *Proc. ACL*. Vancouver, Canada.
- Z. Zhong and H. T. Ng. 2010. [It Makes Sense: A Wide-Coverage Word Sense Disambiguation System for Free Text](#). In *Proc. ACL Demos*. Uppsala, Sweden, pp. 78–83.

<sup>12</sup><https://www.docker.com>