# Chapter 28
# Universal Semantic Annotator

Roberto Navigli, Riccardo Orlando, Cesare Campagnano, and Simone Conia

**Abstract** Explicit semantic knowledge has often been considered a necessary ingredient to enable the development of intelligent systems. However, current state-of-the-art tools for the automatic extraction of such knowledge often require expert understanding of the complex techniques used in lexical and sentence-level semantics and their linguistic theories. To overcome this limitation and lower the barrier to entry, we present the Universal Semantic Annotator (USeA) ELG pilot project, which offers a transparent way to automatically provide high-quality semantic annotations in 100 languages through state-of-the-art models, making it easy to exploit semantic knowledge in real-world applications.

## 1 Overview and Objectives of the Pilot Project

Natural Language Processing (NLP) is the field of Artificial Intelligence (AI) which aims at enabling computers to process, understand and generate text in the same way as we humans do. Although AI systems are nowadays able to process massive amounts of text, they are still far from achieving true Natural Language Understanding (NLU). Indeed, current systems still struggle in explicitly identifying and extracting the meaning or semantics conveyed by a text of interest. Nonetheless, the integration of explicit semantics has already been successfully exploited in a wide array of downstream tasks that span multiple areas of AI from NLP with information retrieval, question answering, text summarisation, and machine translation, to computer vision with visual semantic role labeling and situation recognition. Unfortunately, expert knowledge of lexical semantics, sentence-level semantics and complex deep learning techniques often becomes a roadblock in the integration of explicit semantic information into downstream tasks and real-world applications, especially in multilingual scenarios. To lower the entry point for semantic knowledge integration into multilingual applications, we present the Universal Semantic Anno-

Roberto Navigli · Riccardo Orlando · Cesare Campagnano · Simone Conia
Sapienza University of Rome, Italy, navigli@diag.uniroma1.it, orlando@diag.uniroma1.it,
campagnano@di.uniroma1.it, conia@di.uniroma1.it

tator (USeA) project, the first unified API for three core tasks in NLU: Word Sense Disambiguation (WSD), Semantic Role Labeling (SRL), and Abstract Meaning Representation (AMR) parsing. With USeA, we offer a simple yet efficient way to use state-of-the-art multilingual models within a single framework accessible via REST API, browsers, and programmatically. This will ease the integration of NLU models in NLP pipelines (also for low-resource languages), allowing them to exploit explicit semantic information to improve their performance.

## 2  Methodology

USeA is the first unified set of APIs for high-performance multilingual NLU, supporting 100 languages. USeA employs state-of-the-art multilingual neural networks to provide automatic semantic annotations for WSD, SRL and AMR Parsing.

**Word Sense Disambiguation (WSD)**  is the task of associating a word in context with its most appropriate sense from a sense inventory (Bevilacqua et al. 2021b). USeA provides word sense labels using an improved version of the state-of-the-art WSD model proposed by Conia and Navigli (2021), which, differently from other ready-to-use tools for WSD based on graph-based heuristics (Moro et al. 2014; Scozzafava et al. 2020) or non-neural models (Papandrea et al. 2017), is built on top of a Transformer encoder. Crucially, thanks to BabelNet 5 (Navigli et al. 2021), a multilingual encyclopedic dictionary, USeA is able to disambiguate text in 100 languages.

**Semantic Role Labeling (SRL)**  is the task of answering the question "Who did What, to Whom, Where, When, and How?" (Màrquez et al. 2008), providing a structured and explicit representation of the underlying semantics of a sentence. Differently from other available SRL systems, USeA encapsulates an improved version of the neural model introduced by Conia et al. (2021a), which performs state-of-the-art cross-lingual SRL with heterogeneous linguistic inventories.

**Abstract Meaning Representation (AMR) parsing**  is the task of capturing the semantics of a sentence through a rooted directed acyclic graph, with nodes representing concepts and edges representing their relations (Banarescu et al. 2013). USeA offers a multilingual version of SPRING (Bevilacqua et al. 2021a), a recent state-of-the-art, end-to-end system for Text-to-AMR generation.

## 3  Implementation

The USeA pipeline is organised in five self-contained modules that are transparent to the end user, as shown in Figure 1.

**Orchestrator Module.**  The Orchestrator Module is the core of USeA and serves as an entry point for the semantic API. Being an end-to-end system, the end user
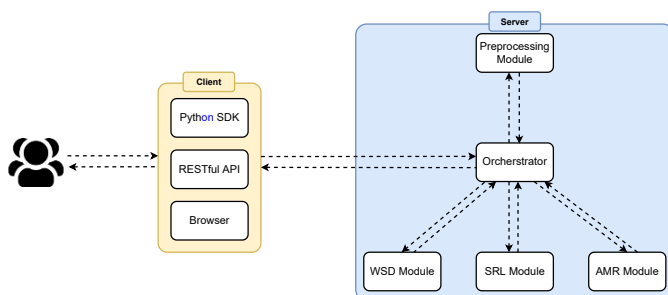
**Fig. 1** USeA architecture: a user sends text to the USeA server and receives semantic information; in the server, the orchestrator processes the input using task-specific modules

is only required to send raw text to our service. The input text is then processed by the Preprocessing Module and the result sent to the WSD, SRL and AMR Parsing modules. In particular, since the SRL and AMR Parsing tasks are more demanding, we offload the WSD module to CPU and run SRL and AMR Parsing requests on GPU to optimise hardware usage. The responses from the three semantic modules are then combined and sent back to the end user.

**Preprocessing Module.** The preprocessing module takes care of producing the preprocessing information that is usually needed by NLP systems, i. e., language identification, document splitting, tokenisation, lemmatisation, and part-of-speech tagging. In order to support as many languages as possible while keeping low hardware requirements, the preprocessing module is built around Trankit (Nguyen et al. 2021) and supports 100 languages with a single model.

**WSD Module.** We developed AMuSE-WSD (Orlando et al. 2021) as our WSD module. Its neural architecture is based on XLM-RoBERTa (Conneau et al. 2020), a multilingual Transformer model. More specifically, given a word in context, the WSD module i) builds a contextualised representation of the word using the hidden states of XLM-RoBERTa, ii) applies a non-linear transformation to obtain a sense-specific representation, and iii) computes the output score distribution over all the possible senses of the input word.

**SRL Module.** InVeRo-XL (Conia et al. 2021b) is the SRL system we developed for USeA. Similarly to the WSD module, the SRL module is also based on XLM-RoBERTa. In particular, given an input sentence, the SRL module i) builds a sequence of contextualised word representations using the hidden states of XLM-RoBERTa, ii) identifies and disambiguates each predicate in the sentence, and iii) for each predicate, produces its arguments and their semantic roles.

**AMR Parsing Module.** The AMR Parsing Module is heavily based on SPRING (Blloshmi et al. 2021), which we extended to support multiple languages. SPRING is a sequence-to-sequence Transformer model that operates as a parser by "translating" an input sentence into a linearised AMR graph. We extend SPRING to support 100 languages by replacing BART with the multilingual version of T5.

| | English datasets | | | | | | Multilingual datasets | | |
|---|---|---|---|---|---|---|---|---|---|
| | Se2 | Se3 | Se07 | Se13 | Se15 | All | Se13 | Se15 | Xl-Wsd |
| Moro et al. (2014) | 67.0 | 63.5 | 51.6 | 66.4 | 70.3 | 65.5 | 65.6 | – | 52.9 |
| Papandrea et al. (2017) | 73.8 | 70.8 | 64.2 | 67.2 | 71.5 | – | – | – | – |
| Scozzafava et al. (2020) | 71.6 | 72.0 | 59.3 | 72.2 | 75.8 | 71.7 | 73.2 | 66.2 | 57.7 |
| USeA (WSD) | **77.8** | **76.0** | **72.1** | **77.7** | **81.5** | **77.5** | **76.8** | **73.0** | **66.2** |

**Table 1** English WSD results in F1 scores on Senseval-2 (SE2), Senseval-3 (SE3), SemEval-2007 (SE07), SemEval-2013 (SE13), SemEval-2015 (SE15), and the concatenation of the datasets (ALL); we also include results on multilingual WSD in SemEval-2013 (DE, ES, FR, IT), SemEval-2015 (IT, ES), and XL-WSD (average over 17 languages, English excluded)

| | Catalan | Czech | German | English | Spanish | Chinese |
|---|---|---|---|---|---|---|
| AllenNLP's SRL demo | – | – | – | 86.5 | – | – |
| InVeRo | – | – | – | 86.2 | – | – |
| USeA (SRL) | **83.3** | **85.9** | **87.0** | **86.8** | **81.8** | **84.9** |

**Table 2** Comparison between USeA and other recent automatic tools for SRL; F1 scores on argument labeling with pre-identified predicates on the CoNLL-2012 English test set and the CoNLL-2009 test sets converted from dependency-based to span-based

## 4 Evaluation

USeA offers state-of-the-art models for multilingual WSD, SRL and AMR Parsing. Here, we report its results on standard gold benchmarks for each task.

**Results in WSD.** We evaluate our WSD Module against other disambiguation tools on gold standard benchmarks for English and multilingual WSD, covering 17 languages. The results (Table 1) show that USeA outperforms its competitors by a wide margin, especially in multilingual WSD (+8.5% in F1 Score on XL-WSD).

**Results in SRL.** We report the performance of our SRL Module on two gold standard benchmarks for SRL, CoNLL-2009[1] and CoNLL-2012, covering six languages. USeA is the first package to provide annotations in languages other than English while also outperforming its competitors in English (Table 2).

**Results in AMR Parsing.** Finally, we examine the performance of our AMR Parsing Module on AMR 3.0[2], which is currently the largest AMR-annotated corpus. Even though USeA supports 100 languages, it is still competitive with other recently proposed English-only AMR parsing systems (Table 3).

---

[1] The CoNLL-2009 dataset was originally intended for dependency-based SRL. We convert dependency-based annotations to span-based annotations using the gold syntactic trees.

[2] https://catalog.ldc.upenn.edu/LDC2020T02

|                                    | SMATCH |
|------------------------------------|--------|
| Lyu et al. (2021)                  | 75.8   |
| Zhou et al. (2021)                 | 81.2   |
| SPRING (Bevilacqua et al. 2021a)   | 83.0   |
| USeA (AMR-Parsing)                 | 80.9   |

**Table 3** SMATCH score obtained by USeA compared with recent literature on AMR 3.0 (English)

# 5  Conclusions and Results of the Pilot Project

We presented the USeA project, providing an overview on its objectives and on how we worked towards achieving them. We hope that USeA will represent a useful tool for the integration of explicit semantic knowledge – word meanings, semantic role labels, and graph-like semantic representations – into real-world applications.

# References

Banarescu, Laura, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider (2013). "Abstract Meaning Representation for Sembanking". In: *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pp. 178–186. URL: https://aclanthology.org/W13-2322.

Bevilacqua, Michele, Rexhina Blloshmi, and Roberto Navigli (2021a). "One SPRING to Rule Them Both: Symmetric AMR Semantic Parsing and Generation without a Complex Pipeline". In: *Proc. of AAAI* 35.14, pp. 12564–12573. URL: https://ojs.aaai.org/index.php/AAAI/article/view/17489.

Bevilacqua, Michele, Tommaso Pasini, Alessandro Raganato, and Roberto Navigli (2021b). "Recent Trends in Word Sense Disambiguation: A Survey". In: *Proc. of IJCAI-21*, pp. 4330–4338. DOI: 10.24963/ijcai.2021/593.

Blloshmi, Rexhina, Michele Bevilacqua, Edoardo Fabiano, Valentina Caruso, and Roberto Navigli (2021). "SPRING Goes Online: End-to-End AMR Parsing and Generation". In: *Proceedings of EMNLP*, pp. 134–142. URL: https://aclanthology.org/2021.emnlp-demo.16.

Conia, Simone, Andrea Bacciu, and Roberto Navigli (2021a). "Unifying Cross-Lingual Semantic Role Labeling with Heterogeneous Linguistic Resources". In: *Proceedings of NAACL*, pp. 338–351. URL: https://www.aclweb.org/anthology/2021.naacl-main.31.

Conia, Simone and Roberto Navigli (2021). "Framing Word Sense Disambiguation as a Multi-Label Problem for Model-Agnostic Knowledge Integration". In: *Proceedings of EACL*, pp. 3269–3275. URL: https://www.aclweb.org/anthology/2021.eacl-main.286.

Conia, Simone, Riccardo Orlando, Fabrizio Brignone, Francesco Cecconi, and Roberto Navigli (2021b). "InVeRo-XL: Making Cross-Lingual Semantic Role Labeling Accessible with Intelligible Verbs and Roles". In: *Proceedings of EMNLP*, pp. 319–328. URL: https://aclanthology.org/2021.emnlp-demo.36/.

Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov (2020). "Unsupervised Cross-lingual Representation Learning at Scale". In: *Proceedings of the 58th*

*Annual Meeting of the Association for Computational Linguistics*. Ed. by Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault. ACL, pp. 8440–8451. DOI: 10.18653/v1/2020.acl-main.747. URL: https://www.aclweb.org/anthology/2020.acl-main.747.

Lyu, Chunchuan, Shay B. Cohen, and Ivan Titov (2021). "A Differentiable Relaxation of Graph Segmentation and Alignment for AMR Parsing". In: *Proc. of EMNLP*, pp. 9075–9091. URL: https://aclanthology.org/2021.emnlp-main.714.

Màrquez, Lluís, Xavier Carreras, Kenneth C. Litkowski, and Suzanne Stevenson (2008). "Semantic Role Labeling: An Introduction to the Special Issue". In: *Comp. Linguistics* 34.2, pp. 145–159. URL: https://aclanthology.org/J08-2001.

Moro, Andrea, Alessandro Raganato, and Roberto Navigli (2014). "Entity Linking meets Word Sense Disambiguation: A Unified Approach". In: *TACL* 2, pp. 231–244. URL: https://aclanthology.org/Q14-1019.

Navigli, Roberto, Michele Bevilacqua, Simone Conia, Dario Montagnini, and Francesco Cecconi (2021). "Ten Years of BabelNet: A Survey". In: *Proc. of IJCAI-21*, pp. 4559–4567. DOI: 10.24963/ijcai.2021/620.

Nguyen, Minh Van, Viet Dac Lai, Amir Pouran Ben Veyseh, and Thien Huu Nguyen (2021). "Trankit: A Light-Weight Transformer-based Toolkit for Multilingual Natural Language Processing". In: *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. ACL, pp. 80–90. DOI: 10.18653/v1/2021.eacl-demos.10. URL: https://aclanthology.org/2021.eacl-demos.10.

Orlando, Riccardo, Simone Conia, Fabrizio Brignone, Francesco Cecconi, and Roberto Navigli (2021). "AMuSE-WSD: An All-in-one Multilingual System for Easy Word Sense Disambiguation". In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. ACL, pp. 298–307. DOI: 10.18653/v1/2021.emnlp-demo.34. URL: https://aclanthology.org/2021.emnlp-demo.34.

Papandrea, Simone, Alessandro Raganato, and Claudio Delli Bovi (2017). "SupWSD: A Flexible Toolkit for Supervised Word Sense Disambiguation". In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Copenhagen, Denmark: ACL, pp. 103–108. DOI: 10.18653/v1/D17-2018. URL: https://www.aclweb.org/anthology/D17-2018.

Scozzafava, Federico, Marco Maru, Fabrizio Brignone, Giovanni Torrisi, and Roberto Navigli (2020). "Personalized PageRank with Syntagmatic Information for Multilingual Word Sense Disambiguation". In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. ACL, pp. 37–46. DOI: 10.18653/v1/2020.acl-demos.6.

Zhou, Jiawei, Tahira Naseem, Ramón Fernandez Astudillo, and Radu Florian (2021). "AMR Parsing with Action-Pointer Transformer". In: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL, pp. 5585–5598. DOI: 10.18653/v1/2021.naacl-main.443. URL: https://aclanthology.org/2021.naacl-main.443.