

# VMWE identification with models trained on GUD (a UDv.2 treebank of Standard Modern Greek)

Stella Markantonatou<sup>1,2</sup>, Vivian Stamou<sup>1</sup>, Stavros Bompolas<sup>1</sup>,  
Katerina Anastasopoulou<sup>3</sup>, Irianna Vasileiadi Linardaki<sup>3</sup>, Konstantinos Diamantopoulos<sup>3</sup>,  
Yannis Kazos<sup>1,4</sup>, Antonios Anastasopoulos<sup>1,5</sup>

<sup>1</sup>Archimedes, Athena Research Center, Greece

<sup>2</sup>ILSP, Athena Research Center <sup>3</sup>Department of Informatics and Telecommunications, NKUA

<sup>4</sup>National Technical University of Athens, NTUA <sup>5</sup>George Mason University, USA

Correspondence: marks@athenarc.gr

## Abstract

UD\_Greek-GUD (GUD) is the most recent *Universal Dependencies* (UD) treebank for *Standard Modern Greek* (SMG) and the first SMG UD treebank to annotate *Verbal Multiword Expressions* (VMWEs). GUD contains material from fiction texts and various sites that use colloquial SMG. We describe the special annotation decisions we implemented with GUD, the pipeline we developed to facilitate the active annotation of new material, and we report on the method we designed to evaluate the performance of models trained on GUD as regards VMWE identification tasks.

## 1 Introduction

*Multiword expressions* (MWEs) pose significant challenges in both linguistic annotation and computational processing due to their semantic and structural idiosyncratic properties. Previous research on MWEs in Modern Greek has explored their theoretical properties (2024) and led to the development of lexical resources documenting their semantic and syntactic behavior (Markantonatou et al., 2019). Lexicographic and annotation studies have examined various semantic, pragmatic, and methodological aspects of MWE (Giouli et al., 2019). Computational approaches have also contributed to MWE processing, focusing on MWE extraction (Stamou et al., 2020b) and multilingual parsing (Michou and Seretan, 2009; Foufi et al., 2019), as well as the evaluation of MWE discovery methods (Stamou et al., 2020a). However, despite these advancements, systematic treatments of MWEs within syntactic parsing adapted to Modern Greek remain relatively rare, with existing studies employing symbolic frameworks such as *Lexical Functional Grammar* (LFG; Samaridi and Markantonatou, 2014).

To address this gap, we introduce UD\_Greek-GUD (GUD)—a new Universal

Dependencies (UD v2) treebank (de Marneffe et al., 2021) for SMG. GUD integrates rich morphological and syntactic annotations with explicit *verbal MWEs* (VMWEs) annotation, in the spirit of the PARSEME guidelines (Savary et al., 2018). We outline specific linguistic decisions regarding tokenization, contractions, functional words, and diminutives/augmentatives, and propose a novel annotation strategy that integrates VMWEs information directly into the syntactic layer of the CoNLL-U format.

Building on this, we experiment with an annotation method that eventually encodes VMWEs as dependency sub-relations, facilitating automatic identification through syntactic parsing. The approach is shown to be promising for computational processing; open issues are the identification of nested, or overlapping expressions and of discontinuous MWEs (Constant et al., 2017).

This paper explores these challenges, evaluates their implications, and outlines ongoing efforts toward improved evaluation and their practical integration into syntactic parsing frameworks, offering new perspectives for linguistic annotation and computational processing in Greek and beyond.

## 2 Materials and annotation method

UD\_Greek-GDT (henceforth GDT; Prokopidis and Papageorgiou (2017)) is the first UD treebank for SMG. Both GUD and GDT have been manually annotated for morphology and syntax, with GUD additionally annotated for VMWEs.

To develop GUD, a total of 1,807 sentences (25,493 tokens) were randomly selected from fiction texts in SMG. Additionally, 723 sentences (13,111 tokens), annotated specifically for VMWEs, were retrieved from IDION (Markantonatou et al., 2019), an open-source web database of SMG VMWEs, resulting in a combined corpus of 2,530 sentences. These VMWE usage examples

have been collected over the past 15 years through Google searches from social media, football sites, and other sources where colloquial SMG is used. The ArboratorGrew tool<sup>1</sup> was used to implement the annotation.

The annotation of GUD was carried out by graduate students in Language Technology (2021-2024) under the supervision of two of the authors. It proceeded in three rounds during this period. In the first round, students edited morphological and syntactic annotations obtained from models trained on GDT, developed morphological guidelines from scratch, and revised and enriched the syntactic guidelines originally produced by the GDT annotators. In the second round, one of the authors reviewed all annotated material and unified the guidelines. In the third round, the authors re-edited GUD and refined the material exemplifying VMWEs based on the established guidelines. Annotation decisions were reached through discussions and consensus among the annotators.

### 3 What is new about GUD

GUD and GDT share the same tokenization and word segmentation guidelines but differ notably in terms of morphological and syntactic annotation.

**Morphological annotation:** The main differences in the morphological annotation of the two treebanks are:

1. *να na* ‘to’, *που pou* ‘that’ (occurring  $\geq 300$  and  $\leq 200$  in GUD, respectively) introduce sentential complements of verbs. Additionally, *pou* introduces relative clauses and certain types of adverbial clauses, while *na* is also used to form periphrastic imperatives, express wishes and curses, and in other constructions, such as pointing to something. In GDT *pou* is tagged as PRON, and *na* as AUX. As shown in example (1a), GUD tags *pou* as SCONJ when it introduces sentential complements of verbs (Joseph and Philippaki-Warburton, 1987; Joseph, 1981) and as PRON (1b) when it introduces relative clauses. For *na*, GUD uses the tag SCONJ when it introduces sentential complements of verbs (2a), the tag AUX when it introduces main clauses expressing orders, wishes, curses, etc. (2b), and the tag PART in clauses with deixis (2c).

2. GUD adheres closely to the UD.v2 morphological guidelines and assigns the DET tag to 39

|      |  |                                 |                                 |
|------|--|---------------------------------|---------------------------------|
| (1a) | Χαίρομαι<br>chairomai<br>be.glad.1SG.PRS | που<br>pou<br>that.SCONJ        | ήρθες<br>irthes<br>come.2SG.PST |
| (1b) | Αυτός<br>aftos<br>he.NOM                 | που<br>pou<br>that.PRON         | ήρθε<br>irthe<br>come.3SG.PST   |
| (2a) | Ελπίζω<br>elpizo<br>hope.1SG.PRS         | να<br>na<br>to.SCONJ            | έρθεις<br>erthis<br>come.2SG    |
| (2b) | Να<br>na<br>to.AUX                       | έρθεις<br>erthis<br>come.2SG    |                                 |
| (2c) | Να<br>na<br>there.PART                   | ήρθαν<br>irthan<br>come.3SG.PST |                                 |

lemmas, whereas GDT assigns it to 17.

3. Unlike GDT, GUD annotates diminutives and augmentatives on nouns, adjectives, and adverbs. As shown in example (3), GUD assigns the lemma without a diminutive (or augmentative) affix to both forms with (3a) and without (3b) a diminutive (or augmentative) affix.

|      |   |      |  |
|------|---|------|--|
| (3a) | λαμπάκι<br>Lemma=λάμπα<br>UPOS=NOUN<br>Case=Nom<br>Gender=Neut<br>Number=Sing<br>Degree=Dim | (3b) | λάμπα<br>Lemma=λάμπα<br>UPOS=NOUN<br>Case=Nom<br>Gender=Fem<br>Number=Sing |
|------|---|------|--|

4. GUD tags passive participles as VERB and GDT as ADJ. In GUD, participles not related to a verb in use are tagged as ADJ.
5. GUD does not use the case DAT tag because the dative belongs to the diachrony of Greek (Anagnostopoulou and Sevdali, 2020); spoken SMG uses the dative only in fixed expressions.
6. GUD tags fossilized forms from the diachrony of Greek with the UPOS X.
7. At the time of GUD’s development, GDT and GUD used different sets of auxiliaries.
8. Unlike GDT, GUD provides an exhaustive annotation of both periphrastic and morphological degrees of comparison in SMG, following the established UD guidelines for comparative constructions.<sup>2</sup>

**Syntactic annotation:** The relations

<sup>2</sup><https://universaldependencies.org/workgroups/newdoc/comparatives.html>

<sup>1</sup><https://arborator.github.io/>

advcl:relcl, dislocated, and nsubj:outer are specific to GUD. Unlike GDT, GUD does not employ the dep relation.

GUD, but not GDT, analyzes the contracted forms *ston*, *stin*, *sto*, *stous*, *stis*, *sta*, which arise from the fusion of two pronouns: one in the genitive case and another in the accusative case. This phenomenon is linked to the broader loss of the dative-genitive distinction in SMG (Anagnostopoulou and Sevdali, 2020), where the genitive has been extended across various functions, while the accusative serves as the direct object. Although these contracted forms are formally identical to those formed by the combination of the adposition *se* and the definite article, they are structurally distinct (see Figure 1).

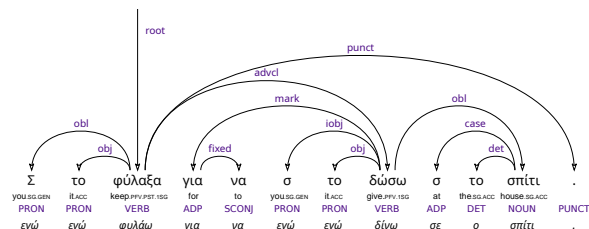


Figure 1: Dependency tree illustrating GUD’s analysis of the Greek sentence ‘I kept it for you so that I could give it to you at home’. The contractions (στο *sto*) differ in their underlying structure: the first two instances represent contracted pronoun forms (σου το *sou to*), combining a genitive pronoun and a direct object pronoun, while the last instance represents a contraction of the adposition (σε *se*) with a definite article (το *to*). For further details, see the main text.

## 4 Verbal MWE annotation

GUD contains material from fiction texts and additional 723 sentences (28% of the total GUD sentences) featuring 100 VMWEs, primarily of the verbal idiom type, along with some light verb constructions and verb-particle combinations. This classification of VMWEs follows the PARSEME typology (Savary et al., 2018). The sentences exemplify flexible usages of VMWEs in terms of morphology, word order permutations of the lexicalized parts of the VMWEs, insertion, and lexical variant pairs.

Our VMWE annotation strategy is in the spirit of the guidelines proposed by Savary et al. (2023) and incorporates suggestions from D. Zeman. In our setup, VMWE annotations are initially integrated into the MISC column (10th column) of the standard CoNLL-U format. Although the PARSEME

project utilizes an additional 11th column in alignment with CoNLL-U Plus specifications,<sup>3</sup> widely-used annotation tools such as ArboratorGrew and parsing frameworks like Stanza currently support only the standard ten-column CoNLL-U format. Notably, the DEPS (9th) and MISC (10th) columns are generally excluded from dependency parsing training procedures (Qi et al., 2020).

To effectively bridge this gap, we created a pre-processing script, `move_mwes.py`, which transfers VMWE annotations from the MISC column to the DEPREL (8th) column by appending them as sub-relation labels to the existing syntactic dependency labels. This transformation allows models such as Stanza to directly predict VMWE subrelations as part of their syntactic parsing task.

The VMWE annotation pipeline is clearly illustrated in Examples 1–3 (see Appendix):

- Example 1 shows the original treebank annotation prior to applying `move_mwes.py`.
- Example 2 illustrates the annotation after applying `move_mwes.py`, with VMWE annotations integrated into the DEPREL column.
- Example 3 presents a correct model prediction from Stanza, precisely identifying tokens that constitute an MWE.

Additionally, the `move_mwes.py` script supports reversing the annotation transformation, enabling the removal of VMWE subrelation labels from the DEPREL column and their restoration back into the MISC column. For systematic evaluation of VMWE predictions, we employ the `evaluate_mwes.py` script, which computes performance metrics detailed in Section 5. In this evaluation procedure, the annotations in the MISC column (10th column) serve as the gold standard reference against which the model predictions (encoded in the 8th column) are compared. Our integrated pipeline—[`move_mwes.py` → Model Prediction → `move_mwes.py` (reversal) → `evaluate_mwes.py`]—facilitates the efficient integration of VMWE predictions into active annotation workflows, thereby promoting continuous improvement in annotation accuracy and model performance.<sup>4</sup>

<sup>3</sup><https://universaldependencies.org/ext-format.html>

<sup>4</sup><https://github.com/JohnKaz/mwes>

## 5 Experiments and evaluation

We trained Stanza models in four experimental settings: three models combined the full GUD corpus (1,807 sentences) with additional subsets of 723, 500, and 300 sentences from IDION, each featuring  $\geq 1$  VMWE; the fourth model was trained exclusively on the 723 IDION sentences, each featuring  $\geq 1$  VMWE, without including the original GUD corpus. (Embeddings: GUD+GreekBert). We used only one test set, consisting of 200 sentences, each featuring  $\geq 1$  VMWE, with a total of 242 VMWE occurrences. Importantly, while many test VMWEs were not identical to those in the training set, a large portion were lexical variants of seen VMWEs. To ensure a diverse test set, sentences were selected to include different morphological forms of the head verb, as well as variations in word order and lexicalized component distance.

Table 1 presents the models’ evaluation results, obtained using standard UD metrics.<sup>5</sup> These metrics assess general syntactic parsing performance. The observed differences between the original GUD and the expanded GUD+723 dataset suggest potential variability. A possible cause of this variability is the difference in annotation quality between GUD and the VMWE material, or the increased structural complexity introduced by the additional VMWE-rich sentences.

Since VMWEs are encoded as subrelations within the syntactic structure, their correct identification depends on the model’s ability to accurately recover syntactic dependencies (as reflected in the UAS and LAS measures).

| Setting <sup>†</sup> | Lemma | UPOS  | UFEATS | UAS   | LAS   |
|----------------------|-------|-------|--------|-------|-------|
| GUD+723              | 90.99 | 94.78 | 87.18  | 88.03 | 81.62 |
| GUD+500              | 90.99 | 94.97 | 87.80  | 87.94 | 82.27 |
| GUD+300              | 90.23 | 94.69 | 86.93  | 88.03 | 81.25 |
| 723                  | 90.12 | 94.01 | 86.39  | 86.67 | 78.94 |

Table 1: Performance metrics for four settings. <sup>†</sup>723/500/300 sentences each one featuring at least one VMWE.

We provide a targeted evaluation of VMWE identification. In CoNLL-U, a sentence is represented as a table with 10 columns and a set of rows numbered from 1 to  $m$ ,  $m > 1$ . The representation of a VMWE with  $l, l \leq m$  lexicalized components in column 8 consists of a set of not necessar-

<sup>5</sup><https://github.com/UniversalDependencies/tools/blob/master/eval.py>

ily contiguous table cells containing information about the VMWE (a sentence may contain  $\geq 1$  VMWE):  $VMWE_x^{C8} = \{r_{i_1}^{C8}, r_{i_2}^{C8}, \dots, r_l^{C8}\}$  and in column 10:  $VMWE_x^{C10} = \{r_{i_1}^{C10}, r_{i_2}^{C10}, \dots, r_l^{C10}\}$ , where in both cases  $i_1 < i_2 < \dots < l \wedge i_n, l, n \in \{1, 2, 3, \dots, m\}$ . These simplified definitions allow to evaluate the model’s ability to discover/identify a VMWE but not its ability to classify it by type (e.g., idiom, light-verb construction, etc.).

We measure recall ( $R=TP/(TP+FN)$ ) and precision ( $P=TP/(TP+FP)$ ) of the model trained in four settings in two ways (see Table 2):

1. **Per-token.** Taking advantage of the tabular format of CoNLL-U, we use the following definitions (see also Savary et al. 2018,38):

TP if  $r_i^{C8} \in VMWE_x^{C8} \wedge r_i^{C10} \in VMWE_x^{C10}$

FP if  $r_i^{C8} \in VMWE_x^{C8} \wedge r_i^{C10} \notin VMWE_x^{C10}$

FN if  $r_i^{C8} \notin VMWE_x^{C8} \wedge r_i^{C10} \in VMWE_x^{C10}$  (for all cases above  $1 \leq i \leq l \leq m$ ).

2. **Per-unit.** A per-VMWE TP occurs if for all  $r_i^{C10} \in VMWE_x^{C10}$  there is a Per-token TP. A Per-VMWE FN occurs when there is at least one  $r_i^{C10} \in VMWE_x^{C10}$  that has a Per-token FN. Per-VMWE FP cannot be defined because we can only identify VMWEs represented in column 10.

| Setting | PTR   | PTP   | PUR   |
|---------|-------|-------|-------|
| GUD+723 | 0.813 | 0.867 | 0.606 |
| GUD+500 | 0.807 | 0.847 | 0.655 |
| GUD+300 | 0.791 | 0.850 | 0.588 |
| 723     | 0.827 | 0.880 | 0.624 |

Table 2: Performance evaluation metrics, including *per-token recall* (PTR), *per-token precision* (PTP), and *per-unit recall* (PUR) for four settings.

It should be noted that our models recognize both contiguous and non-contiguous VMWEs.

The performance analysis of our models, presented in Table 2, reveals interesting patterns regarding the effectiveness of different training configurations. The highest per-token precision (0.88) and recall (0.827) were observed in the 723-only training setting, suggesting that models trained exclusively on VMWE-rich data perform better at accurately identifying multiword expressions. However, the best per-unit recall (0.655) was achieved in the GUD+500 setting, indicating that larger training corpora can improve complete MWE identification, despite minor trade-offs in precision.

The GUD+300 setting consistently underper-

formed, with the lowest per-token recall (0.791) and per-unit recall (0.588), reinforcing the importance of sufficient VMWE-specific training data. Interestingly, while GUD+723 and 723-only performed similarly in precision and recall, the latter showed a slight advantage in correctly predicting token-level VMWE components. Future work should explore larger, more diverse datasets and fine-tune MWE subrelations to further enhance identification accuracy.

## 6 Future plans

We intend to expand our experiments by using larger test sets and corpora that encompass a wider variety of MWE types. Another direction for future research and experimentation is exploring the dissociation of MWE subrelations from syntactic annotation, potentially by encoding them in the (currently empty) XPOS column. Additionally, we aim to develop more informative evaluation metrics to better assess system performance.

The GUD treebank remains a valuable linguistic resource for facilitating knowledge transfer across Greek dialects, contributing to an ongoing contrastive study of low-resource language varieties. Furthermore, the integration of MWE into the treebank could prove beneficial for various downstream applications that rely heavily on idiomatic expressions, such as offensive language detection.

## 7 Limitations

A key limitation of our approach is that the indexes encoded in the MISC column are not interpretable by the model, as they indicate VMWE units rather than POS tags, morphological features, or dependency relations. This results in at least two major consequences:

1. The model cannot distinguish between nested VMWEs, such as those shown in the manually annotated Appendix/Example 1. Additionally, the model itself does not generate VMWE indexes. We are working on a solution to this issue.
2. To integrate VMWE annotation into the active annotation cycle, we have developed a script that transfers VMWE annotations from subrelations in the dependency relations column (8th column) to the MISC column. However, since the model does not generate indexed

VMWE annotations, the resulting MISC column lacks indexes. Consequently, manual annotation is required during the active annotation phase for newly parsed data.

Moreover, per-token evaluation is not entirely informative, as it does not indicate how many lexicalized elements of a VMWE unit are correctly recognized. We are currently exploring evaluation methods that better capture these nuances.

The test set included both seen and partially unseen VMWEs. The unseen instances shared only their fixed components with the seen ones but contained different verbal elements. In other words, the test set included lexical variants of the seen VMWEs. Ideally, our evaluation methods should differentiate between identification and discovery performance; however, this distinction is not currently made. We plan to address this issue in future work.

## Acknowledgements

This work was partially supported by project MIS 5154714 of the National Recovery and Resilience Plan Greece 2.0 funded by the European Union under the NextGenerationEU Program. It also received support from the CA21167 COST action UniDive, funded by COST (European Cooperation in Science and Technology).

## References

- Elena Anagnostopoulou and Christina Sevdali. 2020. [Two modes of dative and genitive case assignment: Evidence from two stages of greek](#). *Natural Language & Linguistic Theory*, 38(4):987–1051.
- Mathieu Constant, Gülşen Eryiğit, Johanna Monti, Lonneke van der Plas, Carlos Ramisch, Michael Rosner, and Amalia Todirascu. 2017. [Multiword Expression Processing: A Survey](#). *Computational Linguistics*, 43(4):837–892. [\\_eprint: https://direct.mit.edu/coli/article-pdf/43/4/837/1808392/coli\\_a\\_00302.pdf](https://direct.mit.edu/coli/article-pdf/43/4/837/1808392/coli_a_00302.pdf).
- Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. [Universal Dependencies](#). *Computational Linguistics*, 47(2):255–308.
- Vasiliki Foufi, Luka Nerima, and Eric Wehrli. 2019. [Multilingual parsing and mwe detection](#). *Representation and parsing of multiword expressions: Current trends*, pages 217–237.
- Voula Giouli, Vasiliki Foufi, and Aggeliki Fotopoulou. 2019. [Annotating greek vmwes in running text: A](#)

- piece of cake or looking for a needle in a haystack? In *Proceedings of the 13th International Conference on Greek Linguistics*. The University of Westminster, St. John's College, University of Cambridge, The H. M. Chadwick Fund, Cambridge University Press.
- Brian D. Joseph. 1981. *On the Synchrony and Diachrony of Modern Greek NA*. *Byzantine and Modern Greek Studies*, 7:139–154.
- Brian D. Joseph and Irene Philippaki-Warbuton. 1987. *Modern Greek*. Routledge Kegan Paul, Oxfordshire, UK.
- Panagiota Kyriazi and Aggeliki Fotopoulou. 2024. *Multiword expressions of greek language: A case study of non-referential clitics in mwes*. In Vojkan Stojičić, Ana Elaković-Nenadović, and Martha Lampropoulou, editors, *Proceedings of the 15th International Conference on Greek Linguistics*. Vol. 2, volume 15 of *International Conference on Greek Linguistics*, chapter 17, pages 292–309. University of Belgrade – Faculty of Philology, Belgrade, Serbia. Available online at [http://doi.fil.bg.ac.rs/volume.php?pt=eb\\_ser&issue=icgl-2024-15-2&i=17](http://doi.fil.bg.ac.rs/volume.php?pt=eb_ser&issue=icgl-2024-15-2&i=17).
- Stella Markantonatou, Panagiotis Minos, George Zakis, Vassiliki Moutzouri, and Maria Chantou. 2019. *IDION: A database for Modern Greek multiword expressions*. In *Proceedings of the Joint Workshop on Multiword Expressions and WordNet (MWE-WN 2019)*, pages 130–134, Florence, Italy. Association for Computational Linguistics.
- Athina Michou and Violeta Seretan. 2009. *A tool for multi-word expression extraction in Modern Greek using syntactic parsing*. In *Proceedings of the Demonstrations Session at EACL 2009*, pages 45–48, Athens, Greece. Association for Computational Linguistics.
- Prokopis Prokopidis and Haris Papageorgiou. 2017. *Universal Dependencies for Greek*. In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 102–106, Gothenburg, Sweden. Association for Computational Linguistics.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. *Stanza: A python natural language processing toolkit for many human languages*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108, Online. Association for Computational Linguistics.
- Niki Samaridi and Stella Markantonatou. 2014. *Parsing Modern Greek verb MWEs with LFG/XLE grammars*. In *Proceedings of the 10th Workshop on Multiword Expressions (MWE)*, pages 33–37, Gothenburg, Sweden. Association for Computational Linguistics.
- Agata Savary, Marie Candito, Verginica Barbu Mititelu, Eduard Bejček, Fabienne Cap, Slavomír Čěplö, Silvio Ricardo Cordeiro, Gülşen Eryiğit, Voula Giouli, Maarten Van Gompel, Yaakov HaCohen-Kerner, Jolanta Kovalevskaitė, Simon Krek, Chaya Liebeskind, Johanna Monti, Carla Parra Escartín, Lonke Van Der Plas, Behrang QasemiZadeh, Carlos Ramisch, Federico Sangati, Ivelina Stoyanova, and Veronika Vincze. 2018. *PARSEME multilingual corpus of verbal multiword expressions*. In Stella Markantonatou, Carlos Ramisch, Agata Savary, and Veronika Vincze, editors, *Multiword expressions at length and in depth: Extended papers from the MWE 2017 workshop*, pages 87–147. Language Science Press, Berlin. 10.5281/ZENODO.1471591.
- Agata Savary, Sara Stymne, Verginica Barbu Mititelu, Nathan Schneider, Carlos Ramisch, and Joakim Nivre. 2023. *Parseme meets universal dependencies: Getting on the same page in representing multiword expressions*. *The Northern European Journal of Language Technology (NEJLT)*, 9(1).
- Vivian Stamou, Marilena Malli, Penny Takorou, Artemis Xylogianni, and Stella Markantonatou. 2020a. *Evaluation of Verb Multiword Expressions discovery measurements in literature corpora of Modern Greek*. In *Lexicography for Inclusion: Proceedings of the 19th EURALEX International Congress, 7-9 September 2021, Alexandroupolis, Vol. 1*, pages 295–301, Alexandroupolis. Democritus University of Thrace.
- Vivian Stamou, Artemis Xylogianni, Marilena Malli, Penny Takorou, and Stella Markantonatou. 2020b. *VMWE discovery: a comparative analysis between literature and Twitter corpora*. In *Proceedings of the Joint Workshop on Multiword Expressions and Electronic Lexicons*, pages 66–72, online. Association for Computational Linguistics.

## A Appendix

|   |        |             |         |
|---|--------|-------------|---------|
| που   | μας    | έχεις       | αφήσει  |
| rou   | mas    | echis       | afisi   |
| because                                       | us     | have.2.SING | left    |
| σύξυλους                                      | στους  | πέντε       | δρόμους |
| sixilous                                      | stous  | pentē       | dromous |
| petrified                                     | in.the | five        | roads   |
| 'because you have astounded and abandoned us' |        |             |         |

```
# text = που μας έχεις αφήσει σύξυλους στους πέντε δρόμους.
18 που που SCONJ __ 21 mark __
19 μας εγώ PRON _ Case=AccINumber=PlurIPerson=1IPronType=Prs
21 obj __
20 έχεις έχω AUX _ Mood=IndINumber=SingIPerson=2ITense=Pres
|VerbForm=Fin|Voice=Act 21 aux __
21 αφήσει αφήνω VERB _ Aspect=PerfIMood=IndI|VerbForm=Inf|Voice=Act 2 advcl _ mwe=1,2:VID
22 σύξυλους σύξυλος ADJ _ Case=AccI|Gender=MascI|Number=Plur
21 xcomp _ mwe=1
23 στους στου ADP _ Case=AccI|Gender=MascI|Number=Plur 25 case
_ mwe=2
24 πέντε πέντε NUM _ Case=AccI|Gender=MascI|Number=Plur
|NumType=Card 25 nummod _ mwe=2
25 δρόμους δρόμος NOUN _ Case=AccI|Gender=MascI|Number=Plur
21 obl _ mwe=2:VID
26 . . PUNCT __ 2 punct _ PunctType=Peri
```

**Example 1: Annotation of 2 conflated VMWEs with the same verb head (afisei) and different lexicalized parts (sixilous, pente dromous).**

|     |       |             |             |     |       |
|-----|-------|-------------|-------------|-----|-------|
| Δεν | βγήκε | ποτέ        | από         | το  | μυαλό |
| den | vgike | pote        | apo         | to  | mialo |
| it  | never | left.3.SING | from        | the | mind  |
| μου | και   | ούτε        | πρόκειται   |     |       |
| mou | kai   | oute        | prokeitai   |     |       |
| my  | and   | neither     | will.happen |     |       |

'It never left my mind and it will not'

```
# text = Δεν βγήκε ποτέ από το μυαλό μου και ούτε πρόκειται.
1 Δεν den PART PtNg Polarity=Neg 2 advmod _ _
2 βγήκε βγαίνω VERB _ _ As-
pect=Perf|Mood=Ind|Number=Sing|Person=3|Tense=Past|VerbForm=Fin|Voice=Act
0 root:vid _ mwe=1:VID
3 ποτέ ποτέ ADV _ _ 2 advmod _ None=Yes
4 από από ADP _ _ 6 case:vid _ mwe=1|None=Yes
5 το ο DET _ Case=Acc|Definite=Def|Gender=Neut|Number=Sing|PronType=Art
6 det:vid _ mwe=1
6 μυαλό μυαλό NOUN _ Case=Acc|Gender=Neut|Number=Sing 2
obl:vid _ mwe=1
7 μου εγώ PRON _ Case=Gen|Number=Sing|Person=1|Poss=Yes|PronType=Prs
6 nmod _ _
8 και και CCONJ _ _ 10 cc _ None=Yes
9 ούτε ούτε PART _ Polarity=Neg 10 advmod _ None=Yes
10 πρόκειται πρόκειται VERB _ _ As-
pect=Impl|Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin|Voice=Pass
2 conj _ _
11 . . PUNCT _ _ 2 punct _ PunctType=Peri
```

**Example 2: Output produced by the script.**

|      |        |      |           |
|------|--------|------|-----------|
| ο    | πήχυς  | για  | φέτος     |
| o    | pichis | gia  | fetos     |
| the  | bar    | for  | this.year |
| έχει | ανέβει | πολύ | ψηλά      |
| echi | anevi  | poli | psila     |
| has  | risen  | very | high      |

'The bar for this year has risen very high'

```
# text = Ο πήχυς για φέτος έχει ανέβει πολύ ψηλά
1 Ο ο DET _ Case=Nom|Definite=Def|Gender=Masc|Number=Sing
|PronType=Art 2 det:vid _ _
2 πήχυς πήχης NOUN _
Case=Nom|Gender=Masc|Number=Sing 6 obj:vid _ _
3 για για ADP _ _ 6 case _ None=Yes
4 φέτος φέτος ADV _ _ 3 fixed _ _
5 έχει έχω AUX _ _
Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin
6 aux _ _
6 ανέβει ανέβω VERB _ _ As-
pect=Perf|Mood=Ind|Number=Sing|Person=3|VerbForm=Fin|Voice=Act
0 root:vid _ _
7 πολύ πολύ ADV _ _ 8 advmod _ None=Yes
8 ψηλά ψηλά ADV _ _ 6 advmod _ _
```

**Example 3: Output produced by the Stanza model.**