

# Building a Contact Language Treebank with UniDive Support: The Asia Minor Greek in Contact (AMGiC) UD Resource

Anonymous ACL submission

001 *Relevant UniDive working groups:* WG1, WG4

## 002 1 Introduction

003 The UniDive COST Action (CA21167) has pro- 040  
004 vided a framework for advancing Universal De- 041  
005 pendencies (UD) resources for typologically di- 042  
006 verse and under-documented languages. In this ab- 043  
007 stract, we report on how UniDive activities—Short- 044  
008 Term Scientific Missions (STSMs), working group 045  
009 discussions, and training schools—have directly 046  
010 contributed to the development of the *Asia Minor*  
011 *Greek in Contact* (AMGiC) treebank, a UD re-  
012 source for annotating contact-induced morphosyn-  
013 tactic phenomena (CIMSP) in Asia Minor Greek  
014 (AMG) varieties.

015 AMG designates a group of Greek dialects spo- 047  
016 ken in Central Anatolia until 1923, when the Treaty 048  
017 of Lausanne mandated a population exchange be- 049  
018 tween Greece and Turkey. These varieties—Silliot, 050  
019 Cappadocian, and Phrasiot Greek—underwent 051  
020 intensive contact with Turkish, resulting in signifi- 052  
021 cant typological restructuring: borrowing of free 053  
022 grammatical morphemes, agglutinative-like noun 054  
023 inflection, copula encliticization, and adoption of 055  
024 head-final syntactic configurations (Janse, 2009; 056  
025 Thomason and Kaufman, 1988). The AMGiC tree- 057  
026 bank was initiated by Sampanis and Prokopidis 058  
027 (2021) and has since expanded through UniDive- 059  
028 supported collaboration. 060

## 029 2 UniDive Contributions to AMGiC

030 **STSMs.** Two UniDive STSMs were central to 063  
031 advancing the treebank. In June 2024, an STSM 064  
032 at the University of Vienna (titled “Extending the 065  
033 AMGiC UD treebank”) focused on developing an- 066  
034 notation infrastructure: the BoAT (Boğaziçi An- 067  
035 notation Tool) web application was extended with 068  
036 metadata editing and sentence management func- 069  
037 tionalities tailored to AMGiC’s needs, and the tool 070  
038 was introduced to the treebank’s annotator, who 071  
039 had previously relied on a text editor. The tool 072

040 and the treebank were also presented in the “His- 041  
042 torical Grammar of Modern Greek” course at the 043  
044 University of Vienna. In September 2024, a re- 045  
046 ciprocal STSM at Boğaziçi University enabled in- 047  
048 tensive collaborative annotation work, resulting in 049  
050 the completion of the first 36 Silliot sentences and 051  
052 their preparation for publication. 053

054 **Working groups.** AMGiC is relevant to WG1 047  
048 (Corpus Annotation), as it develops novel annota- 049  
050 tion conventions for contact phenomena within the 051  
052 UD/CoNLL-U framework, and to WG4 (Quantify- 053  
054 ing and Promoting Diversity), as it contributes a 054

055 **Training schools.** One of the authors partici- 055  
056 pated in both UniDive Training Schools: the 1st 056  
057 (Chişinău, July 2024), which covered dependency 057  
058 syntax, UD annotation, treebank validation and 058  
059 quality control; and the 2nd (Yerevan, January 059  
060 2026), which focused on low-resource and typo- 060  
061 logically diverse languages in NLP. 061

## 062 3 The AMGiC Treebank

063 The initial 36 Silliot sentences were completed 063  
064 as a direct outcome of the two STSMs. The 064  
065 treebank has since been expanded to **72 anno-** 065  
066 **tated sentences** (851 tokens) covering two dialect 066  
067 groups: Silliot (36 sentences; Kostakis, 1968; R. M. 067  
068 Dawkins, 1916) and Cappadocian (36 sentences, 068  
069 Delmesó subdialect; R. M. Dawkins, 1916). The 069  
070 treebank follows the UD annotation framework 070  
071 (Nivre et al., 2020; De Marneffe et al., 2021) and 071  
072 introduces two distinctive features: 072

073 **CIMSP annotation.** A systematic taxonomy 073  
074 of contact-induced phenomena is encoded in 074  
075 the CoNLL-U MISC field using three features: 075  
076 LC=Yes, MorphSynC (broad category), and 076  
077 MorphSynSC (specific subcategory). Table 1 077  
078 shows the distribution across categories. The most 078

Tag	Category	<i>n</i>
FrGrEl	Free Grammatical Elements	37
SynIn	Syntactic Interference	7
MorphIn	Morphological Interference	5
BGrEl	Bound Grammatical Elements	3
L	Lexical	2
ELIs	Embedded Language Islands	1
<b>Total</b>		<b>55</b>

Table 1: CIMSP categories in AMGIC (72 sentences).

frequent category is Free Grammatical Elements (FrGrEl), covering borrowed question particles, subordinating conjunctions, and quantifiers from Turkish.

**Sociodemographic metadata.** Each sentence carries metadata encoding extralinguistic parameters relevant to contact intensity—population coexistence patterns, proximity to urban centres, availability of Greek education, enclave status, and diaspora connections—informed by the framework of Karantzola et al. (2021). This enables correlations between sociolinguistic conditions and CIMSP frequency as the treebank grows.

#### 4 Outcomes and Future Work

The UniDive-supported development of AMGIC has produced the following outcomes: (a) the completion of the initial 36-sentence Silliot portion of the treebank through two STSMs, with subsequent expansion to 72 sentences covering a second dialect group; (b) a replicable CIMSP taxonomy for annotating contact phenomena within the UD framework, applicable to other contact situations; (c) the BoAT collaborative annotation tool, extended with metadata editing capabilities specifically for AMGIC during the Vienna STSM; and (d) sociodemographic metadata enabling quantitative sociolinguistic analysis of contact intensity.

Future directions include expanding the treebank to cover additional Cappadocian subdialects and Pharsiot Greek, and conducting statistical analyses of CIMSP distribution across dialects and sociodemographic conditions. The treebank will be made publicly available through UD (UD\_Cappadocian-AMGIC) under a Creative Commons licence.

#### References

Marie-Catherine De Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. [Universal](#)

<a href="#">Dependencies</a> . <i>Computational Linguistics</i> , pages 1–54.	117 118
Mark Janse. 2009. Greek-Turkish Language Contact in Asia Minor. <i>Études Helléniques/Hellenic Studies</i> , 17:37–54.	119 120 121
Eleni Karantzola, Anatoli Theodoridi, and Konstantinos Sampanis. 2021. <a href="#">The Interplay of External and Sociolinguistic Factors in Contact-Induced Language Change: Cappadocian Greek as a Case Study</a> . <i>Mediterranean Language Review</i> , 28:21.	122 123 124 125 126
Athanasios P. Kostakis. 1968. <i>To glossiko idioma tis Sillis</i> . Centre for Asia Minor Studies, Athens. The dialect of Silli.	127 128 129
Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. <a href="#">Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection</a> . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 4034–4043, Marseille, France. European Language Resources Association.	130 131 132 133 134 135 136 137 138
R. M. Dawkins. 1916. <i>Modern Greek in Asia Minor</i> . Cambridge University Press.	139 140
Konstantinos Sampanis and Prokopis Prokopidis. 2021. <a href="#">Asia Minor Greek in Contact (AMGIC): Towards a dialectal Treebank comprising contact-induced grammatical changes</a> . In <i>Proceedings of the 20th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2021)</i> , pages 86–95, Sofia, Bulgaria. Association for Computational Linguistics.	141 142 143 144 145 146 147
Sarah Grey Thomason and Terrence Kaufman. 1988. <i>Language Contact, Creolization, and Genetic Linguistics</i> , first paperback printing 1991, reprint 2020 edition. University of California Press, Berkeley, CA.	148 149 150 151 152