

HUN-REN Hungarian Research Centre for Linguistics,
Budapest, Hungary, 29-30 January 2025
https://unidive.lisn.upsaclay.fr/







GUD: a new Standard Modern Greek treebank enriched with VMWE annotations

Stella Markantonatou^{1,2},Vivian Stamou², Stavros Bompolas², Katerina Anastasopoulou³, Iriana Vasileiadi-Linardaki^{3,} Konstantinos Diamantopoulos³, Yannis Kazos^{4,2}

¹[Institute for Language and Speech Processing, ²Archimedes]/Athena RC, Greece

³National and Kapodistrian University of Athens (NKUA), Greece

⁴National Technical University of Athens (NTUA), Greece

Introduction

- UD_Greek-GUD: Universal Dependencies (UD).v2 treebank of Standard Modern Greek (SMG)
- The first UD treebank for SMG with annotations for Verbal Multiword Expressions (VMWEs)

Key Features

As compared to previous UD treebanks of SMG, GUD:

- Adheres closer to the UD.v2 morphological and syntactic guidelines
- Annotates more phenomena, including diminutives and augmentatives, adjective degrees and **Verbal MWEs**

Experiments

- Data:
- GUD: 1.807 sentences from fiction texts.
- √ 723 sentences with 100 VMWEs from IDION VMWE database (https://search.idion.athenarc.gr)
- Stanza models trained with GUD + 723/500/300 VMWE sentences and 723 VMWE sentences

Results

Setting	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

Annotation and Processing of Verbal MWEs

Annotation of VMWEs on column 10 (C10) of CONLLU*

text = « Δεν θα βάζω για πάντα πλάτη » (I won't put forever my.back/"I won't keep backing this up forever.")

# text	= « ∆٤V	θα βάζω	για πάντα	α πλάτη :	» (I won't put fo	orever m	ny.back/"I	I won't keep backing this up forever.")
# sent	_id = out	put991						
1	«	«	PUNCT	OPUNCT	_ 4	punct	_	PunctType=Quot
2	Δεν	δεν	PART	OPUNCT	Polarity=Neg	4	advmod	
3	θα	θα	AUX	PtFu	Tense=Fut	4	aux	
4	βάζω	βάζω	VERB	VbMn	Aspect=Imp Mood=	=Ind Num	nber=Sing	Person=1 Tense=Pres VerbForm=Fin Voice=Act
5	για	για	ADP	AsPpSp	_ 6	case	_	None=Yes
6	πάντα	πάντα	DET	AdBa	_ 4	obl	_	None=Yes
7	πλάτη	πλάτη	NOUN	NoCm	Case=Acc Gender=	=Fem Num	nber=Sing	4 obj _ mwe=1
8	>>	>>	PUNCT	NoCm	_ 4	punct	_	PunctType=Quot

Transferring VMWE annotation to the Deprel CONLLU column (C8)

# sent_	id = out	:put991	L					
1	((«	PUNCT	OPUNCT	_ 4	punct	_	PunctType=Quot
2	Δεν	δεν	PART	OPUNCT	Polarity=Neg	4	advmod	
3	θα	θα	AUX	PtFu	Tense=Fut	4	aux	
4	βάζω	βάζω	VERB	VbMn	Aspect=Imp Mood	d=Ind Nun	mber=Sing	Person=1 Tense=Pres VerbForm=Fin Voice=Act 0 root:vid _ mwe=1:VID
5	για	για	ADP	AsPpSp	_ 6	case	_	None=Yes
6	πάντα	πάντα	DET	AdBa	_ 4	obl	_	None=Yes
7	πλάτη	πλάτη	NOUN	NoCm	Case=Acc Gender	=Fem Nun	mber=Sing	4 obj:vid _ mwe=1
8	>>	>>	PUNCT	NoCm	4	punct		PunctType=Quot

Evaluation \rightarrow Definitions:

- TP when a correct annotation occurs on C8 and on C10
- FP when a correct annotation occurs on C8 but not C10
- occur on C8 but occurs on C10

 Recall (R=TP/(TP+FN) and Precision

FN when a correct annotation does not

- Recall (R=TP/(TP+FN) and Precision (P=TP/TP+FP) are measured:
- 1. Per-token
- 2. Per-VMWE
- *Per-VMWE precision: undefinable

	Per-token	Per-VMWE
Precision	0.84	Undefinable
Recall	0.94	0.53

Application of evaluation metrics

# text	= 0 πήχυ	ις για φε	έτος έχει	ι ανέβει	ι πολύ ψηλά (The bar for this.year has risen very high/"The bar has risen very high this year.")	
# sent	_id = out	:put610	2 6			
1	0	0	DET	AtDf	Case=Nom Definite=Def Gender=Masc Number=Sing PronType=Art 2 det:vid _ mwe=1	
2	πήχυς	πήχης	NOUN	NoCm	Case=Nom Gender=Masc Number=Sing 6 obj:vid _ mwe=1	
3	για	για	ADP	AsPpSp	p _ 6 mark _ None=Yes	
4	φέτος	φέτος	ADV	AdBa	_ 3 fixed _	
5	έχει	έχω	AUX	VbMn	Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin Voice=Act 6 aux _	
6	ανέβει	ανέβω	VERB	VbMn	Aspect=Perf Mood=Ind Number=Sing Person=3 VerbForm=Fin Voice=Act 0 root _ mwe=1:	VID
7	πολύ	πολύ	ADV	AdBa	_ 8 advmod _ None=Yes	
8	ψηλά	ψηλά	ADV	AdBa	_ 6 advmod _	

Applications

- Cross-dialectal knowledge transfer for Modern Greek dialects
- Enhancement of downstream NLP tasks,
- e.g., offensive language detection

Per-token recall: 2(TP)/(2(TP)+1(FN)) = 2/3 = 0.66

Per-token/Per-VMWE Recall = 1, Per-token Precision = 1

# text =	= Δεν βγ	ήκε ποτέ	από το	μυαλό μο	υ και ού	τε πρόκε	ιται . (It never	got.3.S	ING out	of my i	mind and	it will n	not/"It	never le	eft my min	d and it	t never v	vill.")
# sent_i	id = out	put61									-					-			·
1	Δεν	δεν	PART	PtNg	Polarit	y=Neg	2	advmod	_	_									
2	βγήκε	βγαίνω	VERB	VbMn	Aspect=	Perf Moo	d=Ind Nu	mber=Sin	g Person	=3 Tens	se=Past	VerbForm:	=Fin Voice	e=Act	0	root:vi	d	_	mwe=1:VID
3	ποτέ	ποτέ	ADV	AdBa	PronTyp	e=Neg	2	advmod	_	None=\	/es								
4	από	από	ADP	AsPpSp	_	6	case:vi	d	_	mwe=1	None=Ye	S							
5	το	0	DET	AtDf	Case=Ac	c Defini	te=Def G	ender=Ne	ut Numbe	r=Sing	PronType	e=Art	6	det:vi	d _	mwe=1			
6	μυαλό	μυαλό	NOUN	NoCm	Case=Ac	c Gender	=Neut Nu	mber=Sin	g	2	obl:v:	id _	mwe=1						
7	μου	εγώ	PRON	PnPo	Case=Ge	n Number	=Sing Pe	rson=1 P	oss=Yes	PronTyp	e=Prs	6	nmod	_	_				
8	και	και	CCONJ	CjCo	_	10	CC	_	None=Ye	S									
9	ούτε	ούτε	CCONJ	CjCo	_	10	CC	_	None=Ye	S									
10	πρόκειτ	αι	πρόκειτ	αι	VERB	VbMn	Aspect=	Imp Mood	=Ind Num	ber=Sir	ng Perso	n=3 Tense	e=Pres Ver	rbForm=F	in Voice	e=Pass	2	conj	
11			PUNCT	PTERMP	_	2	punct	_	PunctTy	pe=Peri	Ĺ								

Future Work

- Further investigate VMWE evaluation issues
- Develop the full suite of tools to support VMWE active annotation
- Experiment with larger datasets/more MWE types

*Annotation done in the spirit of Savary et al. (2023)



EUROPEAN COOPERATION IN SCIENCE & TECHNOLOGY

