









Annotation of MWEs and NEs in the Serbian extension of **ELEXIS-WSD: comparisons, solutions and open questions** Cvetana Krstev⁽¹⁾, Ranka Stanković⁽²⁾, Aleksandra Marković⁽³⁾ 1) Jerteh; 2) UB FMG; 3) SASA ISL; Belgrade, Serbia

Next steps:

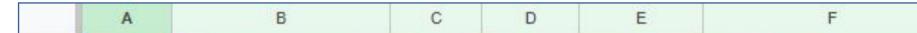
MWEs & NEs annotation

linking with the sense repository.

syntactic annotation

1. The extension of ELEXIS-WSD

ELEXIS-WSD – parallel sense-annotated corpus; N, ADJ, V, ADV with assigned senses for 10 languages: BG, DA, EN, ES, ET, HU, IT, NL, PT, SL. The set of 2024 sentences from WikiMatrix in EN was translated automatically into SR; the translation was checked by native speakers and linguists and proofread afterwards.



= tag p = tag u = lemma = komenar POS

Rece = word

	А	В	С	D	E
1	RB =	En =	Sr =	Google translation =	Corr (ovde stavljati korigovane vrednosti, da original ipak ostane) Pismo, ćirilica
1584	1583.en	The gas is heavier than air.	Течности су вискозније од гасова.	Гас је тежи од ваздуха.	Гас је тежи од ваздуха.
1585	1584.en	The strategic position of the Ardennes has made it a battleground for European powers for centuries.	1	Стратешки положај Ардена чинио их је вековима бојним пољем европских сила.	Стратешки положај који су Ардени имали чинио их је током векова поприштем европских сила.
1586	1585.en	He also established the relationship between barometric pressure and height above sea level.	1	Такође је установио везу између барометарског притиска и висине изнад нивоа мора.	Установио је и везу између барометарског притиска и надморске висине.

2. MWEs & NEs in WSD

NE **MWE** Lang. lemma sense lemma sense 2994652 bg 477440440459 da 179309 3640177 112217145 hu 6 **MWEs**: 41 4237 2733 nl lingua franca appeared in 6 language sets;

e	All MWEs & NEs	Ú
2	from the WSD	
$\frac{9}{1}$	automatically	
8	translated into sR.	
5	Number of exact ma	tches:

and manually corrected.

Set automatically tokenized, lemmatized, POS-tagged,

	Nece .		tag_p .	tag_u .	iemina_ ·	Komenar PO3
9539	1583	je	V	AUX	jesam	
2540	1583	тежи	V	VERB	težiti	
9541	1583	од	PREP	ADP	od	
7542	1583	ваздуха	N	NOUN	vazduh	
7543	1583		SENT	PUNCT		
7544	1583					
7545	1584	<s></s>				
2546	1584	Стратешки	A	ADJ	strateški	
7547	1584	положај	N	NOUN	položaj	
2548	1584	који	PRO	DET	koji	
549	1584	Ардени	N	NOUN	Ardeni	N/PROPN
2550	1584	имају	V	VERB	imati	
9551	1584	чинио	V	VERB	činiti	
9552	1584	ИХ	PRO	PRON	ona	
9553	1584	je	V	AUX	jesam	
2554	1584	током	PREP	ADP	tokom	
9555	1584	векова	N	NOUN	vek	
9556	1584	поприштем	N	NOUN	poprište	
9557	1584	европских	A	ADJ	evropski	
2558	1584	сила	N	NOUN	sila	
7559	1584		SENT	PUNCT		
560	1584					
9561	1585	<s></s>				
9562	1585	Установио	V	VERB	ustanoviti	
2563	1585	je	V	AUX	jesam	
564	1585	И	CONJ	PART	i	
7565	1585	везу	V	VERB	vesti	
2566	1585	између	PREP	ADP	između	
2567	1585	барометарског	A	ADJ	barometarsk	
2568	1585	притиска	N	NOUN	pritisak	
569	1585	И	CONJ	CCONJ	i	
570	1585	надморске	A	ADJ	nadmorski	
9571	1585	висине	N	NOUN	visina	
9572	1585		SENT	PUNCT	-	
9573	1585					

Serbian set automatically annotated using 4 resources & tools:

(1) Dictionary of non-verbal MWEs (653 occurrences): 357 NMWEs

Total	1,710	2,160	606	663
sl	385	451	0	0
pt	113	115	14	15

Table 1: Number of MWEs/NEs in the repository.

3. The comparison of **MWEs & NEs accross** languages

Automatic translation of MWEs sometimes imprecise: ustati (SR) 'get up' as incorrect translation of: издигам се (вс) nastati 'become' holde stand (DA) ------ izdizati 'elevate' stå op (DA) izlaziti 'rise' *üles kasvatama* (Ет)-----> odgajati 'raise'

 high school sR: srednja škola (lit. 'middle school'); вG: *висше училище;* EN: *high school;* SL: *srednja* šola & *visoka šola; IT: scuola media

NEs: Grčka 'Greece' was the most freequent NE, translated from: Grækenland (DA), Grecia (ES), Kreeka (ET), Grécia (PT)

№	1	2	3	4	5	6	Tot.
MWE	92	163	40	14	2	1	1412
NE	453	67	5	1	0	0	526

 Table 2: MWEs and NEs translations into Serbian
obtained by translating from 1 to 6 languages.

Correct equvialents:

4. Open questions & future work

finding the alternative ways of aligning MWEs and NEs across languages;

(134 PROPN), 73 PREP, 52 ADV, 36 CONJ, 1 ADJ.

(2) NER based on e-dictionaries and
rules (2006 occurrences).

(3) Recognition of VMWEs based on e-dictionaries, rules, and the repertoire of VMWEs annotated in the Serbian part of the PARSEME 1.3 (228 occurrences: IRV – 174, LVC.full – 35, VID - 11, and LVC.cause - 8).

Tag	Nº	Tag	N⁰
PERS	329	TIME	372
TOP	448	AMOUNT	169
ORG	126	MEASURE	62
DEMONYM	244	PERCENT	51
ROLE	175	MONEY	12
EVENT	18	Total	2,006

Table 3: Recognized NEs by type.

(4) Recognition of adjectival and verbal similes: not a single simile was retrieved.

In some cases the translation of MWEs was good and used in Serbian set, but not annotated in it because it was missing in the above resources: društvena mreža 'social network'.

The same holds for other languages: *heavy water* (EN) & $mexka \ eoda$ (BG) were not annotated neither. The equivalents in some languages are single tokens: água-pesada (PT) and *nehézvizet* (ни).

Only 93 NEs annotated in SSS were annotated as MWE or PROPN in WSD (maybe due to the poor lemmatization and linking of proper names).

should the set of sentences be enhanced to **(b** capture a more versatile style, e.g. fiction (as the lack of simile figures suggests)?

(c) should the repertoire of NE classes be unique for all languages?

(d) should NEs include numeric and/or temporal expressions?

(e) should the nesting of MWEs/NEs be allowed?

This research was supported by the Science Fund of the Republic of Serbia, #7276, Text Embeddings - Serbian Language Applications - TESLA.

