



Evaluation of taggers for 19th-century fiction

Silvie Cinková – Tomáš Erjavec –
Cláudia Freitas – Ioana Galleron –
Péter Horváth – Christian-Emil Ore –
Pavel Smrž – Balázs Indig



Distant [📄] Reading

19th-century fiction corpora

Language	Texts	Words
cze	23	692936
deu	44	8422194
eng	91	11990064
fra	72	6035694
hun	101	7563930
ita	34	3328244
nor	27	1114092
por	53	3513274
rom	5	200037
slv	72	3894267
spa	17	1981499
srp	24	886794

tokens

The cats sat on mats .

determiner

noun

verb

preposition

noun

plural

past

plural

part of speech (POS) tags

lemmas

the

cat

sit

on

mat



universaldependencies.org



	Afrikaans
	Akkadian
	Amharic
	Ancient Greek
	Arabic
	Armenian
	Assyrian
	Bambara
	Basque
	Belarusian
	Breton
	Bulgarian
	Buryat
	Cantonese
	Catalan
	Chinese
	Classical Chinese
	Coptic
	Croatian
	Czech
	Danish
	Dutch
	English
	Erzya
	Estonian
	Faroese
	Finnish
	French
	Galician
	German
	Gothic
	Greek
	Hebrew
	Hindi
	Hindi English
	Hungarian
	Indonesian
	Irish
	Italian
	Japanese
	Karelian
	Kazakh
	Komi Zyrian
	Korean
	Kurmanji
	Latin
	Latvian
	Lithuanian
	Maltese
	Marathi
	Mbya Guaraní
	Moksha
	Naija
	North Sami
	Norwegian
	Old Church Slavonic
	Old French
	Old Russian
	Persian
	Polish
	Portuguese
	Romanian
	Russian
	Sanskrit

	Serbian
	Skolt Sami
	Slovak
	Slovenian
	Spanish
	Swedish
	Swedish Sign Language
	Tagalog
	Tamil
	Telugu
	Thai
	Turkish
	Ukrainian
	Upper Sorbian
	Urdu
	Uyghur
	Vietnamese
	Warlpiri
	Welsh
	Wolof
	Yoruba

	Assamese
	Bengali
	Bhojpuri
	Cusco Quechua
	Dargwa
	Georgian
	Kannada
	Kyrgyz
	Livvi
	Macedonian
	Pnar
	Romansh
	Scottish Gaelic
	Shipibo Konibo
	Sindhi
	Somali
	Sorani
	Swiss German

form	lemma	upos	features
Her	she	PRON	Gender=Fem Number=Sing Person=3 Poss=Yes PronType=Prs
diamonds	diamond	VERB	Number=Pl
blazed	blaze	VERB	Tense=Past VerbForm=Part
out	out	ADP	_

- universal tagset + dependencies
- over 100 treebanks
- over 70 languages
- regular version releases
- since 2015
- new languages upcoming
- range of tools

universaldependencies.org

Universal POS (parts of speech)		
ADJ	adjective	<i>good</i>
ADP	adposition	<i>before, ago, in</i>
ADV	adverb	<i>fortunately</i>
AUX	auxiliary verb	<i>be, have</i>
CCONJ	coord. conjunction	<i>and, but</i>
DET	determiner	<i>this, a, both</i>
INTJ	interjection	<i>oh, thanks</i>
NOUN	noun, cf. PROPN	<i>woman</i>
NUM	numeral	<i>two, 12, VII</i>
PRON	pronoun	<i>you</i>
PART	particle	<i>indeed, just</i>
PUNCT	punctuation	<i>,;!.</i>
SCONJ	subord. conjunction	<i>because</i>
SYM	symbol	<i>,\$,%</i>
VERB	verb	<i>speak</i>

Feature	Values
Number	Plur, Sing, ...
Animacy	Anim, Inan
Gender	Masc, Fem, ...
VerbForm	Inf, Fin, Part
Definite	Def, Indef
Voice	Act, Pass,...

form	lemma	upos	features
Her	she	PRON	Gender=Fem Number=Sing Person=3 Poss=Yes PronType=Prs
diamonds	diamond	VERB	Number=Pl
blazed	blaze	VERB	Tense=Past VerbForm=Part
out	out	ADP	_

UDPipe

[About](#)[Run](#)[REST API Documentation](#)

UDPipe is a trainable pipeline for tokenization, tagging, lemmatization and dependency parsing of CoNLL-U files. UDPipe is language-agnostic and can be trained given only annotated data in CoNLL-U format. For more information, see the [UDPipe User's Manual](#).

<http://lindat.mff.cuni.cz/services/udpipe/run.php>

Copyright 2016 by Institute of Formal and Applied Linguistics, Faculty of Mathematics and Physics, Charles University in Prague, Czech Republic.

Description of the available methods is available in the [API Documentation](#) and the models are described in the [UDPipe User's Manual](#).

Service

The service is freely available for testing. Respect the [CC BY-NC-SA](#) licence of the models – explicit written permission of the authors is required for any commercial exploitation of the system. If you use the service, you agree that data obtained by us during such use can be used for further improvements of the systems at UFAL. All comments and reactions are welcome.

Model: Universal Dependencies 1.2

czech-ud-1.2-160523

Input: Plain text CoNLL-U Horizontal Vertical

Actions: Tag and Lemmatize Parse

▲ Input Text

📎 Input File

Podle agentury TASS, která se odvolala na očité svědky, neznámý ozbrojenec zasáhl velvyslance do zad.

↓ Process Input ↓

▲ Output Text

📄 Show Table

🌲 Show Trees

📄 Save Tree as SVG

Previous

1

2

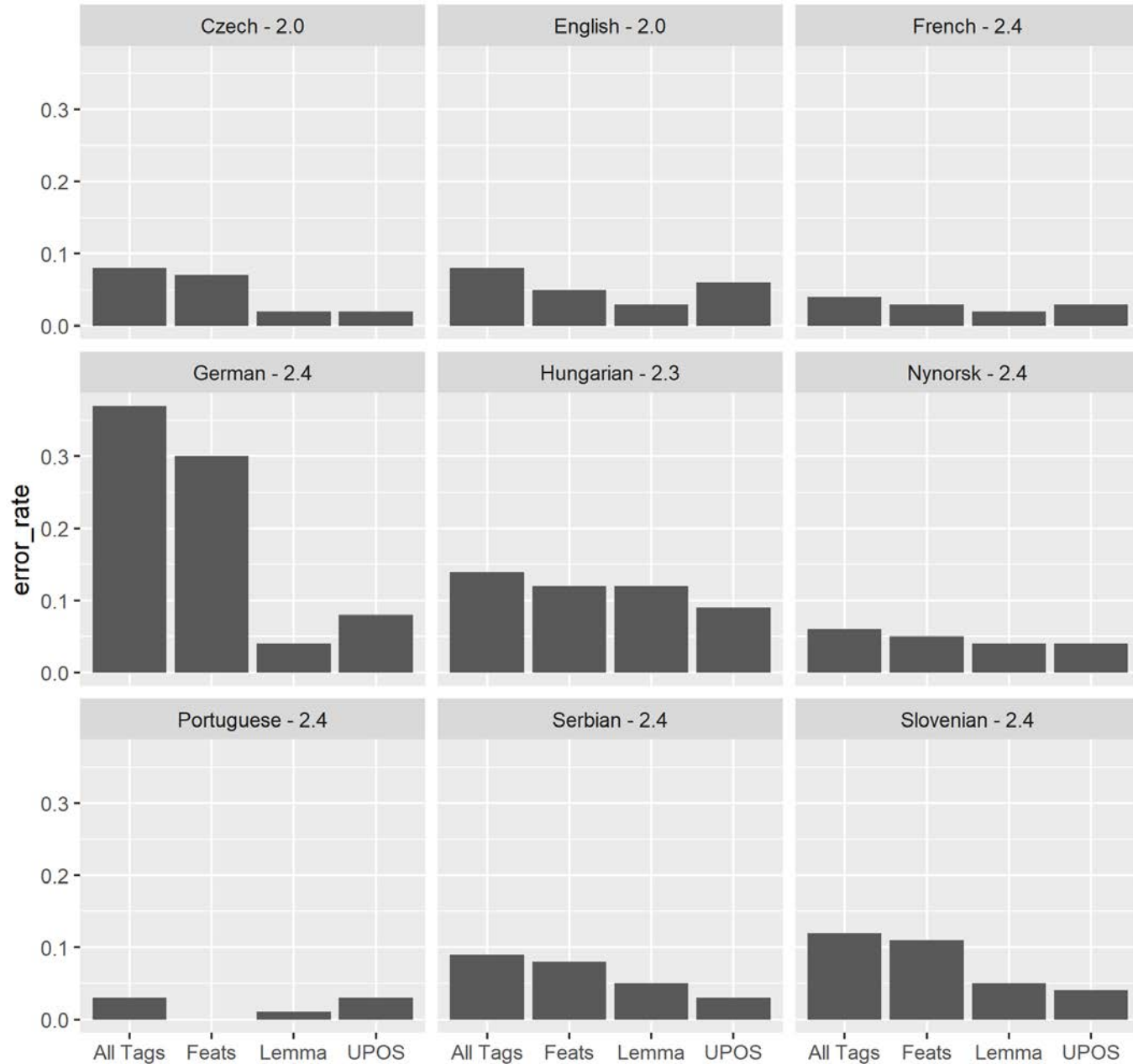
Next

Podle agentury TASS, která se odvolala na očité svědky, neznámý ozbrojenec zasáhl velvyslance do zad.

Conll-u format

Id	Form	Lemma	UPosTag	XPosTag	Feats
# newdoc					
# newpar					
# sent_id = 1					
# text = This is one of the finest of our castles.					
1	This	this	PRON	DT	Number=Sing PronType=Dem
2	is	be	AUX	VBZ	Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin
3	one	one	NUM	CD	NumType=Card
4	of	of	ADP	IN	_
5	the	the	DET	DT	Definite=Def PronType=Art
6	finest	finest	NOUN	NN	Number=Sing
7	of	of	ADP	IN	_
8	our	we	PRON	PRP\$	Number=Plur Person=1 Poss=Yes PronType=Prs
9	castles	castle	NOUN	NNS	Number=Plur
10	.	.	PUNCT	.	_

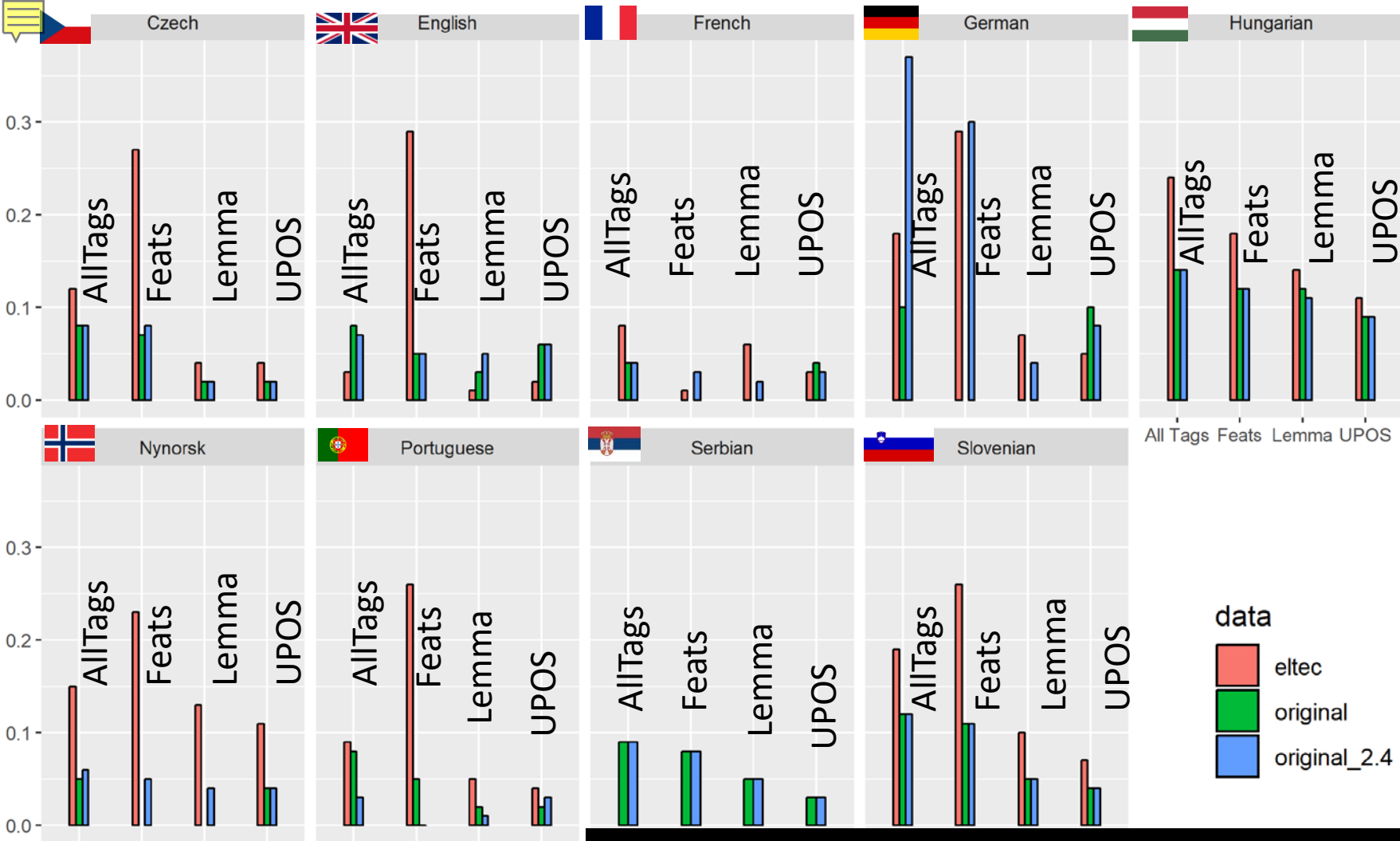
UDPipe error rate in documentation



Manual annotation

Form	Lemma	upos	xpos	feats	EV_lemma	EV_tok	EV_upos	EV_feats_errors
<i>said</i>	say	VERB	VBD	Mood=Ind Tense=Past VerbForm=Fin				
<i>Egremo</i>	Egremo	PROPN	NNP	Number=Sing		F		
<i>nt</i>	not	PART	RB	_	F	F	F	1

- Random sample of 5,000 tokens for each language + tagger
 - across all documents
 - length 3 tokens < complete sentences < length 30 tokens
- Annotation in a spreadsheet editor (MS Excel)
- True/False, number of errors in tokens

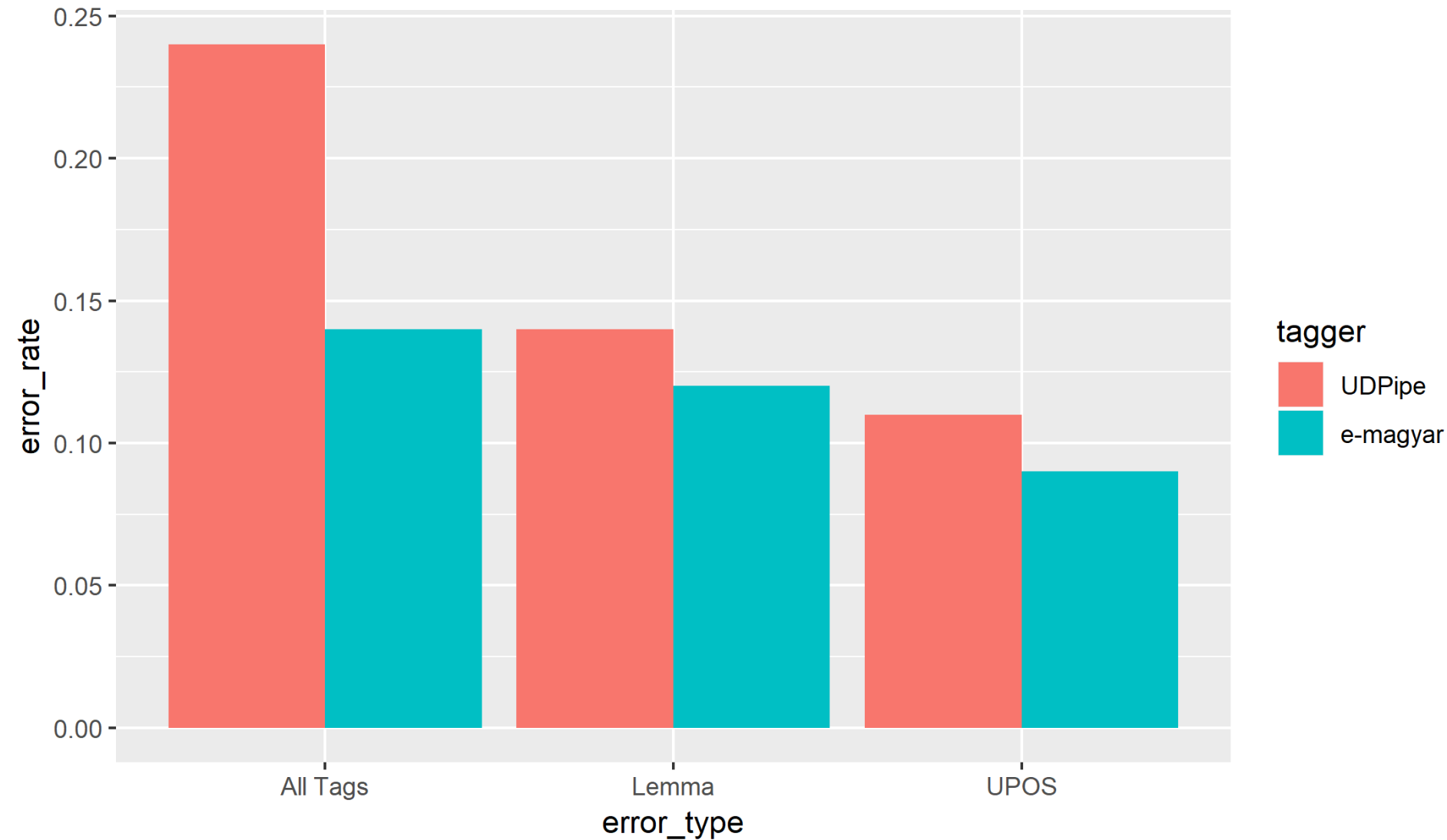


UDPipe reference vs. ELTeC

	reference	ELTeC
AllTags	F-measure all	errors/all except features ⚠
Feats	F-measure	log100k(geom.mean errors)
Lemma	F-measure	errors/all
UPOS	F-measure	errors/all

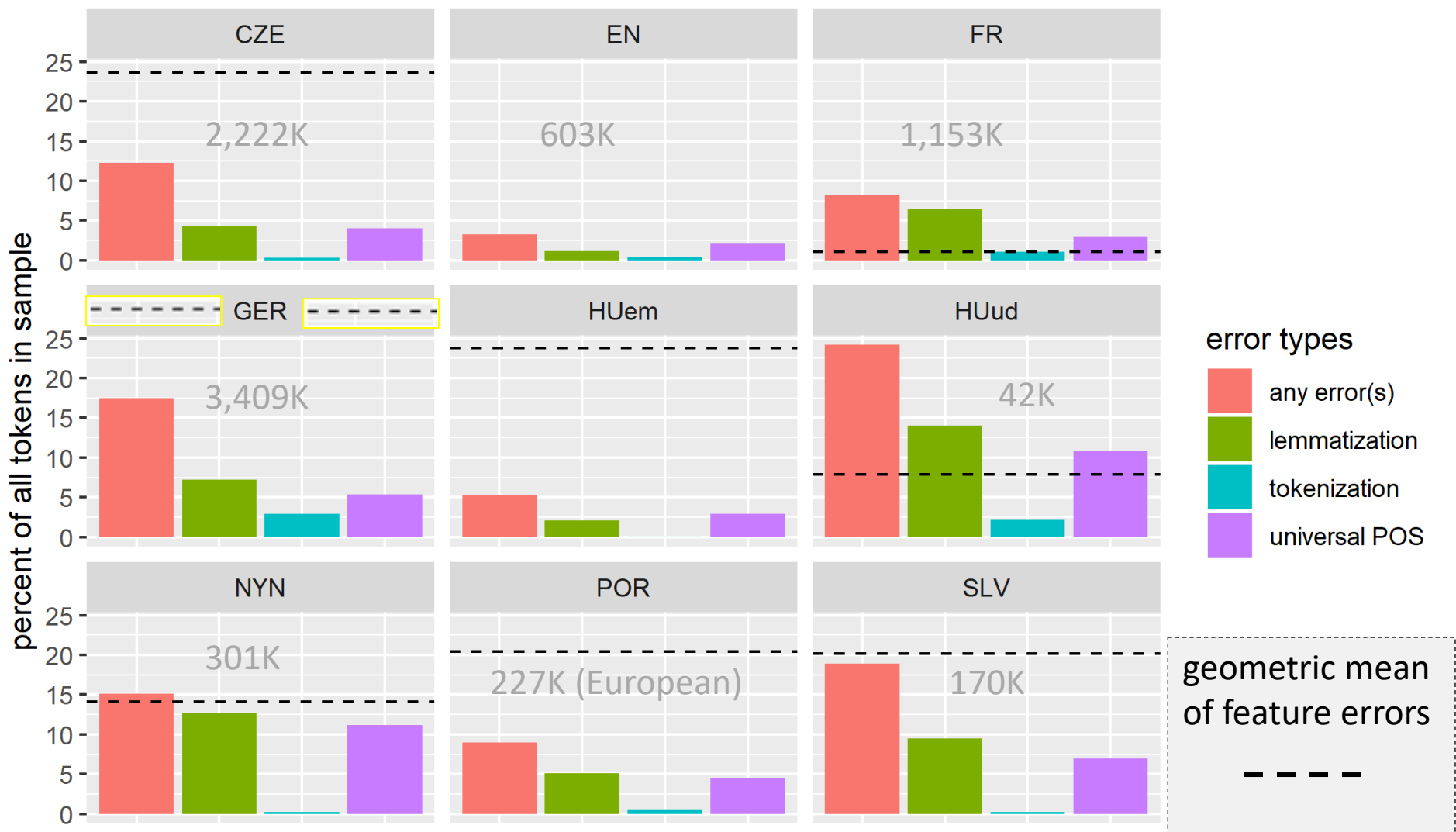


Hungarian taggers on ELTeC texts



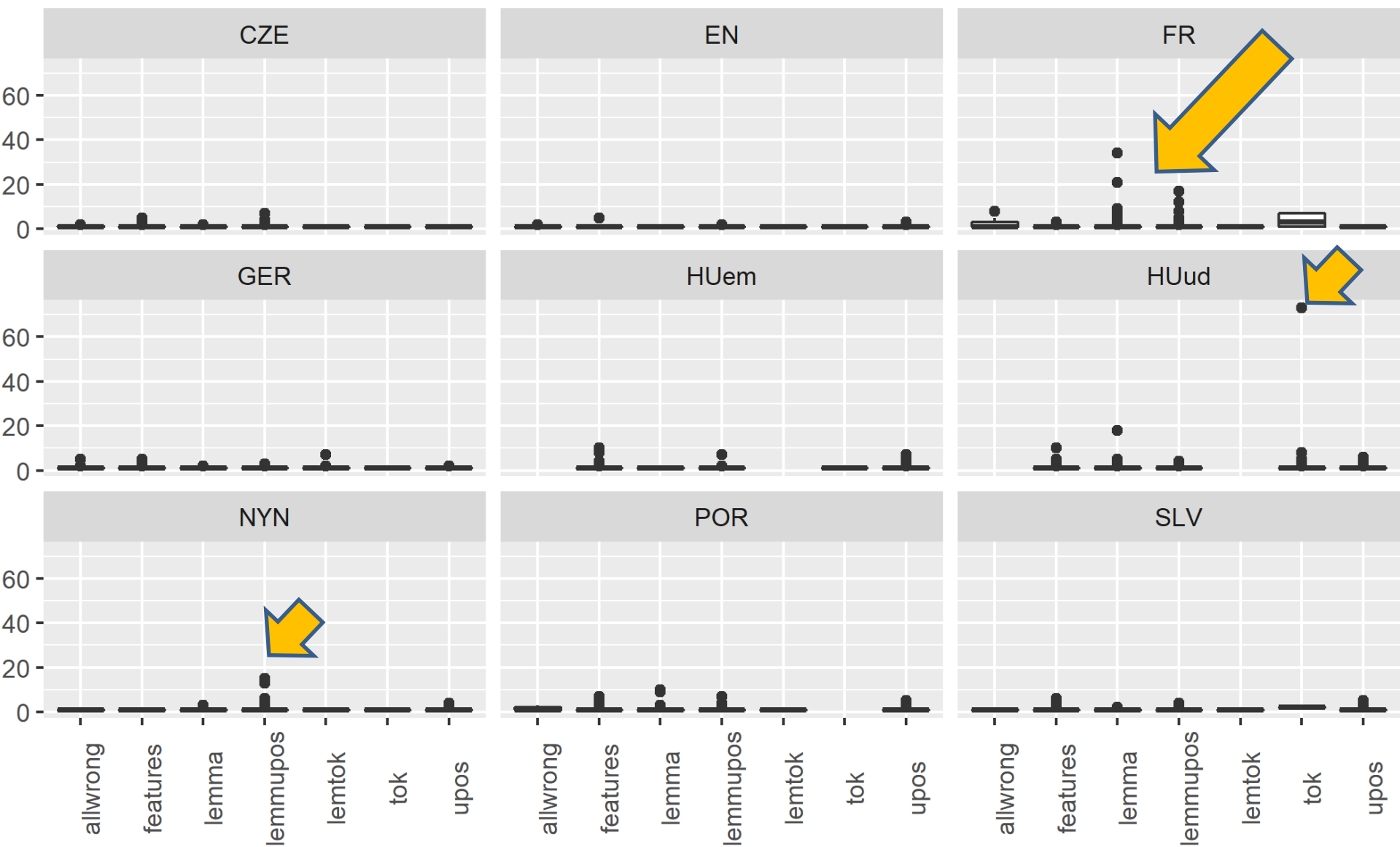


UD tagging errors for individual languages



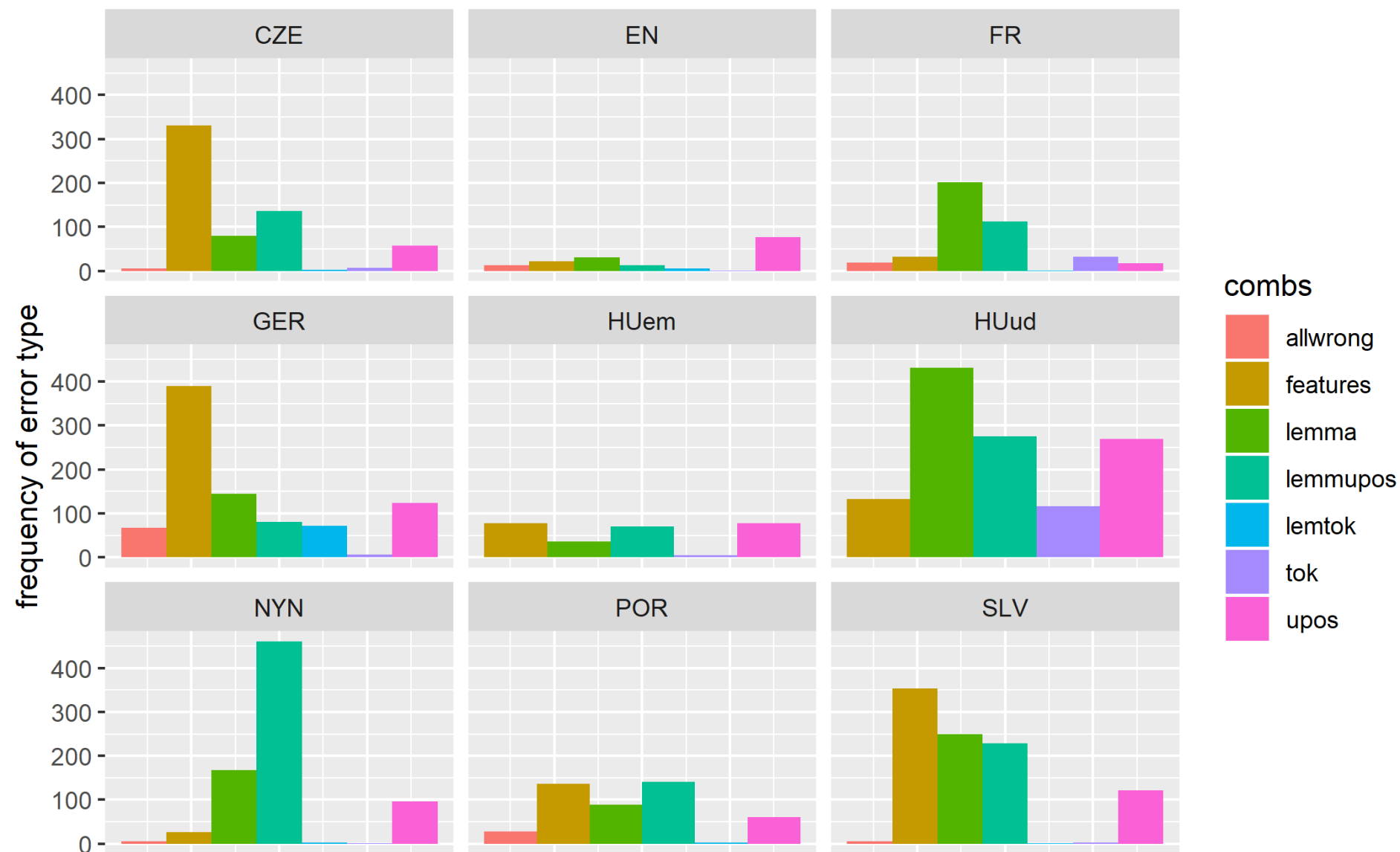


Frequency distribution of erroneous tokens





Distribution of error combinations





allwrong

propose to
 some anglo-
 an nt
 morrow hate—

features

when what terrified
 flashed ye left
 look burned
 mine blazed swelled
 do crept that
 further her where
 waxed

lemma

wharves scorned
 parting ornament invigorating
 mine flouted drawled
 shot bound alarming woven
 throes dared's's' d' escaped
 hers 're bade rushed
 entreated committed
 hooted hesitated
 remembered interferences
 worrying singing slanting

lemmupos

kneeling kindred
 done bedyer
 whom 'ere yer
 're
 browsing
 swathing

lemtok

i one
 ? i
 no—

tok

egremo

freq

a	1
a	2
a	3
a	4
a	5

combs

a	allwrong
a	features
a	lemma
a	lemmupos
a	lemtok
a	tok

EN: wrongly tagged as pronouns



freq

a	1
a	2
a	3
a	4
a	5

combs

a	features
a	lemma
a	lemtok
a	upos



Future work

- Serbian annotations
- manual error analysis - word clouds
- Training data for ELTeC domain necessary? **BAMTINOF** guessing error classification
 - count wrong/correct guesses on words absent in gold standard
 - classify guessing errors: BAMTINOF
 - diachr. word **B**order changes
 - syntactically **A**mbiguous word form
 - archaic or poetic **M**orphology
 - **T**ypo
 - token**I**zation error resulting in a non-existent token
 - **N**ame, proper noun
 - **O**ld spelling or archaic word typically replaced by a modern word
 - **F**oreign-language or heavy dialect