

# From Characters to Words to in Between: Do We Capture Morphology?

Clara Vania and Adam Lopez  
{c.vania@ed.ac.uk, alopez@inf.ed.ac.uk}

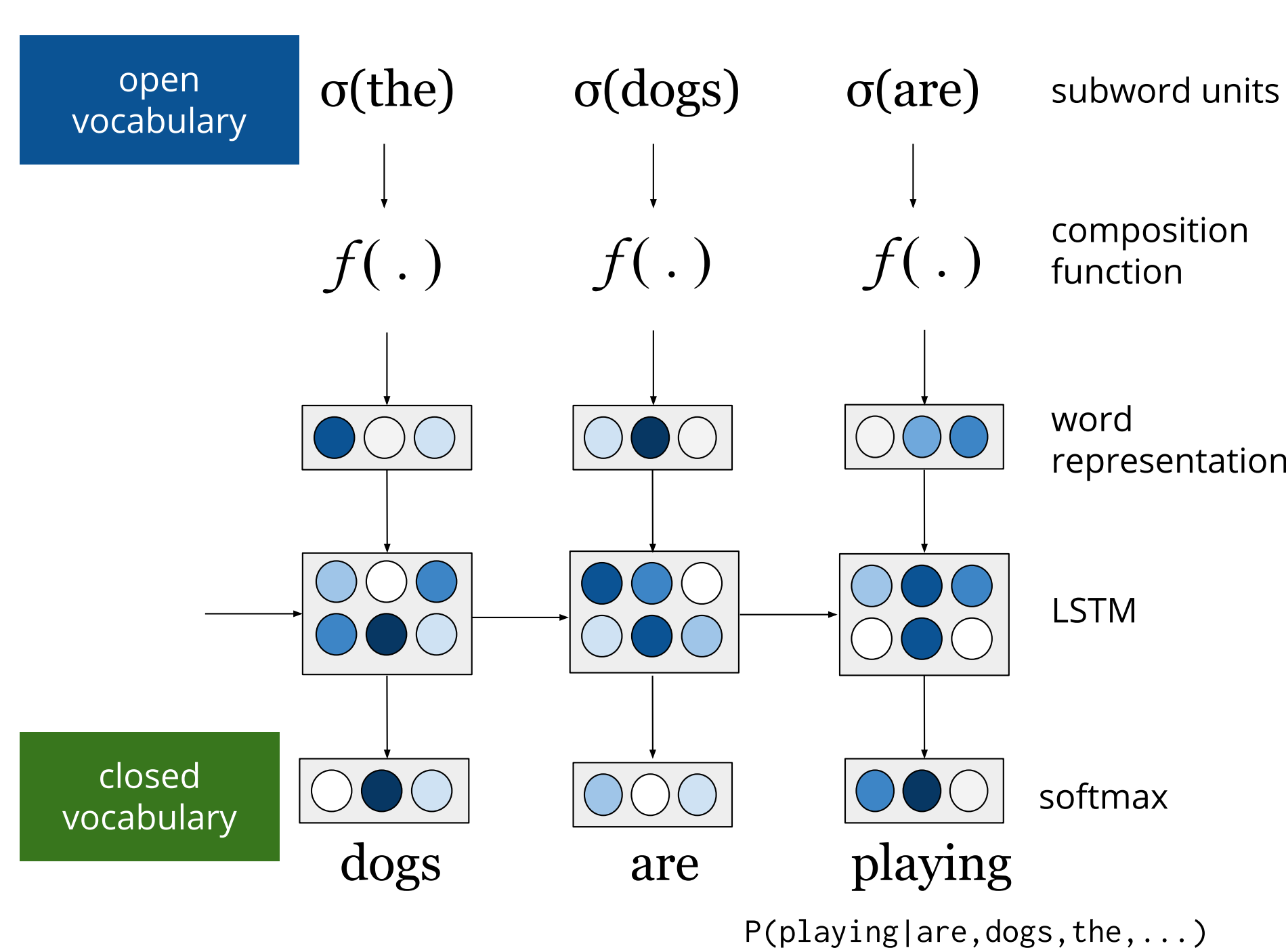
Institute for Language, Cognition, and Computation,  
School of Informatics, University of Edinburgh

## Motivation

Continuous word representations composed from subword representations have shown to be effective for learning the morphological regularities of words. But, some questions remain:

- Type of subword units: characters vs. morphemes?
- How to compose them: addition, bi-LSTM, or CNN?
- Do character-level models capture morphology in terms of predictive utility?
- How do they interact with languages of different morphological typologies?

## Task: Language Modeling



## Variable: Subword Units

Unit	Output of $\sigma(wants)$
Morfessor	want, s
BPE	w, ants
char-trigram	$\hat{w}a, wan, ant, nts, ts\hat{s}$
character	w, a, n, t, s
analysis	want, +VB, +3rd, +SG, +PRS

## Variable: Composition Function

- vector addition, bi-LSTM, CNN.

## Variable: Language Typology

### Concatenative

Agglutinative (Turkish)	Fusional (English)
oku-r-sa-m	read-s
read-AOR.COND.1SG	read-3SG.PRS
'If I read ...'	'reads'

### Non-concatenative

Root&Pattern (Arabic)	Reduplication (Indonesian)
$k\langle a \rangle t\langle a \rangle b\langle a \rangle$	anak~anak
write-PST.3SG.M	child-PL
'he wrote'	'children'

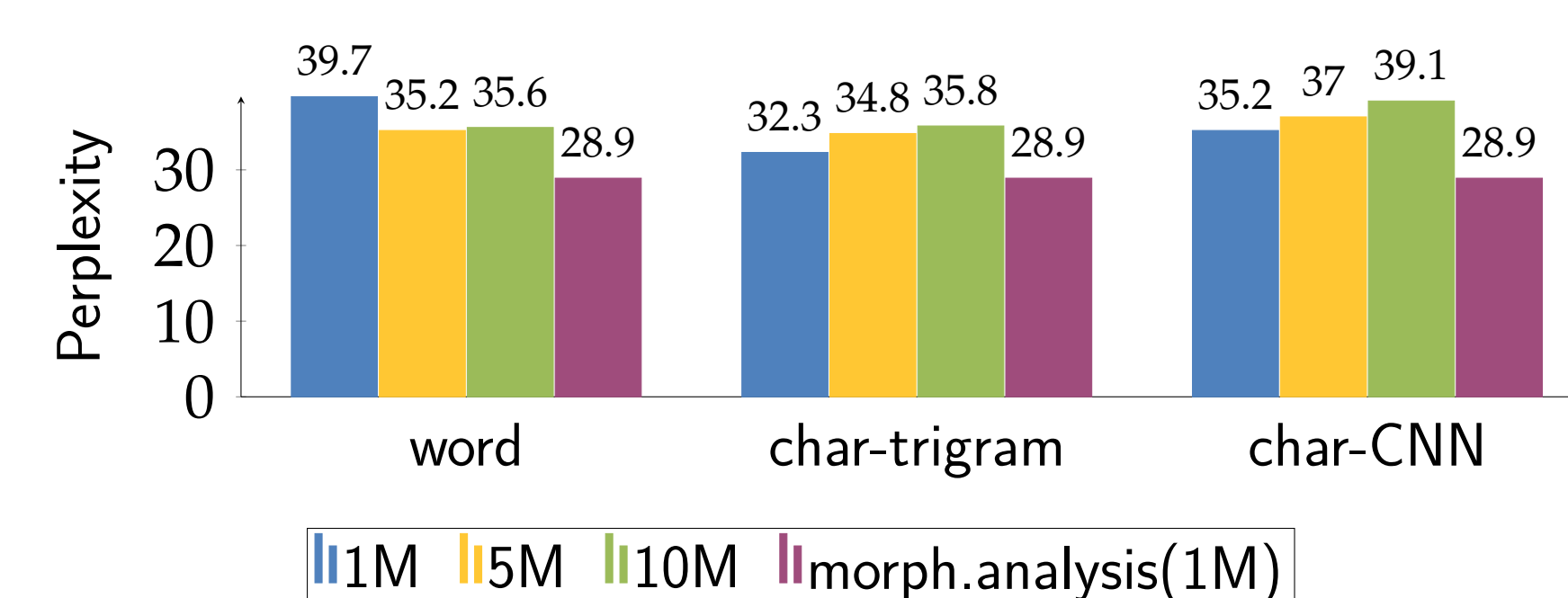
## Qualitative Analysis

Subword Unit	In-Vocabulary	Rare	OOV
	<i>including</i>	<i>unconditional</i>	<i>uploading</i>
BPE	called	unintentional	upbeat
bi-LSTM	involve	ungenerous	uprising
	like	unanimous	handling
	creating	unpalatable	hand-colored
character trigram	include	unconstitutional	drifted
bi-LSTM	includes	constitutional	affected
	undermining	unimolecular	conflicted
	include	medicinal	convicted
character bi-LSTM	inclusion	undamaged	musagète
	insularity	unmyelinated	mutualism
	includes	unconditionally	mutualists
	include	uncoordinated	meursault

## Do character-level models capture morphology in terms of predictive utility?

Language	Addition	bi-LSTM
Czech	51.8	<b>30.1</b>
Russian	41.8	<b>26.4</b>

## How much training data is needed to reach perplexity obtained using annotated data?



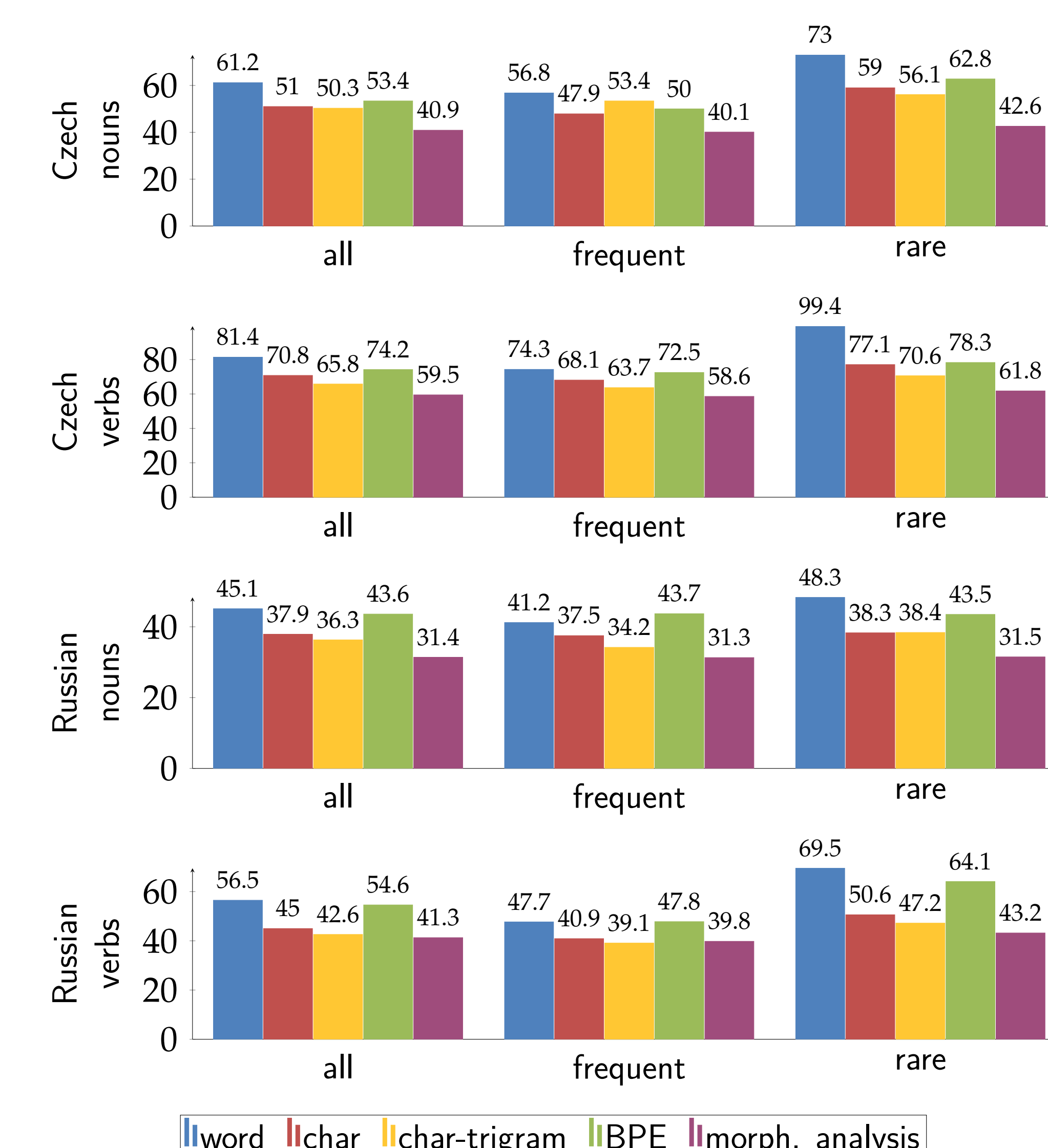
## Perplexity Results

Typology	lang	word	character		char-trigram		BPE		Morfessor		%imp
			bi-lstm	CNN	add	bi-lstm	add	bi-lstm	add	bi-lstm	
Fusional	Czech	41.5	34.2	36.6	42.7	<b>33.6</b>	50.0	33.7	47.7	36.9	19.0
	English	46.4	43.5	44.7	45.4	<b>43.0</b>	47.5	43.3	49.7	49.7	7.4
	Russian	34.9	28.4	29.5	35.2	<b>27.7</b>	40.1	28.5	39.6	31.3	20.6
Agglutinative	Finnish	24.2	20.1	20.3	24.9	<b>18.6</b>	26.8	19.1	27.8	22.5	23.1
	Japanese	98.1	98.1	<b>91.6</b>	102.0	101.1	126.5	96.8	112.0	99.2	6.6
	Turkish	67.0	54.5	55.1	<b>50.1</b>	54.2	59.5	57.3	62.2	62.7	25.2
Root & Pattern	Arabic	48.2	42.0	43.2	50.9	<b>39.9</b>	50.9	42.8	52.9	45.5	17.3
	Hebrew	38.2	31.6	33.2	39.7	<b>30.4</b>	44.2	32.9	44.9	34.3	20.5
Reduplication	Indonesian	46.1	45.5	46.6	58.5	46.0	59.2	<b>43.4</b>	59.3	44.9	5.9
	Malaysian	54.7	53.0	<b>50.6</b>	68.5	50.7	69.0	51.2	68.2	52.5	7.5

## How do we know if these representations actually affect the predictions?

Analyze perplexities when the inflected words of interest are in the most recent history.

"The girl **loves dogs** but the boy **does not**."

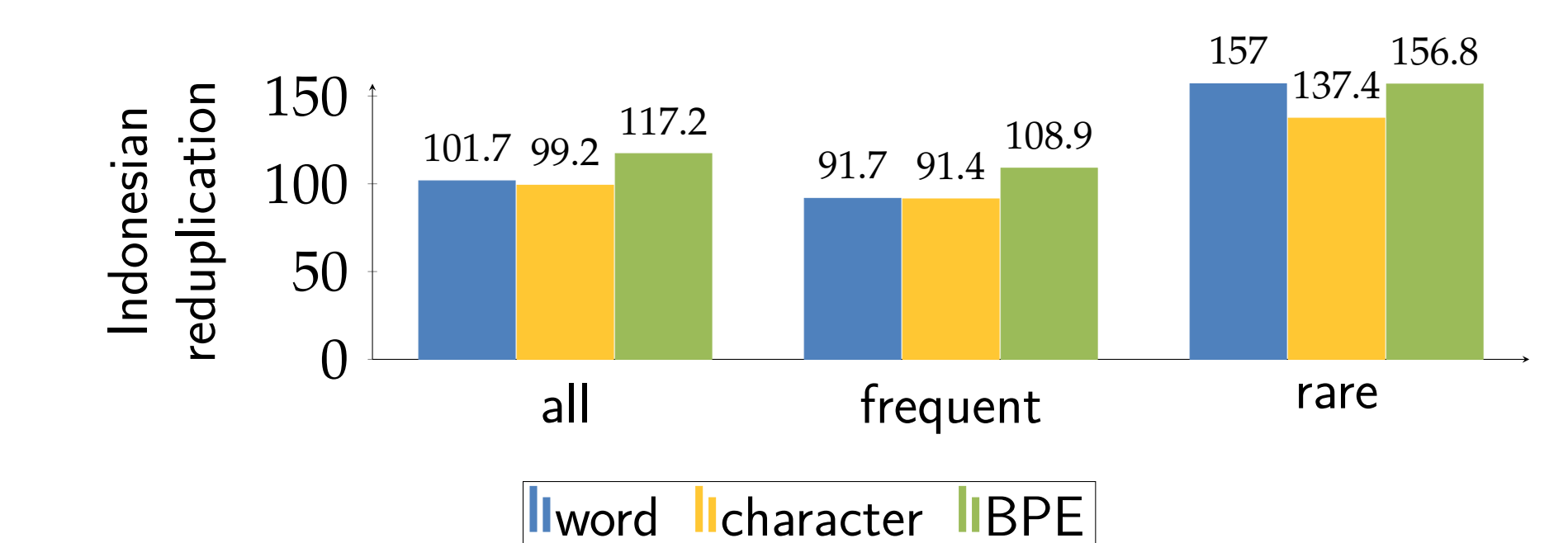


## Which model most effectively captures reduplication?

Percentage of full reduplication on the training data:

Language	type-level (%)	token-level (%)
Indonesian	1.1	2.6
Malay	1.3	2.9

"Saya membeli **buku-buku itu** kemarin."  
I bought **those books** yesterday.



## Conclusion

- Character-level models are effective for many languages, but these models do not match the predictive accuracy of model with explicit knowledge of morphology.
- In this study, a previously unstudied combination of character trigram composed with bi-LSTM outperform most others.
- Our qualitative analysis suggests that they learn orthographic similarity of affixes.
- Other factors such as morphology and orthography affect the utility of these representations.