

# Enriching word vectors with subword information

(paper review)

Khoklin I., Kovrigin A., Shibaev E.

# Motivation

Why embedding words considering each word as a **whole** is not the best idea?

# Motivation: morphologically rich languages

- Finnish has **15 cases for nouns**.
- In French and Spanish most verbs have more than **40 different inflected forms**.

The verb *parler* "to speak", in French orthography and IPA transcription

	Indicative				Subjunctive		Conditional	Imperative
	Present	Simple past	Imperfect	Simple future	Present	Imperfect	Present	Present
<b>Je</b>	parl-e /paʁl/	parl-ai /paʁle/	parl-ais /paʁlɛ/	parl-erai /paʁlɛʁe/	parl-e /paʁl/	parl-asse /paʁlas/	parl-erais /paʁlɛʁɛ/	
<b>tu</b>	parl-es /paʁl/	parl-as /paʁla/	parl-ais /paʁlɛ/	parl-eras /paʁlɛʁa/	parl-es /paʁl/	parl-asses /paʁlas/	parl-erais /paʁlɛʁɛ/	parl-e /paʁl/
<b>il</b>	parl-e /paʁl/	parl-a /paʁla/	parl-ait /paʁlɛ/	parl-era /paʁlɛʁa/	parl-e /paʁl/	parl-ât /paʁla/	parl-erait /paʁlɛʁɛ/	
<b>nous</b>	parl-ons /paʁlɔ̃/	parl-âmes /paʁlam/	parl-ions /paʁljɔ̃/	parl-erons /paʁlɛʁɔ̃/	parl-ions /paʁljɔ̃/	parl-assions /paʁlasjɔ̃/	parl-erions /paʁlɛʁjɔ̃/	parl-ons /paʁlɔ̃/
<b>vous</b>	parl-ez /paʁle/	parl-âtes /paʁlat/	parl-iez /paʁlje/	parl-erez /paʁlɛʁe/	parl-iez /paʁlje/	parl-assiez /paʁlasje/	parl-eriez /paʁlɛʁje/	parl-ez /paʁle/
<b>ils</b>	parl-ent /paʁl/	parl-èrent /paʁlɛʁɛ̃/	parl-aient /paʁlɛ/	parl-eront /paʁlɛʁɔ̃/	parl-ent /paʁl/	parl-assent /paʁlas/	parl-eraient /paʁlɛʁɛ̃/	

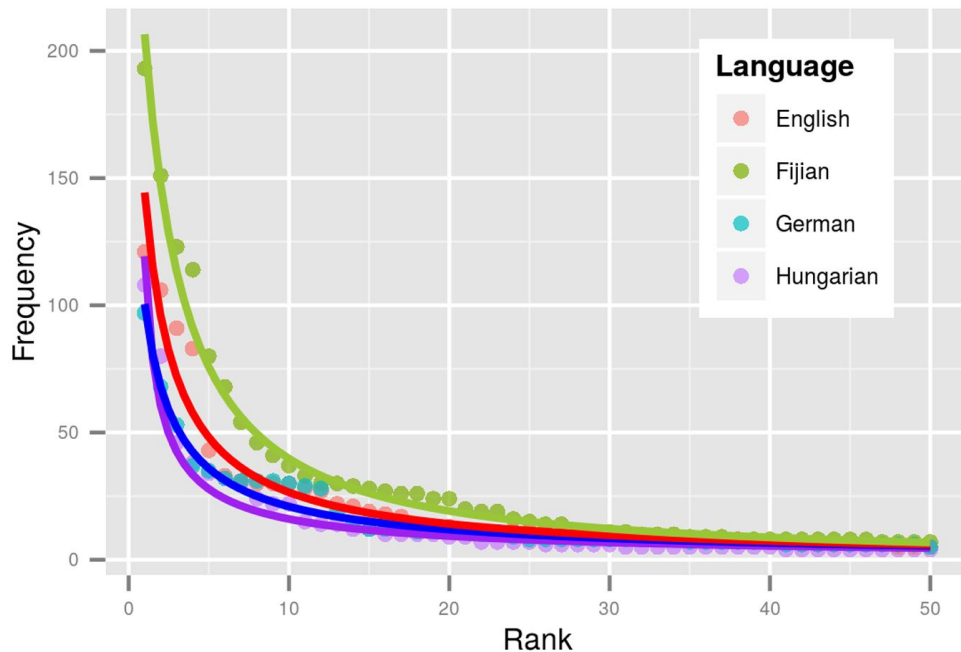
The conjugations of a verb "to speak" in French.

talo 'house'	singular	plural
<b>nominative</b>	talo	talot
<b>genitive</b>	talon	talojen
<b>partitive</b>	taloa	taloja
<b>inessive</b>	talossa	taloissa
<b>elative</b>	talosta	taloista
<b>illative</b>	taloon	taloihin
<b>adessive</b>	talolla	taloilla
<b>ablative</b>	talolta	taloilta
<b>allative</b>	talolle	taloille
<b>essive</b>	talona	taloina
<b>translative</b>	taloksi	taloiksi
<b>instructive</b>	-	taloin
<b>abessive</b>	talotta	taloitta
<b>comitative</b>	-	taloine+POS

All declinations of a word "house" in Finnish

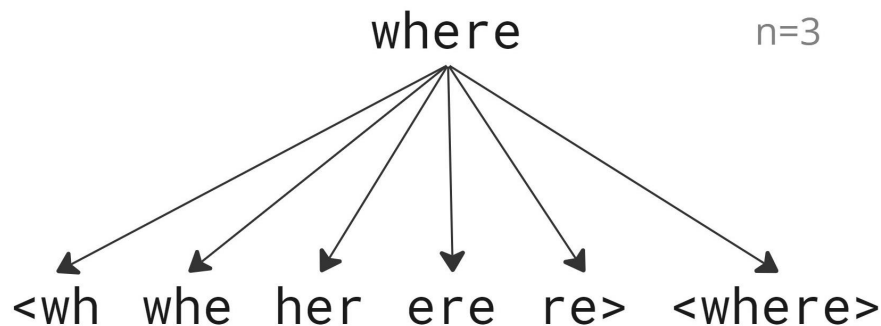
# Motivation: rare words & amount of data for training

- Languages have **big tails of rare words**.
- Hence, “whole word embeddings” require an **immense amount of data** to learn properly and, for rare words, the **embeddings end up being bad**.



# Proposal: Subword model

- Proposal: **represent each word by character n-grams**
- Word embedding is a **sum of n-gram embeddings**
- Technicalities:
  - Boundary symbols < and > added
  - The whole word is added as a special sequence
  - In practice, all n-grams where  $3 \leq n \leq 6$  are taken



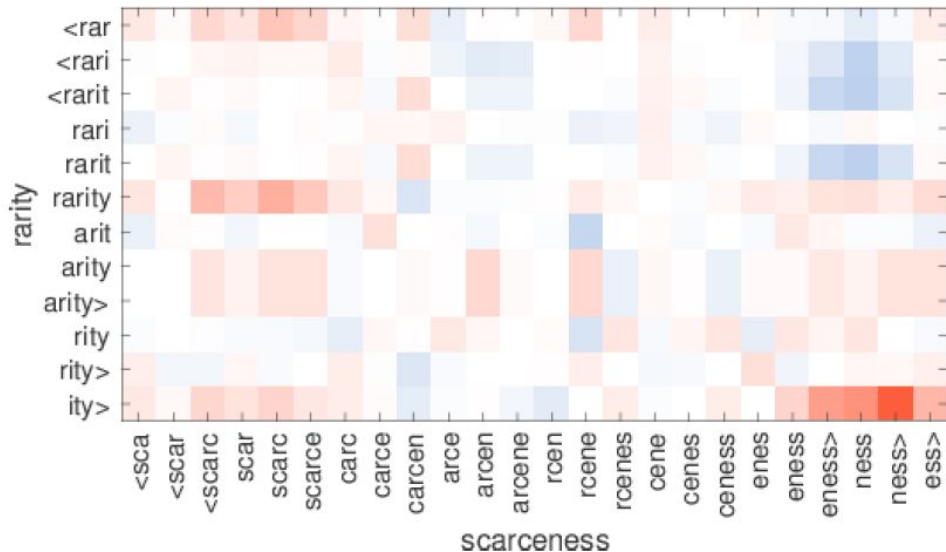
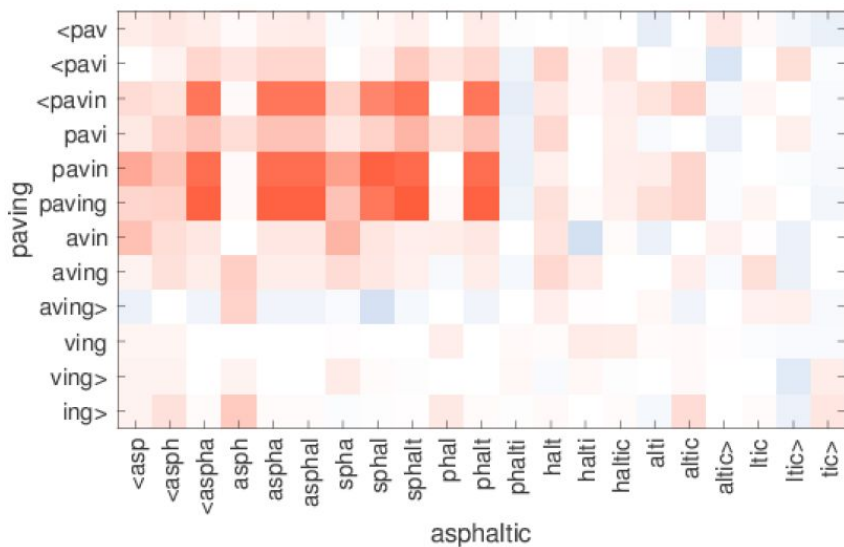
# Qualitative analysis: Subword importance

- Analyze learned embeddings
- **Rank n-grams by importance**
  - Remove n-gram from the sum
  - Compute cosine similarity between word embedding and word embedding without the n-gram
  - **Lower similarity → higher importance**

	word	n-grams		
DE	autofahrer	fahr	fahrer	auto
	freundeskreis	kreis	kreis>	<freun
	grundwort	wort	wort>	grund
	sprachschule	schul	hschul	sprach
	tageslicht	licht	gesl	tages
EN	anarchy	chy	<anar	narchy
	monarchy	monarc	chy	<monar
	kindness	ness>	ness	kind
	politeness	polite	ness>	eness>
	unlucky	<un	cky>	nlucky
	lifetime	life	<life	time
	starfish	fish	fish>	star
	submarine	marine	sub	marin
	transform	trans	<trans	form
FR	finirais	ais>	nir	fini
	finissent	ent>	finiss	<finis
	finissions	ions>	finiss	sions>

# Qualitative analysis: Subword similarity

- Discover which n-grams are considered similar
- Calculate **cosine similarity** between subword embeddings



# Results. Human similarity judgement

The authors compared word similarity through human judgment and cosine similarity of generated embeddings across various languages. They also distinguished between common and rare English words in their datasets.

Their method was compared with CBOW and Skip-gram baselines, showing superior performance on all datasets, with the exception of the English WS353 dataset. Notably, the proposed model's computation of vectors for out-of-vocabulary words (sisg) consistently surpassed the alternative approach of not computing these vectors (sisg-).

		sg	cbow	sisg-	sisg
AR	WS353	51	52	54	<b>55</b>
	GUR350	61	62	64	<b>70</b>
DE	GUR65	78	78	<b>81</b>	<b>81</b>
	ZG222	35	38	41	<b>44</b>
EN	RW	43	43	46	<b>47</b>
	WS353	72	<b>73</b>	71	71
Es	WS353	57	58	58	<b>59</b>
FR	RG65	70	69	<b>75</b>	<b>75</b>
RO	WS353	48	52	51	<b>54</b>
RU	HJ	59	60	60	<b>66</b>

*There are two ways to deal with out of vocabulary words:*

- *sigs – sum of n-grams*
- *sigs- – null vector*



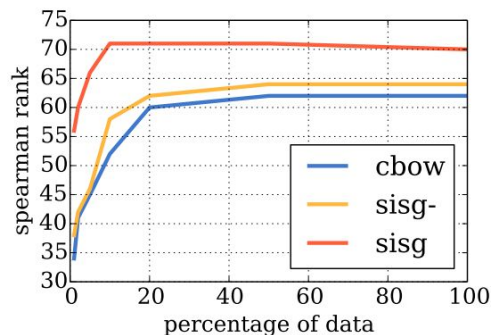
## Results. Word analogy tasks

Word analogy task is questions, of the form A is to B as C is to D, where D must be predicted by the models.

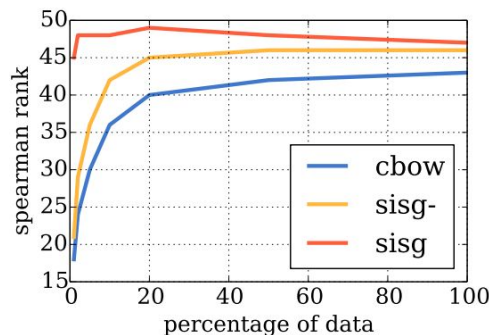
They observe that morphological information significantly improves the syntactic tasks; the approach outperforms the baselines. In contrast, it does not help for semantic questions, and even degrades the performance for German and Italian.

		sg	cbow	sisg
Cs	Semantic	25.7	27.6	27.5
	Syntactic	52.8	55.0	77.8
DE	Semantic	66.5	66.8	62.3
	Syntactic	44.5	45.0	56.4
EN	Semantic	78.5	78.2	77.8
	Syntactic	70.1	69.9	74.9
IT	Semantic	52.3	54.7	52.3
	Syntactic	51.5	51.8	62.7

# Results. Effect of the size of the training data



(a) DE-GUR350



(b) EN-RW

The approach is more robust because it is able to model infrequent words: authors used different proportion of the training dataset and measured the performance on the test dataset.

They noticed that:

- For all datasets, and all sizes, the proposed approach (sisg) performs better than the baseline.
- Proposed approach provides very good word vectors even when using very small training datasets.

# Results. Summary

- ***Morphologically Rich Languages***: The model works really well with languages that have lots of different word forms, outperforming others because it's great at understanding their complexity.
- ***Handling of Rare Words***: Leveraging n-grams, the model effectively captures information from rare words, ensuring robust and meaningful embeddings..
- ***Training on Limited Data***: Demonstrating efficiency, the model achieves impressive results even with a small proportion of training data, making it suitable for tasks with restricted dataset sizes.

Thank you for your attention!