

Better OOV Translation with Bilingual Terminology Mining

Matthias Huck, Viktor Hangya, Alexander Fraser

LMU Munich

31 July 2019

Subword segmentation allows for open-vocabulary translation, but out-of-vocabulary words (OOVs) are still often mistranslated.

Example:

<i>src</i>	A coronary angioplasty may not be technically possible [. . .]
<i>ref</i>	Eine Koronarangioplastie ist wahrscheinlich technisch nicht möglich [. . .]
<i>hyp</i>	Ein Herzinfarkt (<i>heart attack</i>) ist vielleicht technisch nicht möglich [. . .]

“OOVs”:

Source language words that weren't observed in the parallel training corpus

Can adequate translations of OOV words be learned from additional monolingual corpora?

Bilingual word embeddings (BWEs)

- Represent source and target language words in a joint space
- Higher word vocabulary coverage than the parallel corpus

How to best integrate OOV word translation candidates from the BWE space into the NMT system?

- Cross-lingual nearest neighbors in the BWE space are noisy
- Polysemy: Need to disambiguate – choose amongst multiple options depending on context within sentences

① Baseline NMT system

- Trained on parallel corpus (subword-segmented)

② (Unsupervised) BWEs

- Trained on large monolingual data in the two languages

③ Bilingual terminology mining

- Identify test set OOVs & get top-n word translations from BWEs
- In target-language monolingual data, mine sentences that contain the OOV translation candidates

④ NMT fine-tuning

- Backtranslate the mined target-side sentences, force OOV words to be generated in the backtranslations
- Fine-tune NMT model on synthetic data (subword-segmented)

① Baseline NMT system

- Trained on parallel corpus (subword-segmented)

② (Unsupervised) BWEs

- Trained on large monolingual data in the two languages

③ Bilingual terminology mining

- Identify test set OOVs & get top-n word translations from BWEs
- In target-language monolingual data, mine sentences that contain the OOV translation candidates

④ NMT fine-tuning

- Backtranslate the mined target-side sentences, force OOV words to be generated in the backtranslations
- Fine-tune NMT model on synthetic data (subword-segmented)

① Baseline NMT system

- Trained on parallel corpus (subword-segmented)

② (Unsupervised) BWEs

- Trained on large monolingual data in the two languages

③ Bilingual terminology mining

- Identify test set OOVs & get top-n word translations from BWEs
- In target-language monolingual data, mine sentences that contain the OOV translation candidates

④ NMT fine-tuning

- Backtranslate the mined target-side sentences, force OOV words to be generated in the backtranslations
- Fine-tune NMT model on synthetic data (subword-segmented)

① Baseline NMT system

- Trained on parallel corpus (subword-segmented)

② (Unsupervised) BWEs

- Trained on large monolingual data in the two languages

③ Bilingual terminology mining

- Identify test set OOVs & get top-n word translations from BWEs
- In target-language monolingual data, mine sentences that contain the OOV translation candidates

④ NMT fine-tuning

- Backtranslate the mined target-side sentences, force OOV words to be generated in the backtranslations
- Fine-tune NMT model on synthetic data (subword-segmented)

src if you need to take medication for eye health , make sure you take as prescribed and don 't stop without talking to your GP or **optometrist** .

Top-5 word translations from BWEs for the OOV “optometrist”:

- Gesichtsfeldprüfgerät (*visual field checking device*)
- Augenarzt (*eye doctor*)
- Bildanzeigeverfahren (*image display method*)
- Sehtests (*vision test*)
- Sehtestgerät (*eyesight test device*)

Braune et al. (2018): cosine combined with orthography

src if you need to take medication for eye health , make sure you take as prescribed and don 't stop without talking to your GP or **optometrist** .

top-5 **Gesichtsfeldprüfgerät** | **Augenarzt** | **Bildanzeigeverfahren** | **Sehtests** | **Sehtestgerät**

Mine target-language monolingual sentences with OOV translation candidates:

- kompaktes **Gesichtsfeldprüfgerät** nach Anspruch 2 [. . .]
- bei einer Beeinträchtigung des Sehens oder der Augen während der Behandlung wenden Sie sich bitte umgehend an Ihren **Augenarzt** .
- Bildanzeigeeinheit , **Bildanzeigeverfahren** und Bildanzeigeprogramm
- die Erfordernis eines jährlichen Hör- und **Sehtests**
- die Erfindung betrifft ein Verfahren und ein **Sehtestgerät** zur Ermittlung der Notwendigkeit einer Sehhilfe bei Dunkelheit [. . .]

mined bei einer Beeinträchtigung des Sehens oder der Augen während der
Behandlung wenden Sie sich bitte umgehend an Ihren **Augenarzt** .

mined bei einer Beeinträchtigung des Sehens oder der Augen während der
Behandlung wenden Sie sich bitte umgehend an Ihren **OOV** .

mined bei einer Beeinträchtigung des Sehens oder der Augen während der Behandlung wenden Sie sich bitte umgehend an Ihren **OOV** .

bt you are turning straight to your **OOV** in the event of interference in the treatment or the eye during the treatment .

mined bei einer Beeinträchtigung des Sehens oder der Augen während der Behandlung wenden Sie sich bitte umgehend an Ihren **Augenarzt** .

bt you are turning straight to your **optometrist** in the event of interference in the treatment or the eye during the treatment .

mined bei einer Beeinträchtigung des Sehens oder der Augen während der Behandlung wenden Sie sich bitte umgehend an Ihren **Augenarzt** .

bt you are turning straight to your **optometrist** in the event of interference in the treatment or the eye during the treatment .

mined die Erfordernis eines jährlichen Hör- und **Sehtests** (*vision test*) .

bt the requirement for an annual hearing and **optometrist** .

Evaluation: Machine Translation Quality

	BLEU	
	Cochrane	NHS24
baseline	22.4	20.2
with OOV copying	23.4	20.5
fine-tuned with OOV terminology mining	27.2	22.5

Examples: Better OOV Translations

<i>src</i>	[...] without talking to your GP or optometrist
<i>ref</i>	[...] ohne vorherige Rücksprache mit Ihrem Hausarzt oder Optiker (<i>optician</i>)
<i>base</i>	[...] ohne mit Ihrem Arzt oder Ihrem Arzt (<i>physician</i>) zu sprechen
<i>ours</i>	[...] ohne mit Ihrem Arzt oder Augenarzt (<i>eye doctor</i>) zu sprechen

Examples: Better OOV Translations

<i>src</i>	A coronary angioplasty may not be technically possible [...]
<i>ref</i>	Eine Koronarangioplastie ist wahrscheinlich technisch nicht möglich [...]
<i>base</i>	Ein Herzinfarkt (<i>heart attack</i>) ist vielleicht technisch nicht möglich [...]
<i>ours</i>	Eine koronare Angioplastie ist möglicherweise nicht technisch möglich [...]

Examples: Better OOV Translations

<i>src</i>	regular nosebleeds
<i>ref</i>	regelmäßige Nasenbluten
<i>base</i>	regelmäßige Misskredite (<i>discredits</i>)
<i>ours</i>	regelmäßige Nasenbluten

Examples: Better OOV Translations

<i>src</i>	dizziness or lightheadedness
<i>ref</i>	Schwindel oder Benommenheit
<i>base</i>	schwindelerregend (<i>dizzying</i>) oder zurückhaltend (<i>reluctant</i>)
<i>ours</i>	Schwindel oder Schwächegefühl (<i>feeling of faintness</i>)

Examples: Better OOV Translations

-
- src* Four different alpha blockers were tested
(**alfuzosin**, **tamsulosin**, **doxazosin** and **silodosin**).
- ref* Vier verschiedene Alphablocker wurden getestet
(**Alfuzosin**, **Tamsulosin**, **Doxazosin** und **Silodosin**).
- base* Vier verschiedene Alphablocker wurden getestet
(**alfuzos**, **tasulo**, **doxasa** und **silodosin**).
- ours* Vier unterschiedliche Alphablocker wurden untersucht
(**Alfuzosin**, **Tamsulosin**, **Doxazosin** und **Tigecyclin**).
-

BWEs help adequately translate vocabulary which isn't present in parallel training data.

- We've presented a simple approach to effectively integrate BWE-suggested OOV word translation candidates into an NMT system
- **Bilingual terminology mining**
& **backtranslation with forced OOV words**
& **finetuning**
- Multiple candidates provided from the BWEs that the NMT system can choose from

Thank you for your attention

Matthias Huck

mhuck@cis.lmu.de



This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement № 640550).



Braune, F., Hangya, V., Eder, T., and Fraser, A. (2018). Evaluating bilingual word embeddings on the long tail. In *Proc. NAACL-HLT*.

Intrinsic Evaluation: Bilingual Lexicon Induction

		Acc_1	Acc_5
freq	Braune et al. (2018)	38.6	47.4
	EU+UFAL+orth	25.9	40.6
rare	Braune et al. (2018)	26.3	28.2
	EU+UFAL+orth	17.5	28.8

Table: Medical bilingual lexicon induction results showing the quality of the BWE based dictionaries using top-1 and top-5 translations.

Intrinsic Evaluation: Bilingual Terminology Mining

<i>n</i>	UFAL			UFAL+orth		
	<i>P@n</i>	<i>R@n</i>	<i>F₁@n</i>	<i>P@n</i>	<i>R@n</i>	<i>F₁@n</i>
1	58.19	13.58	22.02	58.13	25.28	35.24
5	44.46	26.10	32.89	50.05	43.82	46.73
10	35.80	29.84	32.55	41.04	47.64	44.09
20	29.54	33.58	31.43	34.43	50.16	40.83
<i>n</i>	EU+UFAL			EU+UFAL+orth		
	<i>P@n</i>	<i>R@n</i>	<i>F₁@n</i>	<i>P@n</i>	<i>R@n</i>	<i>F₁@n</i>
1	68.65	37.56	48.55	69.59	41.87	52.28
5	54.33	48.46	51.22	51.13	51.71	51.41
10	42.94	53.41	47.61	44.45	56.34	49.70
20	36.42	58.78	44.98	37.42	61.30	46.47

Table: Mining quality with top-*n* OOV word translation candidates.