

Bilingual Learning of Multi-sense Embeddings with Discrete Autoencoders

Simon Šuster*, Ivan Titov \diamond , and Gertjan van Noord*

*CLCG, University of Groningen

\diamond ILLC, University of Amsterdam

CROSSLINGUAL SUPERVISION

Multilingual learning: more accurate monolingual models

- Ambiguity in L1 often different than ambiguity in L2
[Snyder and Barzilay, 2010]
- Resolving polysemy in L1 by also looking at *translations*
 - well known in WSD (“translation as sense”)
[Diab and Resnik, 2002]
 - but little explored in representation learning
[Guo et al., 2014, Ettinger et al., 2016]
- Availability of parallel corpora

EN-ES INTUITION

track: a course of study; a piece of music, a rough path...

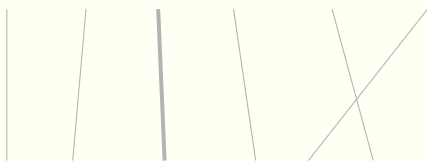
sent_{t1}: *Choose a track that interests you*

EN-ES INTUITION

track: a course of study; a piece of music, a rough path...

sent_{L1}: Choose a track that interests you

sent_{L2}: Pon una canción que te gusta



RELATED EMBEDDING RESEARCH

- ❖ Multi-sense [Neelakantan et al., 2014, Li and Jurafsky, 2015]
 - ❖ deal with polysemy explicitly
 - ❖ monolingual
- ❖ Multilingual
 - ❖ embeddings in the same semantic space [Gouws et al., 2014, Klementiev et al., 2012]
 - ❖ use target-language signal for better source-language embeddings [Hill et al., 2014, Faruqui and Dyer, 2014]
- ❖ Better L1 **multi-sense** embeddings **with L2 signal**?

INSPIRATION FROM AUTOENCODERS

At output, reconstruct input by relying on a latent representation

- Latent sense representation is a categorical variable
- Reconstruct some part of the input (i.e. a word) based on another word and its sense
- Cf. discrete autoencoders

[Marcheggiani and Titov, 2016, Ammar et al., 2014]

MODEL STRUCTURE

- **Encoding:** $p(s|x_i, C_i, C'_i, \theta)$
 - learn a sense mapping with a log-linear model
 - choose the sense of the pivot x_i using the combination of L1 (C_i) and L2 (C'_i) contexts
- **Reconstruction:** $p(x_j|x_i, s, \theta)$
 - learn sense-specific word embeddings with Skip-gram
 - predict a context word x_j based on the pivot and its sense

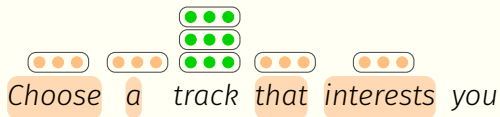
MODEL STRUCTURE


- **Encoding:** $p(s|x_i, C_i, C'_i, \theta)$
 - learn a sense mapping with a log-linear model
 - choose the sense of the pivot x_i using the combination of L1 (C_i) and L2 (C'_i) contexts
- **Reconstruction:** $p(x_j|x_i, s, \theta)$
 - learn sense-specific word embeddings with Skip-gram
 - predict a context word x_j based on the pivot and its sense
- Both components jointly optimized
- Induce a sense mapping that facilitates inferring context words


EN-ES ILLUSTRATION

Encoding:

sense selection $p(s|x_i, C_i, C'_i, \theta)$



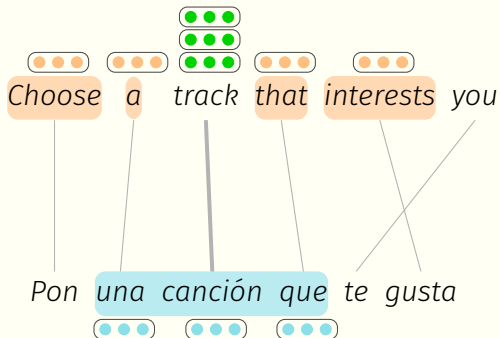
 — sense-specific vector

 — L1 generic vector

EN-ES ILLUSTRATION

Encoding:

sense selection $p(s|x_i, C_i, C'_i, \theta)$



●●● — sense-specific vector

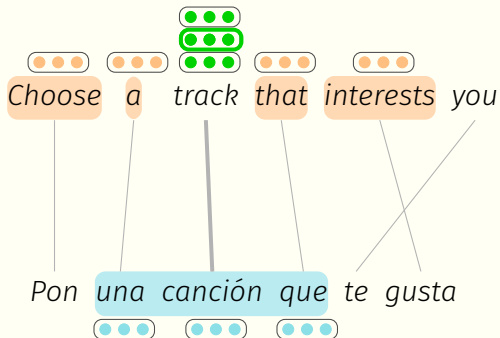
●●● — L1 generic vector

●●● — L2 generic vector

EN-ES ILLUSTRATION

Encoding:

sense selection $p(s|x_i, C_i, C'_i, \theta)$



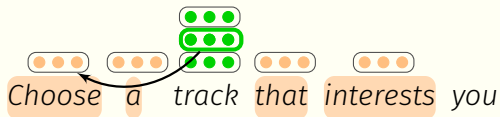
●●● (green) — sense-specific vector

●●● (orange) — L1 generic vector

●●● (blue) — L2 generic vector

EN-ES ILLUSTRATION

Reconstruction:
context-word prediction $p(x_j|x_i, \mathbf{s}, \theta)$



MODEL SET-UP

- **BIMU**: Multi-sense ($n=3$) trained with bilingual signal
- **MU**: Multi-sense trained monolingually
- **SG**: Skipgram

In training of the multi-sense models, use

- entropy regularization to sharpen the encoder posteriors
- or hard updates

EXPERIMENTAL SET-UP

Data

- Use word-aligned parallel corpora in training
- L1: English, paired with
 - French (GigaFrEn)
 - Czech (CzEng)
 - Russian (Yandex)
 - Es, De, Cs, Fr, Ru (NewsCommentary)
- At test time, use only L1 since L2 not available

EXPERIMENTAL SET-UP

Data

- Use word-aligned parallel corpora in training
- L1: English, paired with
 - French (GigaFrEn)
 - Czech (CzEng)
 - Russian (Yandex)
 - Es, De, Cs, Fr, Ru (NewsCommentary)
- At test time, use only L1 since L2 not available

Tasks

	Context?	Representation
Sem. similarity: SCWS	✓	weighted avg.
Sem. similarity: 12 benchmarks	✗	uniform avg.
Qvec	✗	uniform avg.
Neural POS tagger	✓	weighted avg.

NEAREST NEIGHBORS*

≈ 'to follow'

track₁

monitor

keep

analyze

validate

check

manage

evaluate

assess

≈ sports

track₂

jumping

race

scramble

flight

sledge

cycling

rowing

itches

≈ 'railroad line'

track₃

railroad

deck

mainline

wye

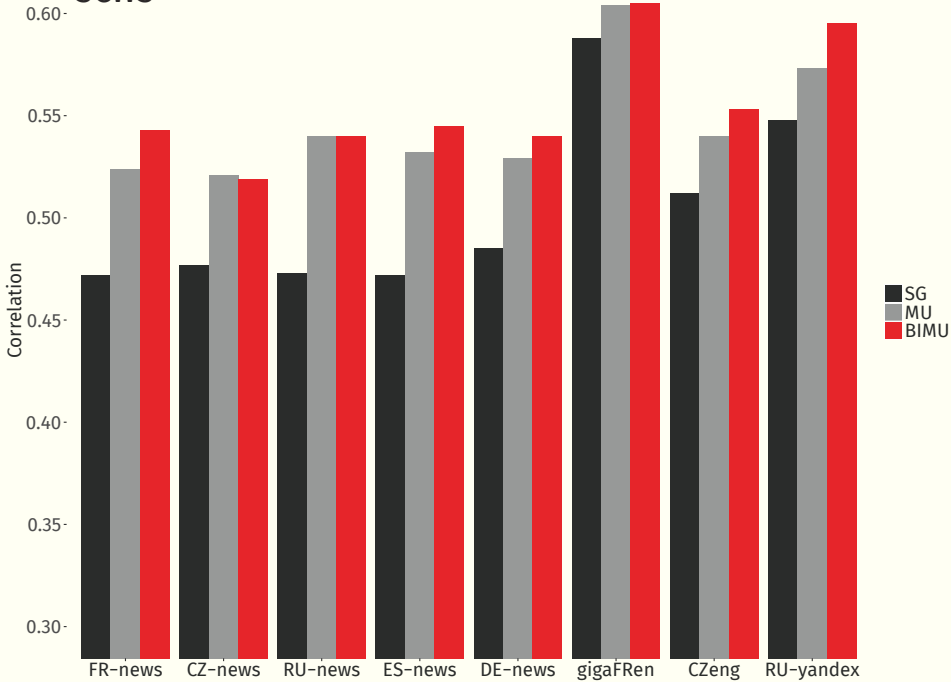
fence

gate

rail

sidings

SCWS



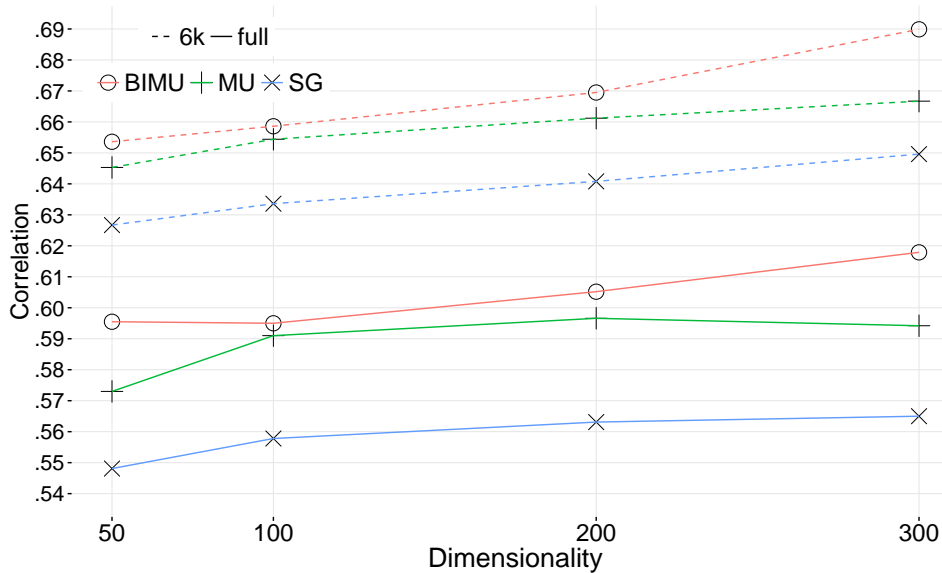
RESULTS ON OTHER EVALUATION TASKS

- Semantic similarity & Qvec:
 - despite the lack of context, large improvements over MU and SG with Russian
 - uniformly averaged sense embeddings might represent rare senses better than SG
- POS tagging:
 - bilingual signal somewhat beneficial
 - overall multi-sense models less robust compared to SG
 - NN might be disentangling the senses in SG embeddings, cf. [Li and Jurafsky, 2015]

OTHER FINDINGS

- Are word alignments necessary?
 - very large L2 context windows (=entire sentence) work well too
- Corpus domain matters
 - e.g. embeddings trained on Yandex (23M) almost as good those trained on GigaFrEn (670M)
- Robustness to increased dimensionality

EFFECT OF EMBEDDING DIMENSIONALITY (SCWS)



Recap

- Bilingual learning affects monolingual quality of English embeddings positively
- Bilingual signal not used at test time
- Some benefits even without word alignments

Recap

- Bilingual learning affects monolingual quality of English embeddings positively
- Bilingual signal not used at test time
- Some benefits even without word alignments

Future directions

- Effect of language choice and language distance
- Bilingual → multilingual signal

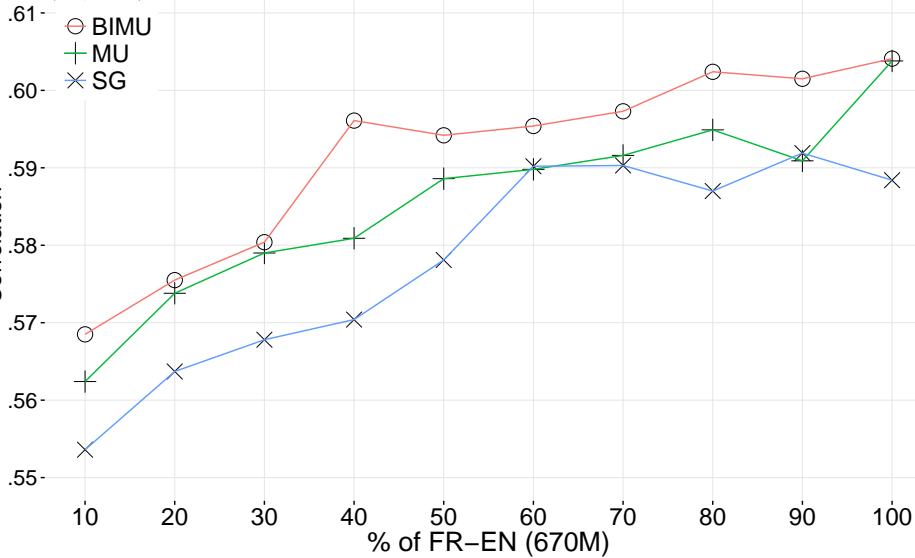
THANK YOU!

Bilingual Learning of Multi-sense Embeddings with Discrete Autoencoders

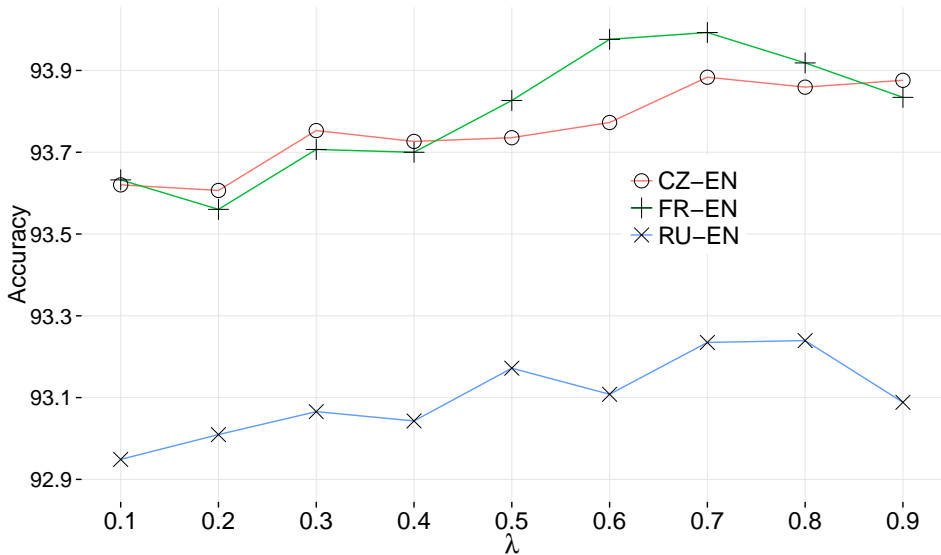
Simon Šuster, Ivan Titov and Gertjan van Noord

Code: github.com/rug-compling/bimu

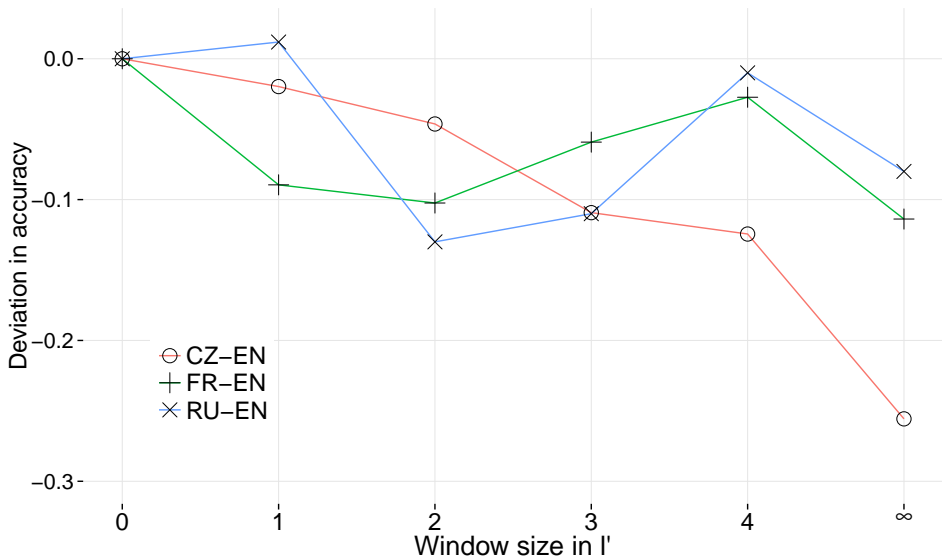
EFFECT OF AMOUNT OF DATA (SCWS)



WEIGHTING L2 (POS TAGGING)



EFFECT OF L2 WINDOW SIZE (POS TAGGING)



Model	RU-EN	CZ-EN	FR-EN
Mu	.63	.59	.64
BiMu, $m = \infty$.66	.62	.64






Table : Comparison of SCWS correlation scores of BiMu trained with infinite l' window to the Mu baseline (vocabulary of top-6000 words).

Model (300-dim.)	SCWS
SG	.65
MU	.66
BIMU	.69
[Chen et al., 2014]	.68
[Neelakantan et al., 2014]	.69
[Li and Jurafsky, 2015]	.70






Table : Comparison to other works (reprinted), for the vocabulary of top-6000 words. Our models are trained on RU-EN, a much smaller corpus than those used in previous work.

Task	Corpus	SG	MU	BIMU
Similarity	RU-EN	37.8	41.2	46.3
	CZ-EN	39.5	36.9	41.9
	FR-EN	46.3	42.0	43.5
	FR-EN (NC)	17.9	26.0	27.6
	RU-EN (NC)	19.3	27.3	28.4
	CZ-EN (NC)	15.8	26.6	25.4
	DE-EN (NC)	20.7	28.4	30.8
	ES-EN (NC)	19.9	27.2	31.2
Qvec	RU-EN	55.8	56.0	56.5
	CZ-EN	56.6	56.5	55.9
	FR-EN	57.5	57.1	57.6
POS	RU-EN	93.5	93.2	93.3
	CZ-EN	94.0	93.7	94.0
	FR-EN	94.1	93.8	94.0

Table : Results, per-row best in bold. SG and MU are trained on the English part of the parallel corpora. In $BIMU-SG$, we report the difference between BIMU and SG, together with the 95% CI of that difference. The Similarity scores are averaged over 12 benchmarks. For POS tagging, we report the accuracy.

-  Ammar, W., Dyer, C., and Smith, N. A. (2014).
Conditional random field autoencoders for unsupervised structured prediction.
In *NIPS*.
-  Chen, X., Liu, Z., and Sun, M. (2014).
A unified model for word sense representation and disambiguation.
In *EMNLP*.
-  Diab, M. and Resnik, P. (2002).
An unsupervised method for word sense tagging using parallel corpora.
In *ACL*.
-  Ettinger, A., Resnik, P., and Carpuat, M. (2016).
Retrofitting sense-specific word vectors using parallel text.
In *NAACL-HLT*.
-  Faruqui, M. and Dyer, C. (2014).
Improving vector space word representations using multilingual correlation.

In *EACL*.

-  Gouws, S., Bengio, Y., and Corrado, G. (2014).
BilBOWA: Fast Bilingual Distributed Representations without
Word Alignments.
arXiv preprint arXiv:1410.2455.
-  Guo, J., Che, W., Wang, H., and Liu, T. (2014).
Learning sense-specific word embeddings by exploiting
bilingual resources.
In *COLING*.
-  Hill, F., Cho, K., Jean, S., Devin, C., and Bengio, Y. (2014).
Embedding word similarity with neural machine translation.
arXiv preprint arXiv:1412.6448.
-  Klementiev, A., Titov, I., and Bhattarai, B. (2012).
Inducing crosslingual distributed representations of words.
In *COLING*.
-  Li, J. and Jurafsky, D. (2015).
Do multi-sense embeddings improve natural language
understanding?

In *EMNLP*.



Marcheggiani, D. and Titov, I. (2016).

Discrete-state variational autoencoders for joint discovery and factorization of relations.

Transactions of the Association for Computational Linguistics, 4.



Neelakantan, A., Shankar, J., Passos, A., and McCallum, A. (2014).

Efficient non-parametric estimation of multiple embeddings per word in vector space.

In *EMNLP*.



Snyder, B. and Barzilay, R. (2010).

Climbing the Tower of Babel: Unsupervised Multilingual Learning.

In *ICML*.