



UPPSALA
UNIVERSITET

Cross-lingual Word Embeddings beyond Zero-shot Machine Translation

Shifei Chen & Ali Basirat

Department of Linguistics and Philology
Uppsala University

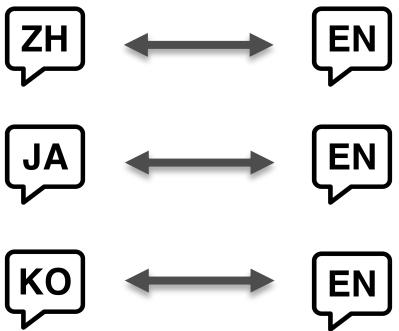
November 27, 2020



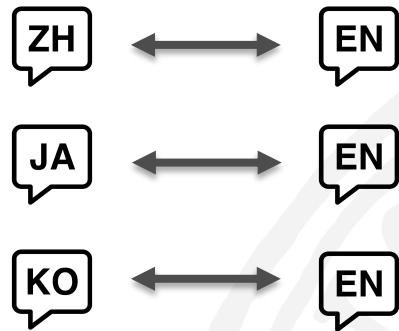
UPPSALA
UNIVERSITET

Introduction

Training



Testing



SV is *never* in the training set
Will it work?





UPPSALA
UNIVERSITET

Research Question

- Is there any *transferability* of a multilingual NMT system from seen languages to *completely unseen* languages?
- Follow up: How *language similarity* works in this scenario?

NMT & Multilingual NMT

- **NMT (Neural Machine Translation)**
 - Uses neural networks to learn the translation relationship between a source and a target language.
 - Comes with an encoder, a decoder and an attention module
- **Multilingual NMT (Johnson et al., 2017)**
 - Same architecture as NMT
 - Additional token indicating the target language
 - Enables transfer learning between languages
- **Even zero-shot translation (Zoph et al., 2016; Nguyen and Chiang, 2017)!**

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
Toan Q. Nguyen and David Chiang. 2017.

Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.

Transfer learning across low-resource, related languages for neural machine translation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 296–301, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Zero-shot Machine Translation

- Translation between unseen language *pairs*
- Benefit low-resource languages
 - By leverage knowledge from high-resource languages
 - No need for a pivot language, or *interlingua*
 - Requires much less data to train
- Our experiment differs
 - Test on *completely unseen* languages

Cross-lingual Word Embeddings

Relates to language similarity (Qi et al., 2018)



“Who” can transfer?

Similar vocabulary distribution exists across languages (Mikolov et al., 2013)



Which part in multilingual NMT is responsible for transferability?

Embedding layers are critical (Kim et al., 2019)



UPPSALA
UNIVERSITET

An Encoder-Decoder LSTM neural network with attention module

Neural Network

TED subtitle corpus (Qi
et al., 2018)

Training: varying from
490k to 1m sentences

Test: varying from 9k to
28k

Corpus

Training: EN+DE+FR

Testing: SV/HE/HU

Languages

fastText(Joulin et al.,
2018; Bojanowski et al.,
2016) aligned word
embeddings

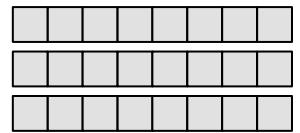
Word Embeddings

Methodology



UPPSALA
UNIVERSITET

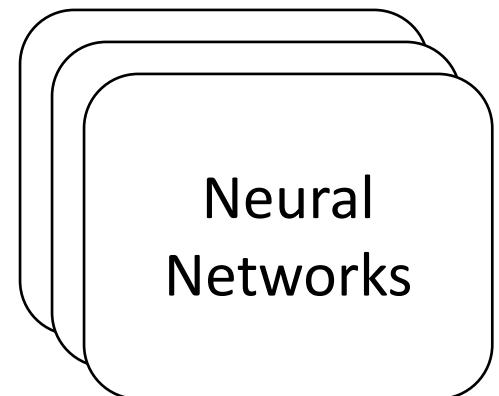
Methodology



Embeddings in the
training languages

EN
DE
FR

Training



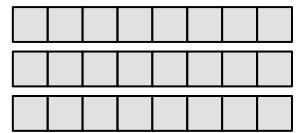
Neural
Networks

SV
EN
DE
FR



UPPSALA
UNIVERSITET

Methodology



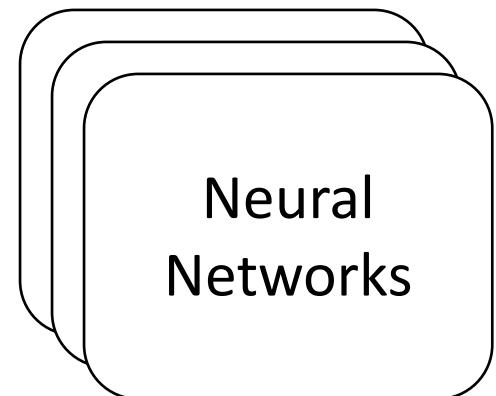
Embeddings in the
training languages



Embeddings in the
unseen language

EN
DE
FR

Training →

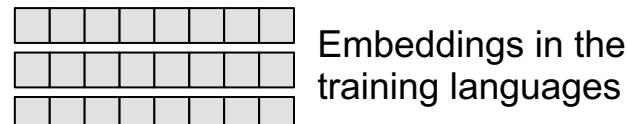


SV
EN
DE
FR



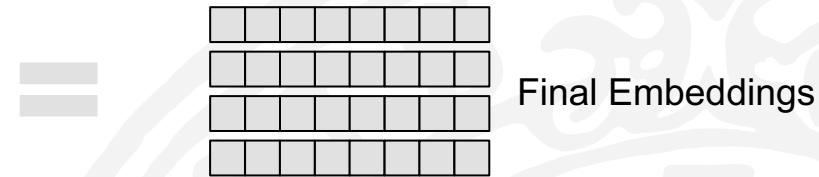
UPPSALA
UNIVERSITET

Methodology



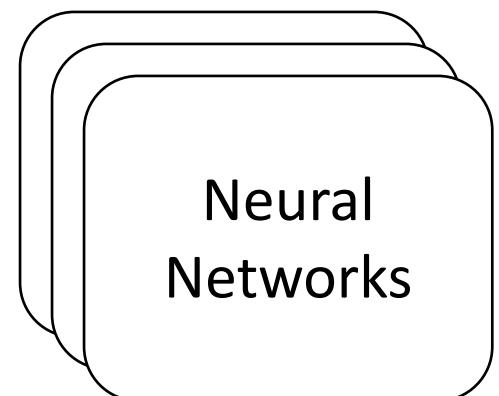
+

Embeddings in the unseen language



EN
DE
FR

Training



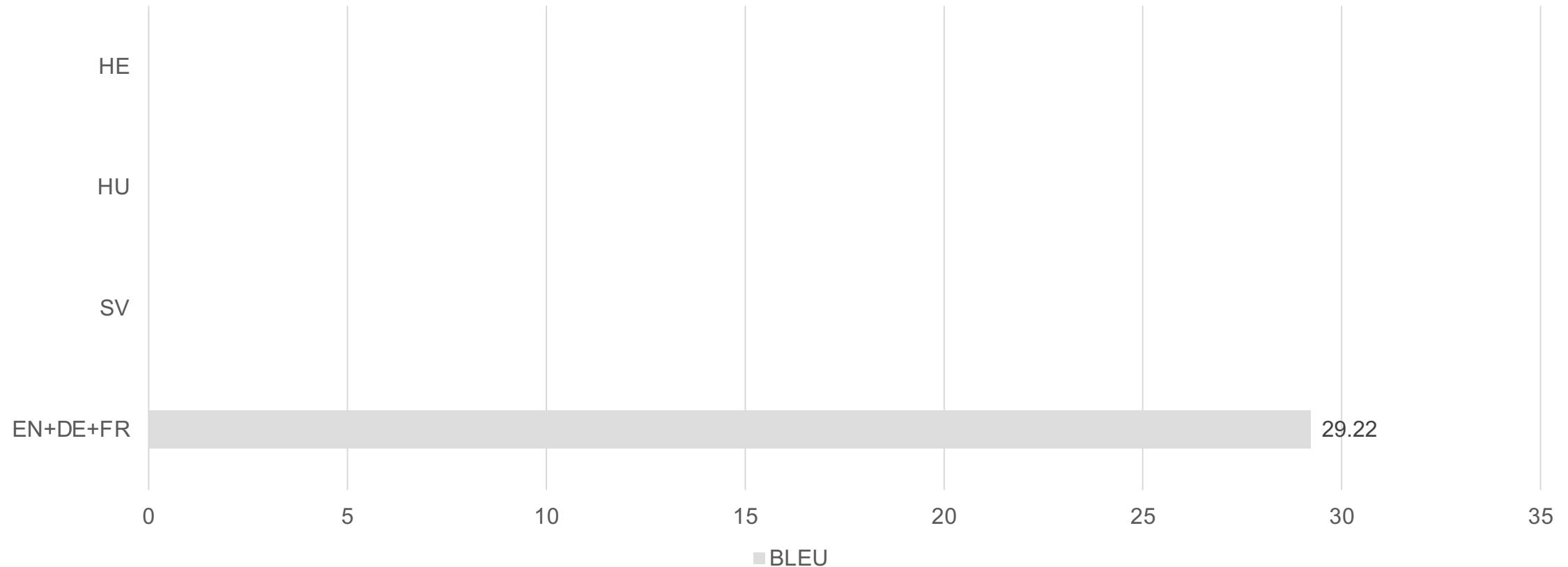
Testing

SV
EN
DE
FR



UPPSALA
UNIVERSITET

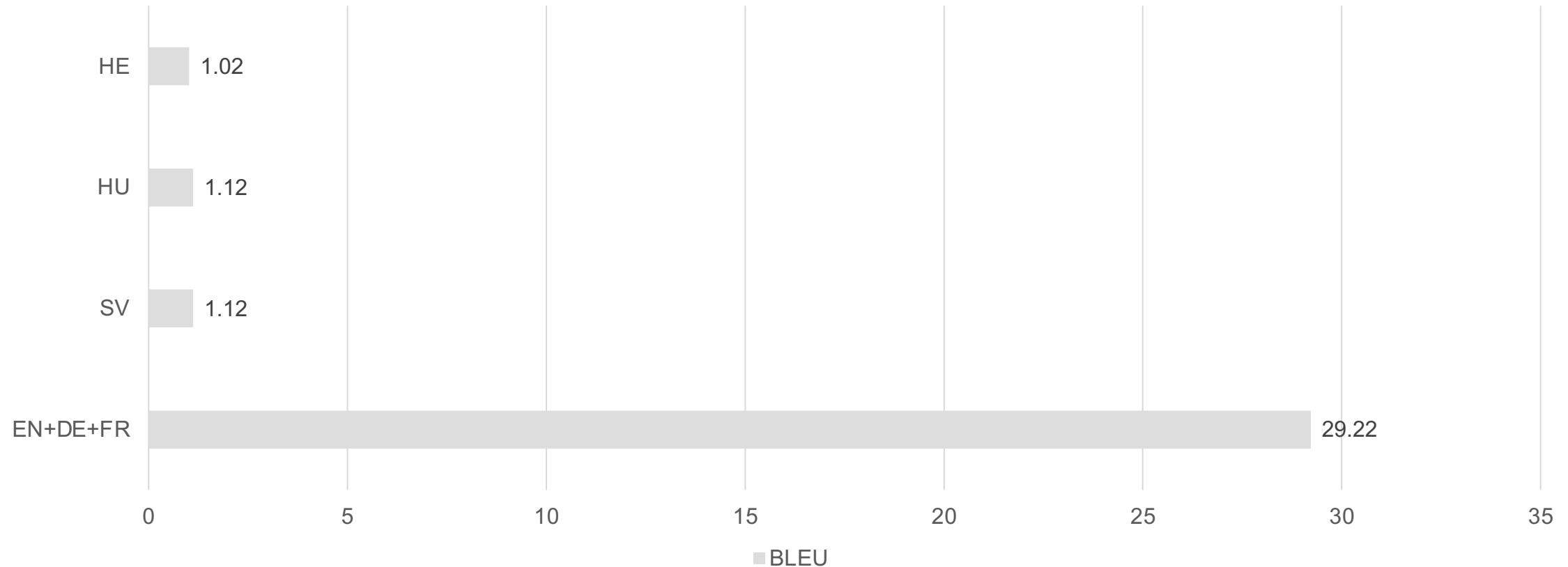
Initial Results





UPPSALA
UNIVERSITET

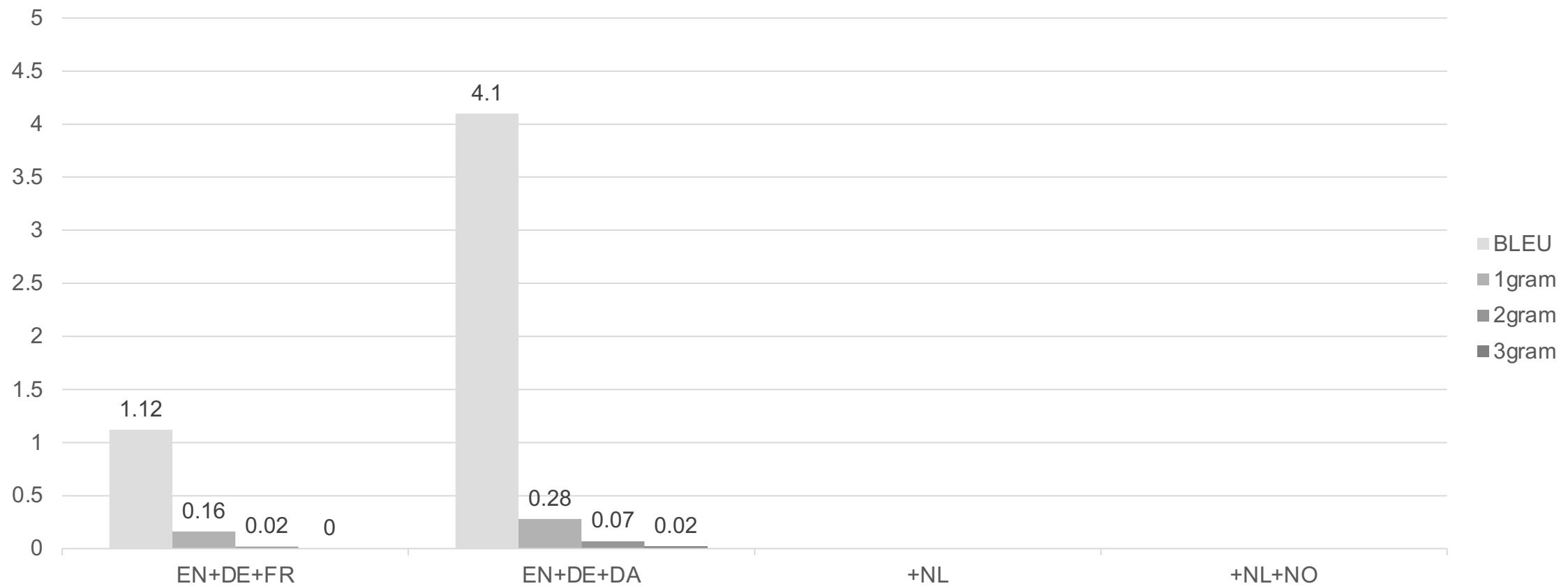
Initial Results





Language Similarity

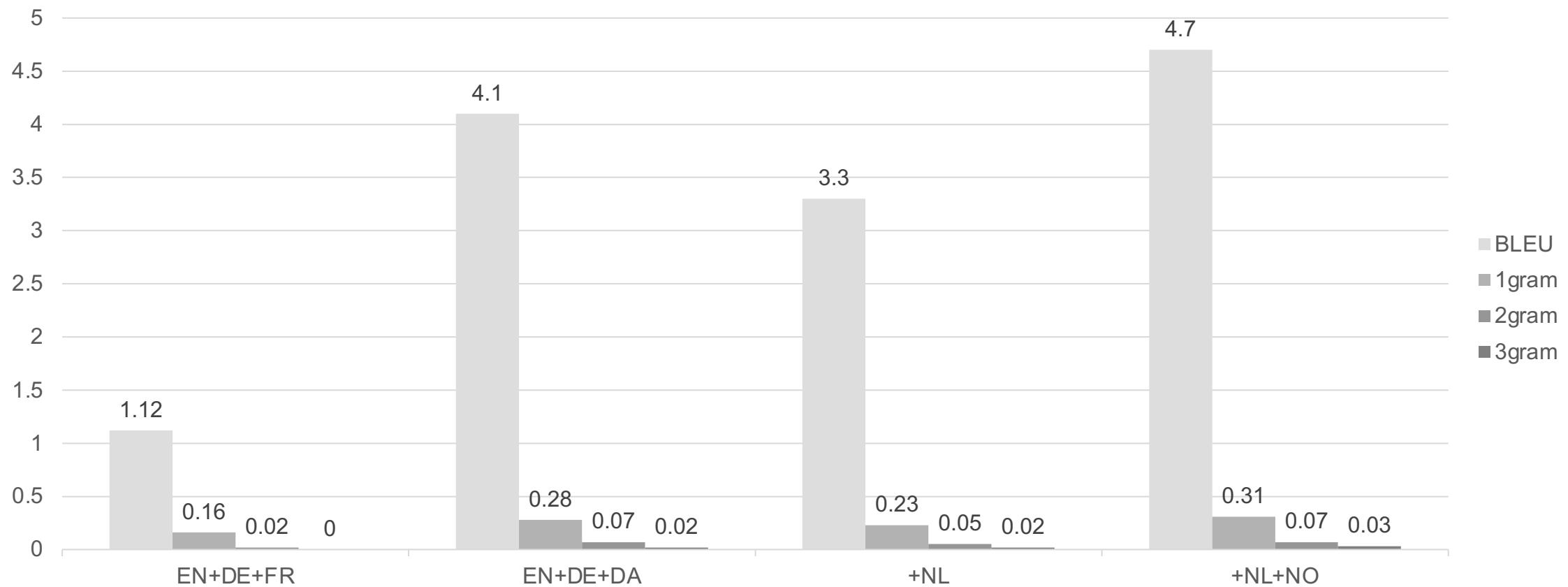
- Remove FR and replace it with DA = better training set homogenization
- Later add NL and NO one by one to the training set





Language Similarity

- Remove FR and replace it with DA = better training set homogenization
- Later add NL and NO one by one to the training set



Source of Language Similarity

Differentiate every token by its origin

__de__ <<sv>>och <<sv>>vi <<sv>>kämpar <<sv>>med <<sv>>dem .



BLEU drops: from 4.1 to 1.7



The system mainly learns lexicon translation

Improvements came from shared vocabularies





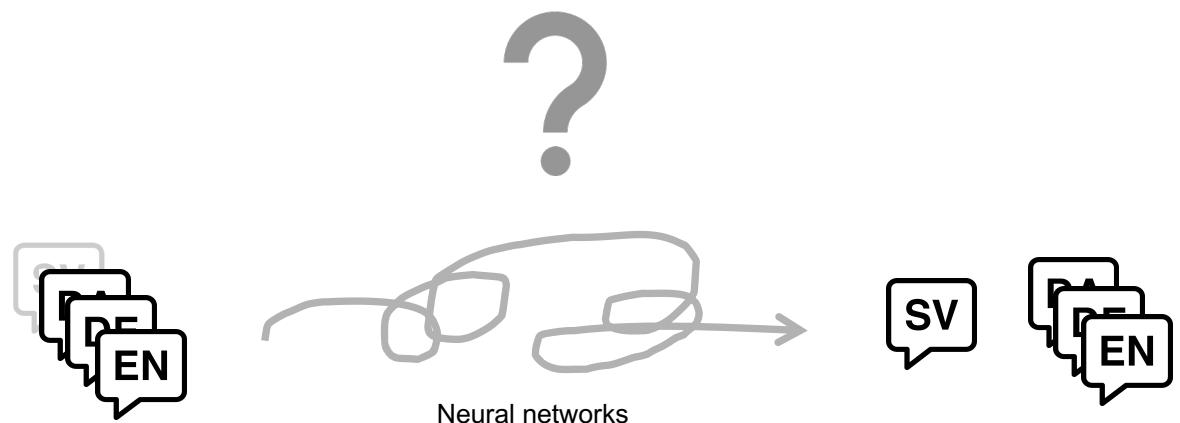
UPPSALA
UNIVERSITET

Gap in Translation Directions

Translation quality differs in different direction:

SV to EN+DE+DA = 6 BLEU scores

EN+DE+DA to SV = 0.65 BLEU score



Transformed Vector Space

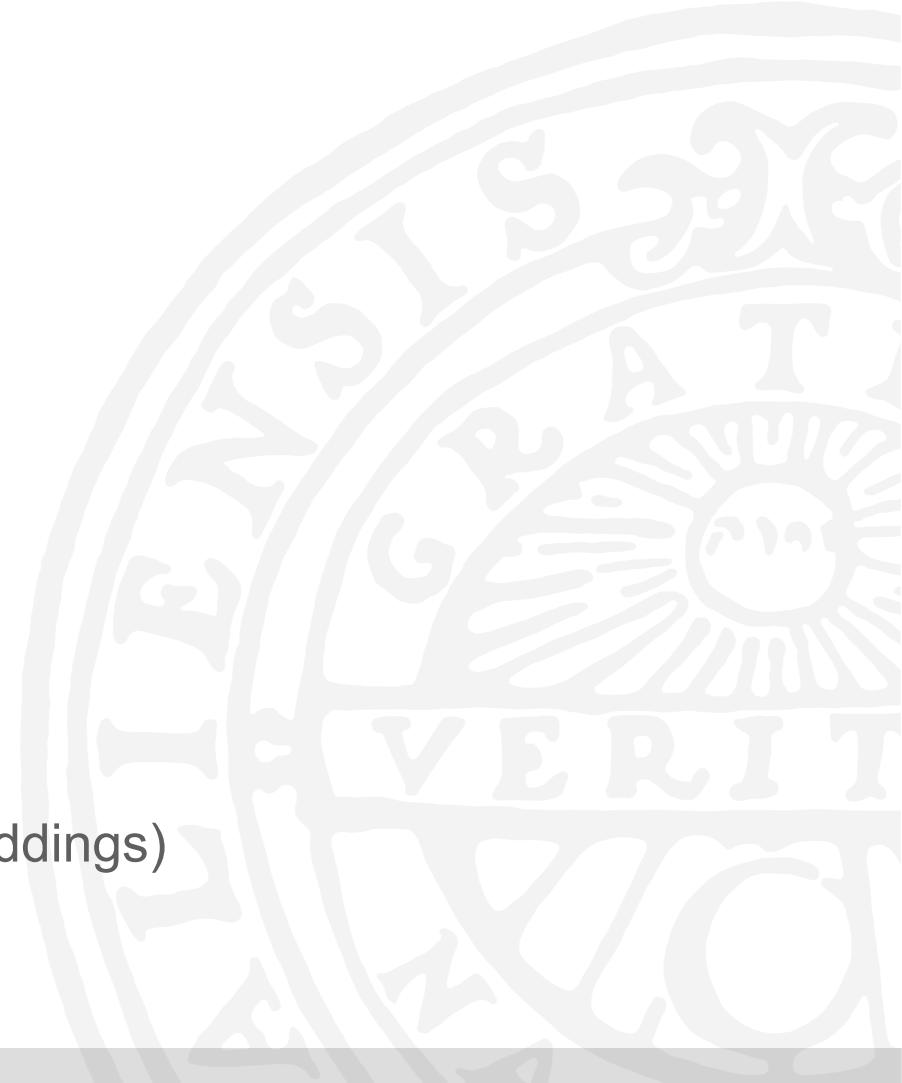
The system *never* sees positive examples



Output vector weight gets *continuously* deducted



Decoder's V_{out} may *no longer* align with the V_{in} (word embeddings)





Conclusion

- Transferability from unseen language to seen language exists
- Language similarity is related to the transferability because of similar vocabularies

We have also observed:

- No positive examples caused the transformed output vector space
- Cross-lingual word embedding alone is not enough for transfer learning (Aji et al., 2020)



Future Work

- Add a regularization layer to the loss function
- Add language information to the model, e.g.,
 - Language embeddings (Littell et al., 2017; Malaviya et al., 2017)
- Explore other NMT architecture, e.g., Transformer (Vaswani et al., 2017)
- Explore other embeddings, e.g.,
 - contextual word embeddings (Devlin et al., 2019)
 - sub-word embeddings (Heinzerling and Strube, 2017)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need.

Littell, P., Mortensen, D. R., Lin, K., Kairis, K., Turner, C., and Levin, L. URIEL and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers (Valencia, Spain, Apr. 2017), Association for Computational Linguistics, pp. 8–14.

Malaviya, C., Neubig, G., and Littell, P. Learning language representations for typology prediction. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (Copenhagen, Denmark, Sept. 2017), Association for Computational Linguistics, pp. 2529–2535.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) (Minneapolis, Minnesota, June 2019), Association for Computational Linguistics, pp. 4171–4186.

Heinzerling, B., and Strube, M. Bpemb: Tokenization-free pre-trained subword embeddings in 275 languages.