

# CUNI Machine Translation Systems for the Covid-19 MLIA Initiative

Ivana Kvapilíková and Ondřej Bojar

Charles University, Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics  
Malostranské náměstí 25, 118 00 Prague, Czech Republic  
<surname>@ufal.mff.cuni.cz

**Abstract.** We participated in the MT task of the Covid-19 MLIA initiative and submitted MT systems translating from English to German, Greek, French, Italian, Italian and Swedish. After the first round, we observe that an efficient training strategy is to use transfer learning to leverage in-domain training data in other languages (both from related and unrelated language families).

## 1 Introduction

A global crisis such as the current Covid-19 pandemic requires information to be spread as efficiently as possible. Working with information from different international resources in multiple languages can resolve possible inconsistencies and prevent misinformation. In an emergency situation, new data are released constantly and are communicated to the general public via national news or government statements, as well as international reports and scientific journals. There are extensive data resources written in English which are not accessible for non-English speakers. In order to quickly access the information in a foreign language, machine translation (MT) can be of great help. However, Covid-related texts are a part of a specific domain and MT models are known to struggle outside of the general domain.

The machine translation task of the MLIA @ Eval initiative consists of translating from English to German, Greek, French, Italian, Italian and Swedish. Trained MT systems can be used to make Covid-related texts accessible to speakers of these six languages. Our team participates in all six language tracks.

## 2 Data

### 2.1 Round 1

In Round 1 we only submitted constrained systems trained on the data provided by the task organizers. The data for Round 1 are summarized in Table 1. The

---

Copyright © 2020-2021 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

Covid-19 MLIA @ Eval Initiative, <http://eval.covid19-mlia.eu/>.

development test set used for final model selection was obtained by cutting 500 sentences off of either the train set or the development set, depending on the original development set size.

	<b>de</b>	<b>el</b>	<b>es</b>	<b>fr</b>	<b>it</b>	<b>sv</b>
Train	925,647	834,240	1,028,287	1,004,215	900,472	806,425
Dev	528	3,378	1,973	728	3,245	723
Dev Test	500	500	500	500	500	500
Blind Test	2,000	2,000	2,000	2,000	2,000	2,000

Table 1: Data Summary

### 3 Methodology

#### 3.1 Round 1

In Round 1 we experimented with three training approaches:

1. standard NMT training with back-translation (*BASE*);
2. transfer learning (*TRANSFER*);
3. multilingual training (*MULTILING*).

The first approach relies on one bidirectional model (sharing the encoder and decoder for both translation directions) which constantly switches between the training and the inference mode to produce batches of synthetic sentence pairs and learn from both authentic and synthetic training samples using online back-translation (BT) [4]. The models are trained on BPE units [5] with a vocabulary of 30k items.

The second transfer learning approach was proposed by Kocmi and Bojar [3] who fine-tune a low-resource child model from a pre-trained high-resource parent model for a different language pair. The method requires a shared subword vocabulary generated from the concatenation of corpora of both the child and the parent language pair. The training procedure consists of first training an NMT model on the parent parallel corpus until it converges, then replace the training data with the child corpus. We experiment with repeating this procedure several times with the child becoming the parent for either a completely new language (e.g. German  $\rightarrow$  English  $\rightarrow$  Spanish  $\rightarrow$  ...) or for the original parent (e.g. German  $\rightarrow$  English  $\rightarrow$  German  $\rightarrow$  ...). When adding a new language, the joint BPE vocabulary has to be modified by replacing the original parent vocabulary entries with the new child's.

The multilingual approach uses the same architecture as described above and trains one MT model to translate from English into three languages (French, Italian and Spanish). During inference, the target language is determined from

indicated language embeddings of the target sentence. We selected these three languages for their similarity which could help the model reuse and share some knowledge. The BPE vocabulary was extracted from the concatenation of all four corpora, using only unique English sentences to reach a comparable corpus size.

For all our MT models we use a 6-layer Transformer [7] architecture with 8 heads, embedding dimension of 1024 and GELU [1] activations. The training is performed using the XLM<sup>1</sup> toolkit. The translation models were trained on 4 GPU<sup>2</sup> with 2-step gradient accumulation to reach an effective batch size of  $8 \times 3400$  tokens. Effective batch size has a significant impact on the training and we observe that the models converge on lower BLEU scores for smaller batch sizes. We used Adam [2] optimizer with inverse square root decay ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ ,  $lr = 0.0001$ ). Beam search with the beam size of 4 was used during final decoding; greedy decoding was used for back-translation. The vocabulary size was set to 30k. Using larger vocabulary leads to a performance drop.

## 4 Results

### 4.1 Round 1

For each language pair we trained a bidirectional back-translation model described in Section 3 and compared it to a standard unidirectional model without back-translation. We experimented with a dropout of 0.1 and 0.2 and concluded that higher dropout helps in most settings. This observation is in line with Senrich and Zhang [6] who emphasize the role of higher dropout when working with low- to medium- sized resources. The results are summarized in Table 2. We submitted the best model for each language pair under the name *BASE*.

	de	el	es	fr	it	sv
bidirectional w\BT	<b>21.52</b>	22.30	<b>40.94</b>	<b>38.46</b>	<b>33.17</b>	<b>20.61</b>
unidirectional w\o BT	20.76	<b>22.70</b>	40.46	35.57	30.97	19.13

Table 2: Translating from English using the *BASE* models: BLEU scores on dev set.

We used the best-performing *BASE* models as the parent models and continued with unidirectional training (foreign language  $\rightarrow$  English) for our transfer learning experiments. To our surprise, it often helped to use the transfer several times, having the model converge on one parallel corpus, switch the target language, wait for convergence and switch again. For example transferring from

<sup>1</sup> <https://github.com/facebookresearch/XLM>

<sup>2</sup> Quadro P5000, 16GB of RAM

German to Spanish to Italian (32.10 BLEU) performs better than transferring directly from Spanish to Italian (31.68 BLEU). The best combination is to even repeat the Spanish-Italian transfer twice (33.07 BLEU).

When translating from English to German, fine-tuning the en-de *BASE* model on English→Spanish (or English→Swedish) and switching back to English → German adds around 1 BLEU on top of the original *BASE* model. The language combinations used in our transfer learning experiments are described in Table 3.

We observe that transfer learning improves the performance in all cases but French, where the *BASE* model with BT reaches 38.46 BLEU, which is  $\sim 3$  BLEU points more than transfer learning.

Transfer Combination	de	el	es	fr	it	sv
en-es → en-de	21.26					
en-de → en-es			41.28			
en-de → en-es → en-de	<b>22.60</b>					
en-de → en-es → en-fr				<b>35.10</b>		
en-de → en-es → en-it					32.10	
en-de → en-es → en-it → en-es			<b>41.34</b>			
en-de → en-es → en-it → en-es → en-it					<b>33.07</b>	
en-es → en-fr				32.43		
en-es → en-it					31.68	
en-de → en-el		<b>23.29</b>				
en-es → en-el		20.91				
en-de → en-sv						<b>21.69</b>
en-de → en-sv → en-de	22.55					
en-de → en-sv → en-de → en-sv						20.56
en-de → en-sv → en-de → en-sv → en-de	22.50					

Table 3: Translating from English using the *TRANSFER* models: BLEU scores on dev set.

Table 4 shows the comparison of the *TRANSFER* models with a multilingual model trained jointly for French, Italian and Spanish. We observe that transfer learning is a more effective way to leverage multilingual data than joint multilingual training. However, there is an advantage of a joint model in terms of the training and storage cost. After three days of training, the multilingual model can be used for translation into all three languages. The initial *BASE* models can take between one (without BT) and five (with BT) days to train and fine-tuning on a child language pair adds around 6 hours.

Table 5 lists our task submissions and compares all approaches on the official blind test test.

	de	el	es	fr	it	sv
multilingual	-	-	40.15	36.07	32.76	-
best transfer	<b>22.60</b>	<b>23.29</b>	<b>41.34</b>	35.10	33.07	<b>21.69</b>
best base	21.52	22.7	40.94	<b>38.46</b>	<b>33.17</b>	20.61

Table 4: Translating from English using the best *BASE*, *TRANSFER* and *MULTILING* models: BLEU scores on dev set.

	de	el	es	fr	it	sv
multilingual	-	-	47.3	48.0	<b>28.3</b>	-
best transfer	<b>31.6</b>	<b>24.7</b>	<b>47.9</b>	47.1	<b>28.3</b>	<b>30.1</b>
best base	31.4	24.1	47.3	<b>48.4</b>	-	28.5

Table 5: Translating from English using the best *BASE*, *TRANSFER* and *MULTILING* models: BLEU scores on blind test set.

## 5 Conclusion

We experimented with three training approaches and conclude that there is not a universal winner that would defeat the other models in all language directions. However, transfer learning brings promising results across the board. In this setting, transferring knowledge is a more effective way to leverage multilingual data than joint training. For English→German, we observe that a transfer learning detour via Spanish or Swedish improves the parent model itself. For English→Greek, transfer learning via German works well, despite the unrelatedness of the two languages. For English→French on the other hand, a bidirectional model with back-translation beats the *TRANSFER* models. In the following rounds we would like to continue analyzing the transfer combinations on the final translation quality. We would also like to train an unconstrained model using a large pretrained model from the general domain.

## Acknowledgments

This study was supported in parts by the grants 19-26934X (Ivana Kvapilíková) and 18-24210S (Ondřej Bojar) of the Czech Science Foundation, SVV 260 575 and GAUK 1050119 of the Charles University. This work has been using language resources and tools stored and distributed by the LINDAT/CLARIN project of the Ministry of Education, Youth and Sports of the Czech Republic (LM2018101).

## References

- [1] Hendrycks, D., Gimpel, K.: Bridging nonlinearities and stochastic regularizers with gaussian error linear units. arXiv [e-Print archive] [abs/1606.08415](https://arxiv.org/abs/1606.08415) (2017), URL <https://arxiv.org/abs/1606.08415>
- [2] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Proceedings of the 3rd International Conference for Learning Representations (2015), URL <http://arxiv.org/abs/1412.6980>
- [3] Kocmi, T., Bojar, O.: Trivial transfer learning for low-resource neural machine translation. In: Proceedings of the Third Conference on Machine Translation: Research Papers, pp. 244–252, Association for Computational Linguistics, Brussels (Oct 2018), <https://doi.org/10.18653/v1/W18-6325>, URL <https://www.aclweb.org/anthology/W18-6325>
- [4] Lample, G., Denoyer, L., Ranzato, M.: Unsupervised machine translation using monolingual corpora only. In: 6th International Conference on Learning Representations (ICLR 2018) (2018), URL <http://arxiv.org/abs/1711.00043>
- [5] Sennrich, R., Haddow, B., Birch, A.: Neural machine translation of rare words with subword units. In: Proceedings of the 54th Annual Meeting of the ACL, pp. 1715–1725, Association for Computational Linguistics, Berlin (Aug 2016), <https://doi.org/10.18653/v1/P16-1162>, URL <https://www.aclweb.org/anthology/P16-1162>
- [6] Sennrich, R., Zhang, B.: Revisiting low-resource neural machine translation: A case study. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 211–221, Association for Computational Linguistics, Florence, Italy (Jul 2019), <https://doi.org/10.18653/v1/P19-1021>, URL <https://www.aclweb.org/anthology/P19-1021>
- [7] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30, pp. 6000–6010, Curran Associates, Inc. (2017), URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>