

# Breaking Down Multilingual Machine Translation

Anonymous ACL submission

## Abstract

While multilingual training is now an essential ingredient in machine translation (MT) systems, recent work has demonstrated that it has different effects in different multilingual settings, such as many-to-one, one-to-many, and many-to-many learning. These training settings expose the encoder and the decoder in a machine translation model with different data distributions. In this paper, we examine how different varieties of multilingual training contribute to learning these two components of the MT model. Specifically, we compare bilingual models with encoders and/or decoders initialized by multilingual training. We show that multilingual training is beneficial to encoders in general, while it only benefits decoders for low-resource languages (LRLs). We further find the important attention heads for each language pair and compare their correlations during inference. Our analysis sheds light on how multilingual translation models work and also enables us to propose methods to improve performance by training with highly related languages. Our many-to-one models for high-resource languages and one-to-many models for LRL outperform the best results reported by Aharoni et al. (2019).<sup>1</sup>

## 1 Introduction

Multilingual training regimens (Dong et al., 2015; Firat et al., 2016; Ha et al., 2016) are now a key element of natural language processing, especially for low-resource languages (LRLs) (Neubig and Hu, 2018; Aharoni et al., 2019). These algorithms are presumed to be helpful because they leverage syntactic or semantic similarities between languages, and transfer processing abilities across language boundaries.

In general, English is used as a central language due to its data availability, and three different multilingual training settings are considered: (1) *one-to-many*: training a model with languages pairs from

English to many other languages. (2) *many-to-one*: training a model with languages pairs from many languages to English (3) *many-to-many*: training a model with the union of the above two settings' data. (1) and (3) can be used for English to other (En-X) translation, while (2) and (3) can be used for other to English (X-En) translation.

However, multilingual training has not proven equally helpful in every setting. Arivazhagan et al. (2019) showed that many-to-one training improves performance over bilingual baselines more than one-to-many does. In this paper we consider this result from the point of view of the components of the MT model. In the many-to-one setting, inputs of the model are from different language distributions so the encoder can be considered a multi-domain model, whereas the decoder is trained on a single distribution. In the one-to-many setting, it is the opposite: the encoder shares data, and the decoder is multi-domain. While there are recent studies analyzing multilingual translation models (Kudugunta et al., 2019; Voita et al., 2019a; Aji et al., 2020; Mueller et al., 2020), in general they do not (1) examine the impact of different multilingual training settings such as one-to-many and many-to-one, and (2) they do not examine the different components such as encoder and the decoder separately.

This motivates us to ask “*how do various types of multilingual training interact with learning of the encoder and decoder?*” To answer this question, we set up controlled experiments that decouple the contribution to the encoder and the decoder in various training settings. We first train multilingual models using many-to-one, one-to-many, or many-to-many training paradigms. We then compare training bilingual models with and without initializing the encoder or the decoder with parameters learnt by multilingual training. We find that, for LRLs, multilingual training is beneficial to both the encoder and the decoder. However, surprisingly, for high-resource languages (HRL), we found mul-

<sup>1</sup>We will release our scripts once accepted.

Lang.	az	be	gl	sk	ar	de	he	it
Size (K)	6	5	10	61	214	168	212	205

Table 1: Training data size.

083 tilingual training only beneficial to encoder but not  
084 to the decoder.

085 To further analyze the result, we examine "*to*  
086 *what degree are the learnt parameters shared*  
087 *across languages?*". We use the head importance  
088 estimation method proposed by Michel et al. (2019)  
089 as a tool to identify the important attention heads  
090 in the model, and measure the consistency between  
091 the heads sets that are important for different lan-  
092 guage pairs. The results suggest that the encoder  
093 does share parameters across different languages  
094 in all settings. On the other hand, the decoder  
095 can treat the representation from the encoder in a  
096 language-agnostic way for X-En translation, and  
097 less parameter sharing is observed for En-X trans-  
098 lation. Our analyses on parameter sharing also  
099 provides a possible explanation of Kudugunta et al.  
100 (2019)'s observation that the representation from  
101 the encoder is target-language-dependent .

102 Our investigation of how multilingual training  
103 works leads us to a method for improving MT mod-  
104 els. With the comprehensive experiments in mul-  
105 tilingual settings, for translation in HRL (Ar-En,  
106 De-En, He-En, It-En), we discover that fine-tuning  
107 multilingual model with target bilingual data out-  
108 performs the best results in Aharoni et al. (2019)  
109 by 2.99 to 4.63 BLEU score . With the analy-  
110 sis on the parameter sharing in the decoder, we  
111 are able to identify related languages. Fine-tuning  
112 jointly with the identified related languages boosts  
113 low-resource translation (En-Az, En-Be, En-Go,  
114 En-Sk) over the best results in Aharoni et al. (2019)  
115 by 1.66 to 4.44 BLEU score. Compared to Neubig  
116 and Hu (2018), our method does not require lin-  
117 guist knowledge, and thus may be more useful for  
118 less-studied low-resource languages.

119 In sum, our contributions are in three-fold. First,  
120 our experiments can be used as a diagnostic tool  
121 for multilingual translation to investigate how an  
122 encoder and a decoder benefit from multilingual  
123 training. Second, our results provide insights into  
124 how multilingual translation works. Third, we im-  
125 prove the translation models based on the findings  
126 from our analysis, showing a promising path for fu-  
127 ture research on multilingual machine translation.

## 2 Experimental Settings for Multilingual Training

128 Before stepping into our analysis, we first explain  
129 our experimental setup. The publicly available  
130 TED Talks Dataset (Qi et al., 2018) is used to train  
131 all our machine translation models. Following Neu-  
132 big and Hu (2018), we break words into subwords  
133 with BPE jointly learnt over all source languages  
134 using the `sentencepiece` toolkit. The vocabu-  
135 lary size is 32000. We perform experiments with  
136 the Transformer architecture (Vaswani et al., 2017)  
137 using the hyper parameters same as in (Arivazha-  
138 gan et al., 2019)<sup>2</sup>. All models are implemented  
139 and trained using *Fairseq* 0.10.0 (Ott et al., 2019).  
140 We trained multilingual translation models with 60  
141 different languages on the TED Talks Dataset with  
142 the three settings described in Section 1: *one-to-*  
143 *many*, *many-to-one* and *many-to-many*. For *one-to-*  
144 *many* and *many-to-many* settings, we add a special  
145 language token to the input of the encoder to indi-  
146 cate the target language. Following Aharoni et al.  
147 (2019), we evaluate our models with BLEU score  
148 (Papineni et al., 2002; Post, 2018) on the selected  
149 8 languages. They are representative for different  
150 language families (Qi et al., 2018). The size of the  
151 training is shown in Table 1.

## 3 How Multilingual Training Benefits Each Component

152 Previous studies have shown that the multilingual  
153 training results are generally stronger than the bilin-  
154 gual training (Arivazhagan et al., 2019). To under-  
155 stand how multilingual training benefits NMT, we  
156 analyze the effect of multilingual training on dif-  
157 ferent components of an NMT model, specifically,  
158 the encoder and decoder.

### 3.1 Experiments Design

159 To study how multilingual training benefits each  
160 component, we train models on bilingual data with  
161 components initialized differently as follows:  
162

- 163 • **Bilingual Only:** Models trained from scratch  
164 with no components initialized with param-  
165 eters learnt from multilingual training.
- 166 • **Load encoder/decoder:** Models with train-  
167 able parameters of either encoder or decoder  
168

<sup>2</sup>6 layers in both the encoder and the decoder, 8 atten-  
tion head, state dimension=512, ffn dimension=2048, label  
smoothing=0.1

Model		→ en							
		az	be	gl	sk	ar	de	he	it
	All-En	9.1	15.2	27.4	25.4	23.9	28.3	27.9	31.5
	All-All	8.1	12.6	22.8	24.6	21.7	27.1	26.1	31.1
Bilingual Only		2.1	1.4	2.8	18.5	28.5	32.0	34.8	35.7
All-En	Load Enc.	2.8	1.8	5.9	18.1	30.6	35.5	36.9	35.7
	Load Dec.	2.5	1.8	5.7	17.8	27.2	30.3	33.2	35.7
	Freeze Enc.	5.0	6.0	19.3	26.3	28.4	33.0	33.6	36.4
	Freeze Dec.	3.4	4.1	16.9	24.7	28.1	31.4	33.4	33.6
	Load Both	<b>11.5</b>	19.0	29.9	28.00	30.4	33.1	36.2	36.7
All-All	Load Enc.	5.4	7.0	20.6	28.0	30.9	35.7	37.1	38.1
	Load Dec.	1.4	0.5	0.9	20.4	28.9	32.2	34.0	35.3
	Freeze Enc.	3.3	5.0	9.3	23.8	25.9	32.4	32.2	34.2
	Freeze Dec.	2.0	6.2	20.1	26.9	30.1	34.4	35.9	36.8
	Load Both	11.3	<b>19.4</b>	<b>31.8</b>	<b>29.6</b>	<b>31.3</b>	<b>36.0</b>	<b>37.8</b>	<b>38.7</b>

Table 2: Results of translating into English. **All** in the model name refers to using all 60 languages.

initialized with parameters learnt from multi-lingual training.

- **Load both:** Models with parameters of both encoder and decoder initialized with parameters learnt from multilingual training. This can be seen as fine-tuning the multilingual model on bilingual data.

The motivation for this paradigm is that if multilingual training is beneficial to a component, then initializing the parameters of that component should result in improvements over random initialization and training on only bilingual data. If *load encoder* outperforms *bilingual only*, then we can say that multilingual training is beneficial for the encoder, and if *load decoder* outperforms we can make the analogous conclusion for the decoder. Thus comparing these models reveals how each component benefit from multilingual training.

We also consider a *load and freeze* setting (Thompson et al., 2018), where we initialize a component from a multilingual model and freeze its weights when fine-tuning on bilingual data. For example, in the *load decoder* setting, we train the loaded decoder with a randomly initialized encoder. We suspect that learning with randomly initialized component might ruin the other component which is well-trained with multilingual data, especially in the beginning of the training. Thus, we additionally experiment with this *load and freeze* setting to ensure the multilingual-trained component is not

deteriorated.

## 3.2 Results and Discussion

The overall results of X-En and En-X are shown in Table 2 and Table 3, respectively. Because they are highly dependent on the training data size (Table 1), we discuss the results in two groups: HRL (HRL; referring to *ar*, *de*, *he*, and *it*) and LRL (LRL; referring to *az*, *be*, *gl*, *sk*).<sup>3</sup>

### 3.2.1 Low-Resource Language Results

For LRLs, we find that multilingual training is generally beneficial to both the encoders and the decoders in all of the three multilingual models. Both *load encoder* and *load and freeze decoder* can achieve performance better than the bilingual baseline. This suggests that the parameters in the encoder and the decoder learnt by multilingual training do contain information that is not effectively learnt from the smaller bilingual data.

The results also suggest that multilingual training is more beneficial for the encoders than for the decoders. In all cases, either *load encoder* or *freeze encoder* outperforms both *load decoder* and *load and freeze decoder*. However, multilingual training of the encoder and the decoder are complementary; loading both the encoder and the decoder can usually improve the performance over loading only one component.

<sup>3</sup>*sk* has intermediate size, and its behavior is not always consistent with the other LRL.

Model	en→								
	az	be	gl	sk	ar	de	he	it	
En-All	4.9	9.0	24.2	21.9	15.1	27.9	24.1	33.3	
All-All	3.1	6.2	20.5	18.4	12.7	24.5	21.1	30.5	
Bilingual Baseline	1.3	1.9	3.9	13.1	15.6	27.1	25.4	32.0	
En-All	Load Enc.	3.0	5.6	16.7	21.7	<b>17.2</b>	30.0	27.5	34.6
	Load Dec.	1.3	2.0	8.1	17.4	16.0	26.7	25.8	32.6
	Freeze Enc.	2.7	4.6	14.7	21.1	9.7	24.4	22.6	33.4
	Freeze Dec.	1.9	3.7	14.5	17.6	16.2	28.0	25.9	33.3
	Load All	<b>6.4</b>	<b>14.7</b>	<b>26.9</b>	<b>23.5</b>	17.1	<b>31.1</b>	<b>28.2</b>	<b>34.9</b>
All-All	Load Enc.	2.4	5.0	16.9	21.4	16.9	29.8	27.4	34.4
	Load Dec.	1.1	2.2	7.0	17.5	16.0	28.1	25.6	32.5
	Freeze Enc.	2.1	0.5	12.6	19.4	10.2	24.4	24.3	33.1
	Freeze Dec.	0.9	4.7	15.0	18.8	15.1	27.5	24.9	32.4
	Load All	6.1	13.0	26.4	23.2	17.0	30.3	27.9	34.6

Table 3: Results of translating from English. **All** in the model name refers to using all 60 languages.

### 3.2.2 High-Resource Language Results

On HRLs, we find that multilingual training is generally beneficial to the encoders in all of the three multilingual models, while it is not beneficial for the decoders in some settings. *Load encoder* always outperform the baseline models, but for the All-En model on X-En translation, and the All-All model on En-X translation, neither *load decoder* nor *load and freeze decoder* outperform the baseline model.

We also observe that multilingual training is generally more beneficial to the encoders than to the decoders. In all of the cases, *load encoder* can achieve performance competitive to *load both* (better or less by within 1 BLEU score). However, in all of the cases, both *load decoder* and *load and freeze decoder* have performance worse than *load both*. Therefore, multilingual training is not as beneficial to the decoders as to the encoders.

### 3.3 Discussion

For LRL, because the size of bilingual training data is small, it is not surprising that multilingual training is beneficial for both the encoder and the decoder. However, our results are somewhat more surprising for HRL — it is not trivial that multilingual training is not as beneficial. In the next section, we focus on explaining the phenomena observed on HRL by investigating how parameters are shared across languages.

## 4 How Multilingual Parameters are Shared in Each Component

Given the previous results, we are interested in exactly *how* parameters are shared among different language pairs. Given that we are using the Transformer architecture, for which multi-head attention is a fundamental component, we use the attention heads as a proxy to analyze how multilingual models work differently when translating between different languages. Specifically, we analyze our models by identifying the attention heads that are important when translating a language pair. Measuring the consistency between the sets of important attention heads for two language pairs gives us hints on the extent of parameter sharing.

### 4.1 Head Importance Estimation

First, we provide some background on head importance estimation, specifically the method proposed by Michel et al. (2019).

Given a set of multi-head attention modules, each of which can be written as

$$\text{MHAtt}(x) = \sum_{h=1}^{N_h} \xi_h \text{Att}_{W_q^{(h)}, W_k^{(h)}, W_v^{(h)}}(x), \quad (1)$$

where  $N_h$  is the number of attention heads, and  $\xi_h = 1$  for all  $h$ .

The importance of a head can be estimated as

$$\tilde{I}_h = \mathbb{E}_{x \sim X} \left| \frac{\partial L(x)}{\partial \xi_h} \right|. \quad (2)$$



given a loss function  $L$  and input  $X$ . Then, the importance score of each head in an attention module is normalized

$$I_h = \frac{\tilde{I}_h}{\sqrt{\sum_i^{N_h} I_h^2}}. \quad (3)$$

Note that when the input  $X$  is different, the estimated importance score can be different. Therefore, when different language pairs are fed in, the important heads identified can be different. We denote the set of attention head scores estimated on translation from language  $l_a$  to language  $l_b$  as  $H(l_a, l_b)$ . We denote the scores of attention heads in a component by using superscript. For example,  $H^{enc}$  represents the scores of the heads in an encoder.

## 4.2 Measuring Parameter Sharing by Correlation of Head Scores

With the attention head importance scores estimated by Equation 3, we can investigate how parameters are shared across languages. For each of the En-All, All-En, All-All multilingual models, we estimated a set of head-importance scores  $H(l_a, l_b)$  for each language pair  $(l_a, l_b)$  in the training setting. We calculate the head scores with the training loss function (MLE with label smoothing) and 100K randomly sampled sentences in the training set.

To investigate how much parameters are shared by two pairs of languages  $(l_a, l_b)$  and  $(l_c, l_d)$ , we measure the agreement between  $H(l_a, l_b)$  and  $H(l_c, l_d)$ . If a head is important for both of  $(l_a, l_b)$  and  $(l_c, l_d)$ , then important parameters for translating are shared. Thus high agreement suggests high parameter sharing.

To quantify the agreement between two score sets, we use Spearman’s rank correlation (Spearman, 1987). A rank-based correlation metric is used because the importance estimation was originally proposed to order attention heads in a model. Higher correlation implies higher agreement and thus implies higher parameter sharing. For each of the En-All, All-En, All-All models, we calculate the correlation between  $H(l_a, l_b)$  and  $H(l_c, l_d)$  for all language pairs  $(l_a, l_b)$  and  $(l_c, l_d)$  that are used to train the model. The detailed correlation computation process can be found in Appendix A. We plot the correlation matrices of the head scores (included in appendix) and summarize them in Table 4.

Model	Lang. Pair	$H^{enc}$	$H^{dec}$
All-En	X-En	.871 (.086)	.973 (.023)
En-All	En-X	.806 (.153)	.720 (.150)
All-All	X-En	.898 (.073)	.967 (.029)
All-All	En-X	.813 (.126)	.762 (.141)

Table 4: Correlation between the attention head scores when estimated using different language pairs.

## 4.3 How Multilingual Translation Models Share

Results in Table 4 combined with Section 3 provides the insights into how multilingual translation models work with respect to cross-lingual sharing:

**Encoder for En-X:** It is natural that the encoder from En-X likely benefit from multilingual training because it can generate representations tailored for different target languages with shared parameters. En-X is a set of language pairs where the source language is always English. Therefore, if the prepended target language token is ignored, the inputs of the encoders for all pairs in En-X are from one identical distribution. This is in contrast to X-En pairs, where the inputs are in different languages. However, for the encoders, we observe from Table 4 that the average correlation scores of En-X pairs (0.806 and 0.813), are lower than the correlation scores of X-En pairs (0.871 and 0.898). Kudugunta et al. discovers that the representation of the encoder is target-language-dependent. Thus we conjecture that some parameters may be used to generate representation tailored for the target languages. At the same time, since the inputs are from a single distribution (English) for different target languages, a large portion of parameters may still be shareable across target languages. Therefore, in this case, multilingual training is beneficial.

**Encoder for X-En:** For X-En language pairs, the input of the encoder is multilingual, which means the input from different X-En language pairs has distinct distribution. However, the correlation between different source languages is still high. It shows that high parameters sharing in the encoder is possible.

**Decoder for En-X:** The decoders for En-X have the lowest correlation. From the correlation matrix, we do see some parameter sharing between some language pairs. However, larger model capacity might be required for a model to be proficient in

all the languages.

**Decoder for X-En:** The decoder have average correlation as high as 0.973 and 0.967 for All-En and All-All models respectively. This suggests that to decode intermediate representation encoded by the encoder, the decoder use almost the same set of parameters. However, [Kudugunta et al.](#) shows that the representation encoded by the encoder is not language-agnostic. A possible explanation is that the important parameters of the decoder are highly determined by the target output, which is always in English. Therefore, even though the encoder representation is not language-agnostic, it is still difficult to learn parameters reflecting the difference. It suggests why multilingual training does not benefit the decoder in the X-En setting. The set of English sentences is almost the same for all the HRL pairs in the TED Talks dataset, so multilingual training can hardly provide more unique English sentences than bilingual training does. If the decoder is dedicated for generation, multilingual training cannot expose the decoder to more diverse data. Therefore the multilingually trained decoder does not perform better than the bilingual one.

## 5 Improving Translation Based on the Degree of Parameter Sharing

Insights from the previous section provide us with a new way to choose languages for multilingual training. In previous work ([Lin et al., 2019](#)), choosing on languages with similar linguistic properties is a popular practice. However, [Mueller et al. \(2020\)](#) found the effect is highly language-dependent. Sometimes training with similar languages might be worse than training on a set of unrelated languages. Here we otherwise propose an entirely model-driven way to find related languages to improve multilingual translation models. We explore choosing languages where parameters can be better shared.

### 5.1 Improving X-En by Related En-X Pairs

In the All-All model, we notice low parameter sharing between En-X and X-En pairs. The average correlation between  $H^{enc}(En, X)$  and  $H^{enc}(X, En)$  is 0.44 (std: 0.17). The average correlation between  $H^{dec}(En, X)$  and  $H^{dec}(X, En)$  is 0.49 (std: 0.13). It provides a possible explanation why training with both the En-X and the X-En pairs only brings little

improvement over training with only En-X alone or with X-En alone.

The low correlation combined with results in Section 3 motivate us to experiment on improving X-En with related En-X pairs. Section 3 shows that the multilingual decoder has less advantage than the encoder. This may suggest the inefficiency of parameter sharing in the decoder. Therefore we experiment on choosing a set of related languages based on the degree of parameter in the decoder. We choose the language set  $L$  such that for all  $l \in L$ , the average correlation  $\frac{1}{60} \sum_{l_i=1}^{60} \text{Corr}(H^{dec}(En, l), H^{dec}(l_i, En))$  is higher than 0.60.

Results are shown in Table 5. Even though fine-tuning on related languages improves the overall performance, it is not better than fine-tuning on the All-En pairs only. Also, the average correlation between  $H^{dec}(En, l_a)$  and  $H^{dec}(l_b, En)$  is not improved. Our experiment demonstrates the difficulty of sharing parameters between All-En pairs and En-All pairs. We leave this problem for future work.

### 5.2 Improving En-X by Language Clusters

The low correlation between attention head scores of language pairs motivates us to improve the performance of En-X using related language pairs. As shown in Table 4, the decoders have the lowest correlation scores. We conjecture that it is due to the difficulty of sharing parameters between distant languages. Thus, we seek for finding related language sets, in each of which parameters can be shared.

Again, we resort to the attention head importance scores to find the related languages. Our intuition is that related languages would share many parameters in between and training a model on related languages would be helpful. As a sanity check of our idea, we first use t-SNE ([Maaten and Hinton, 2008](#)) to reduce the dimension of head-importance scores  $H(l_a, l_b)$ . We only focus on heads in the decoders, because the correlation score between  $H_{(En, l_c)}$  and  $H_{(En, l_d)}$  is lower in average for the decoders. The result visualized in Figure 1 illustrates that, the distance between  $H_{(En, l_c)}$  and  $H_{(En, l_d)}$  tend to be shorter if languages  $l_c$  and  $l_d$  are linguistically related. Hence, determining related languages with head score  $H_{(En, l)}$  should be reasonable.

We then fine-tune multilingual models on related language clusters. Related languages clusters are determined by k-mean++ ([Arthur and Vassilvitskii,](#)

Model	az	be	gl	sk	ar	de	he	it
All-All	8.1	12.6	22.8	24.6	21.7	27.1	26.1	31.1
+ f.t. on All-En	10.5	17.5	29.7	28.1	25.9	31.3	30.5	34.0
+ f.t. on All-En & related	10.5	17.4	28.3	27.0	25.1	30.0	29.9	32.7

Table 5: Performance of All-All model fine-tuned on All-En pairs and fine-tuned on the union of All-En pairs and related En-All languages.

Model	az	be	gl	sk	ar	de	he	it
Bilingual Baseline	1.3	1.9	3.9	13.1	15.6	27.1	25.4	32.0
All-All	3.1	6.2	20.5	18.4	12.7	24.5	21.1	30.5
All-All w/ f.t. on related clusters	7.9	12.8	27.5	24.9	-	30.2	27.0	35.4
All-All w/ f.t. on random groups	6.9	13.3	22.5	24.3	-	-	27.5	35.2
En-All	4.9	9.00	24.2	21.9	15.1	27.9	24.1	33.3
En-All w/ f.t. on related clusters	7.9	13.9	21.0	26.2	16.7	30.4	27.1	35.4
En-All w/ f.t. on random groups	7.0	13.1	23.1	24.7	-	-	27.6	35.2
Load En-All w/ f.t. on closest	7.8	15.2	28.6					

Table 6: Performance of En-All model without and with fine-tuning on language clusters.

2007) with  $k = 5$ . We consider clusters that cover all of the four low-resource languages. For the All-All model, one of the cluster we consider contains Be, Gl, De, He, It, and the other one contains Az. For the En-All model, we also experiment with two clusters. One includes Ar, De, He, It, and the other includes Az, Be, Gl, Sk. As a baseline, we also experiment with random groups. They are groups generated by randomly splitting the 59 target languages.

The results are shown in Table 6. For both the En-All and the All-All model, except En-Gl, fine-tuning on clusters can improve performance on all the considered language pairs consistently. For LRLs, fine-tuning on related language clusters is also better than fine-tuning on random groups in general. To verify whether this improvement is brought by increased parameter sharing in the decoders, we check the correlation between  $H^{dec}$  after fine-tuning. The results shown in Table 7 shows improvements after fine-tuning on the clusters.

For low-resource language pairs En-Az, En-Be, En-Sk on the En-All model, we notice that only few languages are highly correlated with them (with correlation  $> 0.80$ ). Therefore, we also experiment with fine-tuning the En-All model with only the language pairs with high correlation scores ( $> 0.80$ ) for each of the three pairs, which boosts the performance of En-Be to 15.2 and En-Sk to 28.6.

Model	$H^{dec}$ w/o f.t.	$H^{dec}$ w/ f.t.
All-All	.762 (.141)	.894 (.069)
En-All (HL)	.855 (.066)	.866 (.065)
En-All (LL)	.826 (.096)	.834 (.091)

Table 7: Correlation between the decoder attention head scores when estimated using the language pairs in the cluster. HL and LL represent the cluster that includes HRL and the one that includes LRL respectively.

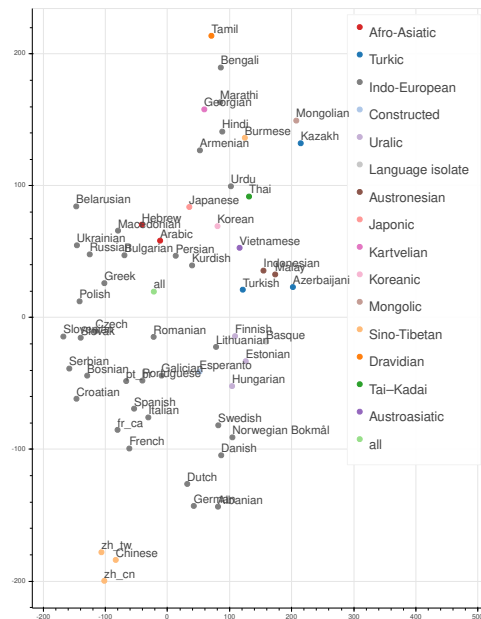


Figure 1: Visualization of the En-All decoder head scores of languages by t-SNE.

## 6 Related Work

The early attempts of multilingual training for machine translation use a single model to translate between multiple languages (Dong et al., 2015; Firat et al., 2016; Ha et al., 2016). Those works find multilingual NMT models are appealing because they not only give us a simple paradigm to handle mapping between multiple languages, but also improve performance on low and zero-resource languages pairs (Gu et al., 2018). However, how multilingual training contributes to components in the translation model still remains unknown.

There are some attempts at analyzing and explaining the translation models. Thompson et al. (2018) analyze the contribution of different components of NMT model to domain adaptation by freezing the weights of components during continued training. Arivazhagan et al. (2019) provide an comprehensive study on the state-of-the-art multilingual NMT model in different training and testing scenarios. Sachan and Neubig (2018) experiment with different parameter sharing strategies in Transformer models, showing that sharing parameters of embedding, key and query performs well for *one-to-many* settings. Artetxe et al. (2020) shows the strong transferability of monolingual representation to different languages. The intermediate representation of BERT can be language-agnostic if we freeze the embeddings during training. The deficiency of the *one-to-many* setting is explored in (Johnson et al., 2017). They find only the *many-to-one* setting consistently improves the performance across languages. Wang et al. (2018) also explore problems of the *one-to-many* setting, and show language-specific components are effective to improve the performance. Voita et al. (2019a) analyzes how generated sentences of NMT models are influenced by context in the encoder and decoder. The attempt to investigate encoder and decoder separately is similar to our work.

Multi-head attention has been shown effective in different NLP tasks. Beyond improving performance, multi-head attention can help with subject-verb agreement (Tang et al., 2018), and some heads are predictive of dependency structures (Raganato and Tiedemann, 2018). Htut et al. (2019) and Clark et al. (2019) report that heads in BERT attend significantly more to words in certain syntactic position. They show some heads seem to specialize in certain types of syntactic relations. Michel et al. (2019), Voita et al. (2019b), and Behnke and

Heafield (2020) study the importance of different attention heads in NMT models, and suggest that we can prune those attention heads which are less important. Brix et al. (2020) also shows pruning NMT models can improve the sparsity level to optimize the memory usage and inference speed.

However, all previous works do not directly investigate how encoder and decoder of NMT models benefit from multilingual training, which is the key question of why multilingual training works. To our best knowledge, we are the first to tackle the question, and our analysis can be used to further improve multilingual NMT models.

## 7 Conclusion

In this work, we have the following findings: 1) In Section 3, we examine how multilingual training contributes to each of the components in a machine translation model. We discover that, while multilingual training is beneficial to the encoders, it is less beneficial to the decoders. 2) In Section 4, our analysis of important attention heads provides insight into the behavior of multilingual components. Results suggest that the encoder in the En-All model may generate target-language-specific representation, while the behavior of the decoder of the All-En model may be source-language-agnostic. In addition, in the All-All model, we observe indications of lower parameter sharing between X-En pairs and En-X pairs. 3) In Section 5, we explore approaches to improve the model based on our findings. On En-X translation, we outperform the best results in (Aharoni et al., 2019). With our proposed analysis as diagnostic tools, future work may further improve the multilingual systems.

Our findings provide some possible future directions. First, parameter sharing between En-X and X-En pairs in the All-All model seems low. Improving the sharing may improve the performance. Second, the decoder in the All-En model seems to behave in a source-language-agnostic way. It may not be optimal since the representation from the encoder is not source-language-agnostic (Kudugunta et al., 2019). To mitigate this issue, either the encoder is required to encode inputs into language-agnostic representation, or the decoder should behave in different ways according to the input representation. We leave the exploration in future work.





708	Paul Michel, Omer Levy, and Graham Neubig. 2019.	Charles Spearman. 1987. The proof and measurement	765
709	Are sixteen heads really better than one? In <i>Ad-</i>	of association between two things. <i>The American</i>	766
710	<i>advances in Neural Information Processing Systems</i> ,	<i>journal of psychology</i> , 100(3/4):441–471.	767
711	pages 14014–14024.		
712	Aaron Mueller, Garrett Nicolai, Arya D. McCarthy,	Gongbo Tang, Mathias Müller, Annette Rios, and Rico	768
713	Dylan Lewis, Winston Wu, and David Yarowsky.	Sennrich. 2018. <b>Why self-attention? a targeted</b>	769
714	2020. <b>An analysis of massively multilingual neural</b>	<b>evaluation of neural machine translation architec-</b>	770
715	<b>machine translation for low-resource languages.</b>	<b>tures.</b> In <i>Proceedings of the 2018 Conference on</i>	771
716	In <i>Proceedings of the 12th Language Resources</i>	<i>Empirical Methods in Natural Language Processing</i> ,	772
717	<i>and Evaluation Conference</i> , pages 3710–3718, Mar-	pages 4263–4272, Brussels, Belgium. Association	773
718	seille, France. European Language Resources Asso-	for Computational Linguistics.	774
719	ciation.		
720	Graham Neubig and Junjie Hu. 2018. <b>Rapid adapta-</b>	Brian Thompson, Huda Khayrallah, Antonios Anasta-	775
721	<b>tion of neural machine translation to new languages.</b>	sopoulos, Arya D. McCarthy, Kevin Duh, Rebecca	776
722	In <i>Proceedings of the 2018 Conference on Empiri-</i>	Marvin, Paul McNamee, Jeremy Gwinnup, Tim An-	777
723	<i>cal Methods in Natural Language Processing</i> , pages	derson, and Philipp Koehn. 2018. <b>Freezing subnet-</b>	778
724	875–880, Brussels, Belgium. Association for Com-	<b>works to analyze domain adaptation in neural ma-</b>	779
725	putational Linguistics.	<b>chine translation.</b> In <i>Proceedings of the Third Con-</i>	780
726	Myle Ott, Sergey Edunov, Alexei Baevski, Angela	<i>ference on Machine Translation: Research Papers</i> ,	781
727	Fan, Sam Gross, Nathan Ng, David Grangier, and	pages 124–132, Brussels, Belgium. Association for	782
728	Michael Auli. 2019. fairseq: A fast, extensible	Computational Linguistics.	783
729	toolkit for sequence modeling. In <i>Proceedings of</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	784
730	<i>NAACL-HLT 2019: Demonstrations.</i>	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	785
731	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Kaiser, and Illia Polosukhin. 2017. Attention is all	786
732	Jing Zhu. 2002. <b>Bleu: a method for automatic eval-</b>	you need. <i>Advances in neural information process-</i>	787
733	<b>uation of machine translation.</b> In <i>Proceedings of</i>	<i>ing systems</i> , 30:5998–6008.	788
734	<i>the 40th Annual Meeting of the Association for Com-</i>	Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-	789
735	<i>putational Linguistics</i> , pages 311–318, Philadelphia,	nrich, and Ivan Titov. 2019a. <b>Analyzing multi-head</b>	790
736	Pennsylvania, USA. Association for Computational	<b>self-attention: Specialized heads do the heavy lift-</b>	791
737	Linguistics.	<b>ing, the rest can be pruned.</b> In <i>Proceedings of the</i>	792
738	Matt Post. 2018. <b>A call for clarity in reporting BLEU</b>	<i>57th Annual Meeting of the Association for Com-</i>	793
739	<b>scores.</b> In <i>Proceedings of the Third Conference on</i>	<i>putational Linguistics</i> , pages 5797–5808, Florence,	794
740	<i>Machine Translation: Research Papers</i> , pages 186–	Italy. Association for Computational Linguistics.	795
741	191, Brussels, Belgium. Association for Computa-	Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-	796
742	tional Linguistics.	nrich, and Ivan Titov. 2019b. <b>Analyzing multi-head</b>	797
743	Ye Qi, Devendra Sachan, Matthieu Felix, Sarguna Pad-	<b>self-attention: Specialized heads do the heavy lift-</b>	798
744	manabhan, and Graham Neubig. 2018. <b>When and</b>	<b>ing, the rest can be pruned.</b> In <i>Proceedings of the</i>	799
745	<b>why are pre-trained word embeddings useful for neu-</b>	<i>57th Annual Meeting of the Association for Com-</i>	800
746	<b>ral machine translation?</b> In <i>Proceedings of the 2018</i>	<i>putational Linguistics</i> , pages 5797–5808, Florence,	801
747	<i>Conference of the North American Chapter of the</i>	Italy. Association for Computational Linguistics.	802
748	<i>Association for Computational Linguistics: Human</i>	Yining Wang, Jiajun Zhang, Feifei Zhai, Jingfang Xu,	803
749	<i>Language Technologies, Volume 2 (Short Papers)</i> ,	and Chengqing Zong. 2018. <b>Three strategies to im-</b>	804
750	pages 529–535, New Orleans, Louisiana. Associa-	<b>prove one-to-many multilingual translation.</b> In <i>Pro-</i>	805
751	tion for Computational Linguistics.	<i>ceedings of the 2018 Conference on Empirical Meth-</i>	806
752	Alessandro Raganato and Jörg Tiedemann. 2018. <b>An</b>	<i>ods in Natural Language Processing</i> , pages 2955–	807
753	<b>analysis of encoder representations in transformer-</b>	2960, Brussels, Belgium. Association for Computa-	808
754	<b>based machine translation.</b> In <i>Proceedings of the</i>	tional Linguistics.	809
755	<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>		
756	<i>and Interpreting Neural Networks for NLP</i> , pages		
757	287–297, Brussels, Belgium. Association for Com-		
758	putational Linguistics.		
759	Devendra Sachan and Graham Neubig. 2018. <b>Parame-</b>		
760	<b>ter sharing methods for multilingual self-attentional</b>		
761	<b>translation models.</b> In <i>Proceedings of the Third Con-</i>		
762	<i>ference on Machine Translation: Research Papers</i> ,		
763	pages 261–271, Brussels, Belgium. Association for		
764	Computational Linguistics.		

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850

## A Correlation of Head Scores

Here we detail the computation of the correlation of head scores for two pairs of languages  $(l_a, l_b)$  and  $(l_c, l_d)$ . The steps are as follow:

1. The the two language pairs' head importance scores  $H(l_a, l_b)$  and  $H(l_c, l_d)$  are estimated with Equation 3. Since there are many heads in a Transformer model, both  $H(l_a, l_b)$  and  $H(l_c, l_d)$  are vectors.
2. We flatten the scores in  $H(l_a, l_b)$  and  $H(l_c, l_d)$  into two arrays of scalars. We treat the two arrays as the observations of two variables. Then we use Spearman correlation to compute the correlation between the two variables. In other words, the input of the Spearman correlation function is the two arrays.

## B Related Related Language Pairs

The related language pairs used in Section 5 are: en-zh\_cn en-it en-es en-vi en-zh\_tw en-nl en-fr en-fr\_ca en-th en-pt\_br en-ru.

## C Language Clusters

En-All model:

- en-ja en-ko en-zh en-zh-cn en-zh-tw
- en-az en-be en-bs en-cs en-da en-eo en-et en-eu en-fi en-gl en-hr en-hu en-lt en-mk en-nb en-pl en-sk en-sl en-sq en-sr en-sv en-tr en-uk
- en-bn en-hi en-hy en-ka en-ku en-mr en-my en-ta en-th en-ur
- en-ar en-bg en-de en-el en-es en-fa en-fr en-fr-ca en-he en-id en-it en-ms en-nl en-pt en-pt-br en-ro en-ru en-vi
- en-kk en-mn

All-All:

- en-be, en-bg, en-bs, en-cs, en-de, en-el, en-es, en-fr, en-fr-ca, en-gl, en-he, en-hr, en-it, en-lt, en-mk, en-pl, en-pt, en-pt-br, en-ro, en-ru, en-sk, en-sl, en-sq, en-sr, en-uk
- en-ar, en-fa, en-ja, en-ko, en-th, en-vi, en-zh, en-zh-cn, en-zh-tw
- en-bn, en-hi, en-hy, en-ka, en-ku, en-mr, en-my, en-ur

- en-az, en-da, en-eo, en-et, en-fi, en-hu, en-id, en-ms, en-nb, en-nl, en-sv, en-tr 851  
852

- en-eu, en-kk, en-mn, en-ta 853

## D Random Clusters 854

- en-pt en-fa en-fr en-kk en-hi en-da en-hu en-de en-nl en-ar en-hy en-zh-cn 855  
856

- en-sr en-fi en-be en-ko en-ru en-ur en-it en-id en-el en-eu en-sq en-zh en-bs en-bn en-sv en-bg en-my en-ro en-ta en-sl en-et en-ku en-mn en-uk en-he en-tr 857  
858  
859  
860

- en-mk en-mr 861

- en-ms en-pl en-pt-br en-cs en-zh-tw en-es 862

- en-vi en-eo en-hr en-nb en-fr-ca en-az en-sk en-ka en-lt en-th en-ja en-gl 863  
864

Theses random clusters are generated by (1) shuffling the 59 languages, (2) randomly selecting positions. The results 5 segments separated by the 4 positions are the 5 clusters. 865  
866  
867  
868

## E Closest Languages 869

The closest languages used in Section ?? are: 870

- Az: en-az en-eu en-fi en-tr 871
- Be: en-be en-it en-uk 872
- Gl: en-gl en-pt en-es en-lt en-it en-pt\_br 873

## F Experimental Details 874

- Infrastructure: All the experiments can be conducted on one single RTX 2080Ti GPU. 875  
876

- Evaluation: We report the BLEU score calculated by FairSeq. 877  
878

- Version of FairSeq: We use v0.10.0 (<https://github.com/pytorch/fairseq/tree/v0.10.0>) 879  
880  
881

- Dataset: It can be downloaded from <https://github.com/neulab/word-embeddings-for-nmt>. 882  
883  
884

Figure 2: Correlation matrix between language pairs. The top-left corner is the correlation between the encoder head scores  $H^{enc}$ , while the bottom-right corner is the correlation between the decoder head scores  $H^{dec}$ . The top matrix is the correlation matrix of the All-All model, while the bottom-left and the bottom-right ones are the correlation matrices of the All-En and the En-All models respectively.

Figure 3: Correlation matrix between language pairs after fine-tuning on related languages. The top-left corner is the correlation between the encoder head scores  $H^{enc}$ , while the bottom-right corner is the correlation between the decoder head scores  $H^{dec}$ .

Figure 4: Correlation matrix between language pairs after fine-tuning on the languages clusters. The first figure is the matrix of the fine-tuned All-All model. The second and the third ones are the matrix of the En-All model fine-tuned on the language clusters containing the high-resource and the LRL respectively. The top-left corner is the correlation between the encoder head scores  $H^{enc}$ , while the bottom-right corner is the correlation between the decoder head scores  $H^{dec}$ .