

Improving Both Domain Robustness and Domain Adaptability in Machine Translation

Wen Lai¹ and Jindřich Libovický² and Alexander Fraser¹

¹ Center for Information and Language Processing, LMU Munich, Germany

² Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic
{lavine, fraser}@cis.lmu.de libovicky@ufal.mff.cuni.cz

Abstract

We consider two problems of NMT domain adaptation using meta-learning. First, we want to reach domain *robustness*, i.e., we want to reach high quality on both domains seen in the training data and unseen domains. Second, we want our systems to be *adaptive*, i.e., making it possible to finetune systems with just hundreds of in-domain parallel sentences. We study the domain adaptability of meta-learning when improving the domain robustness of the model. In this paper, we propose a novel approach, **RML-NMT** (Robust Meta-Learning Framework for Neural Machine Translation Domain Adaptation), which improves the robustness of existing meta-learning models. More specifically, we show how to use a domain classifier in curriculum learning and we integrate the word-level domain mixing model into the meta-learning framework with a balanced sampling strategy. Experiments on English→German and English→Chinese translation show that RMLNMT improves in terms of both domain robustness and domain adaptability in seen and unseen domains¹.

1 Introduction

The success of Neural Machine Translation (NMT; Bahdanau et al., 2015; Vaswani et al., 2017) heavily relies on large-scale high-quality parallel data, which is difficult to obtain in some domains. We study two major problems in NMT domain adaptation. First, models should work well on both seen domains (the domains in the training data) and unseen domains (domains which do not occur in the training data). We call this property *domain robustness*. Second, with just hundreds of in-domain sentences, we want to be able to quickly adapt to a new domain. We call this property *domain adaptability*. Previous work on NMT domain adaptation has usually focused on only one aspect of domain

adaptation at the expense of the other one, and our motivation is to consider both of the two properties.

There are a few works attempting to solve domain adaptability. The most basic approach is *fine-tuning*, in which an out-of-domain model is continually trained on in-domain data (Freitag and Al-Onaizan, 2016; Dakwale and Monz, 2017). Although fine-tuning is effective, it can suffer from so-called catastrophic forgetting (French, 1999), resulting in deteriorated model performance in general domains (Thompson et al., 2019). Another efficient method is *Meta-Learning* (Hospedales et al., 2021), which trains models which can be later rapidly adapted to new scenarios using only a small amount of data. It works for many natural language processing (NLP) tasks (Gu et al., 2018; Qian and Yu, 2019; Yu et al., 2020; Bansal et al., 2020; Wang et al., 2021; Du et al., 2021), especially in low-resource scenarios (Dou et al., 2019; Yin, 2020). As a result, meta-learning is often used for NMT domain adaptation. For example, Sharaf et al. (2020) and Li et al. (2020) fast adapt NMT models to new domains with meta-learning using a small amount of training data. Zhan et al. (2021) improve meta-learning-based NMT models with a curriculum-based (Bengio et al., 2009) sampling strategy. Meta-learning works well for adapting to new domains, however, previous work tends to neglect the problem of robustness towards domains unseen at training time.

Müller et al. (2020) defined the concept of domain robustness and propose to improve the domain robustness by subword regularization (Kudo, 2018), defensive distillation (Papernot et al., 2016), reconstruction (Tu et al., 2017) and neural noisy channel reranking (Yee et al., 2019). Jiang et al. (2020) proposed using individual modules for each domain with a word-level domain mixing strategy, which they showed has domain robustness on seen domains. The work on domain robustness, however, tends to neglect the adaptability of the models

¹Our source code is available at <https://github.com/lavine-lmu/RMLNMT>

for new domains.

To address both domain adaptability and domain robustness at the same time, we propose RMLNMT (robust meta-learning NMT), a more robust meta-learning-based NMT domain adaptation framework. We first train a word-level domain mixing model to improve the robustness on seen domains, and show that, surprisingly, this improves robustness on unseen domains as well. Then, we train a domain classifier based on BERT (Devlin et al., 2019) to score training sentences; the score measures similarity between out-of-domain and general-domain sentences. This score is used to determine a curriculum to improve the meta-learning process. Finally, we improve domain adaptability by integrating the domain-mixing model into a meta-learning framework with the domain classifier using a balanced sampling strategy.

In summary, we make the following contributions: i) we propose RMLNMT, which shows better domain robustness and domain adaptability than all previous baseline systems; ii) we show that unseen domains can be very effectively handled with domain-robust models, even though post-hoc adaptation with domain-specific data still delivers the best overall translation quality; iii) Experiments on English→German and English→Chinese translation tasks show the effectiveness of RMLNMT. To the best of our knowledge, this is the first work that considers both domain adaptability and domain robustness in NMT domain adaptation, a combination which we suggest the community pay more attention to.

2 Preliminaries

Neural Machine Translation. The goal of the NMT model is to model the conditional distribution of translated sentence $y = (y_1, \dots, y_n)$ given a source sentence $x = (x_1, \dots, x_m)$. Current state-of-art NMT models (Transformers; Vaswani et al., 2017) model the multi-head attention mechanism to focus on information in different representation subspaces from different positions

$$\text{MultiHead}(Q, K, V) = \text{Concat}(h_1, \dots, h_h) W^O$$

$$h_i = \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right),$$

where $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d \times d/m}$ and $W^O \in \mathbb{R}^{d \times d}$. For the i -th head h_i , m is the number of heads, and d is the dimension of the model output. In some of our experiments (see Section 3.1),

we modify the multihead attention to do domain mixing (Jiang et al., 2020).

Meta-learning for NMT. The goal of Meta-Learning is training a teacher model that using previous experience can be better finetuned for new tasks, including handling different domains in NMT domain adaptation (Gu et al., 2018; Sharaf et al., 2020; Zhan et al., 2021). The idea of NMT domain adaptation with meta-learning is to use a small set of source tasks $\{\mathcal{T}_1, \dots, \mathcal{T}_n\}$ (which correspond to domains) to find the initialization of model parameters θ from which finetuning for task \mathcal{T}_0 would require only a small number of training examples. These meta-learning algorithms consist of three main steps: (i) split the seen domain corpus into small tasks \mathcal{T} containing a small amount of data as $\mathcal{D}_{\text{meta-train}}$ and $\mathcal{D}_{\text{meta-test}}$ to simulate the low-resource scenarios. Data for each task \mathcal{T}_i is decomposed into two sub-sets: a support set $\mathcal{T}_{\text{support}}$ used for training the model and a query set $\mathcal{T}_{\text{query}}$ used for evaluating the model; (ii) leverage a meta-learning policy to adapt model parameters to different small tasks using $\mathcal{D}_{\text{meta-train}}$ datasets. We use MAML, proposed by Finn et al. (2017), to create adaptable NMT systems which will be useful for different domains; (iii) finetune the model using the support set of $\mathcal{D}_{\text{meta-test}}$.

3 Method

In our initial experiments, we observed that the standard meta-learning approach for NMT domain adaptation sacrifices the domain robustness on seen domains in order to improve the domain adaptability on unseen domains. To address these issues, we propose a novel approach, RMLNMT, which combines meta-learning with a word-level domain-mixing system (for improving domain robustness) in a single model. RMLNMT consists of three parts: Word-Level Domain Mixing, Domain Classification, and Online Meta-Learning. Figure 1 illustrates RMLNMT.

3.1 Word-level Domain Mixing

In order to improve the robustness of NMT domain adaptation, we follow the approach of Jiang et al. (2020) and train a word-level layer-wise domain mixing NMT model.

Domain Proportion. From a sentence-level perspective (i.e., the classifier-based curriculum step), each sentence has a domain label. However, the

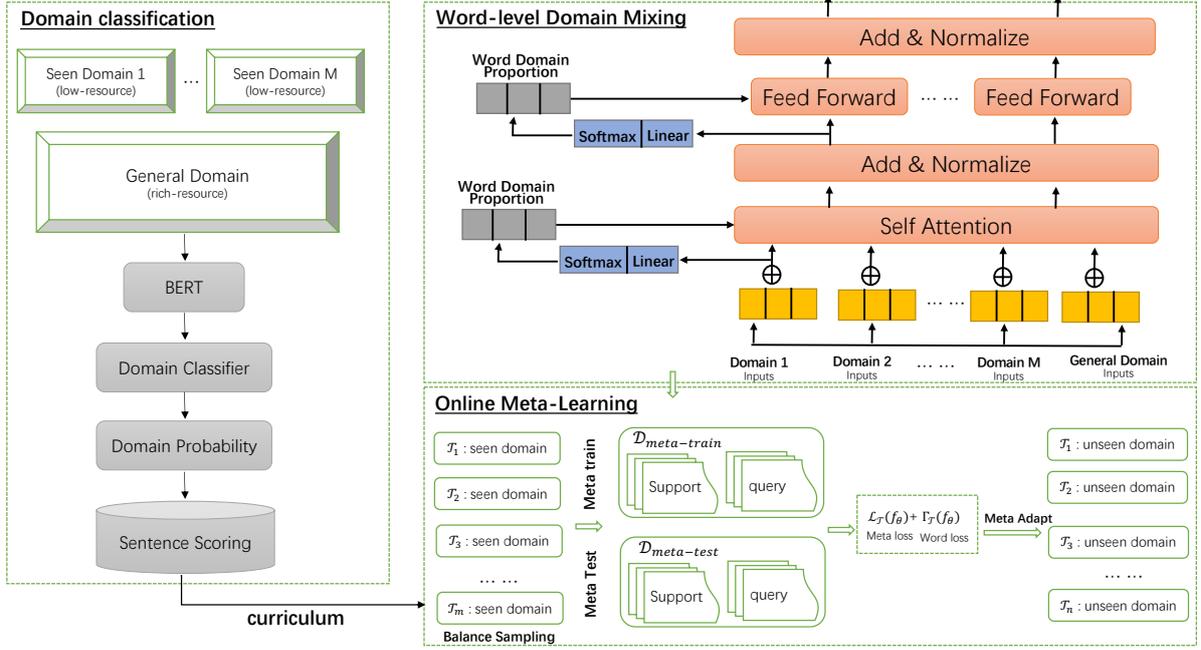


Figure 1: Method overview. The whole procedure mainly consists of three parts: domain classification, word-level domain mixing and online meta-learning.

domain of a word in the sentence is not necessarily consistent with the sentence domain. E.g., the word *doctor* can have a different meaning in the medical domain and the academic domain. More specifically, for k domains, the embedding $\mathbf{w} \in \mathbb{R}^d$ of a word, and a matrix $R \in \mathbb{R}^{k \times d}$, the domain proportion of the word is represented by a smoothed softmax function as:

$$\Phi(\mathbf{w}) = (1 - \epsilon) \cdot \text{softmax}(R\mathbf{w}) + \epsilon/k,$$

where $\epsilon \in (0, 1)$ is a smoothing parameter to prevent the output of $\Phi(\mathbf{w})$ from collapsing towards 0 or 1.

Domain Mixing. Following Jiang et al. (2020), each domain has its own multi-head attention modules. Therefore, we can integrate the domain proportion of each word into its multi-head attention module. Specifically, we take the weighted average of the linear transformation based on the domain proportion Φ . For example, we consider the point-wise linear transformation $\{W_{i,V,j}\}_{j=1}^k$ on the t -th word of the input, V_t , of all domains. The mixed linear transformation can be written as

$$\bar{V}_{i,t} = \sum_{j=1}^k V_t^\top W_{i,V,j} \Phi_{V,j}(V_t),$$

where $\Phi_{V,j}(V_t)$ denotes the j -th entry of $\Phi_V(V_t)$, and Φ_V is the domain proportion layer related to

V . For other linear transformations, we apply the domain mixing scheme in the same way for all attention layers and the fully-connected layers.

Training. The model can be efficiently trained by minimizing a composite loss:

$$L^* = L_{\text{gen}}(\theta) + L_{\text{mix}}(\theta),$$

where θ contains the parameter in encoder, decoder and domain proportion. $L_{\text{gen}}(\theta)$ denotes the cross-entropy loss over training data $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ and $L_{\text{mix}}(\theta)$ denotes the cross-entropy loss over the words/domain labels. For $L_{\text{mix}}(\theta)$, we compute the cross-entropy loss of its domain proportion $\Phi(\mathbf{w})$ as $-\log(\Phi_J(\mathbf{w}))$, which take J as the domain label. Hence, $L_{\text{mix}}(\theta)$ is computed as the sum of the cross-entropy loss over all such pairs of word labels of the training data.

3.2 Domain Classification

Domain similarity has been successfully applied in NMT domain adaptation. Moore and Lewis (2010) calculate cross-entropy scores with a language model to represent the domain similarity. Rieß et al. (2021) leverage simple classifiers to compute similarity scores; these scores are more effective than scores from language models for NMT domain adaptation. Motivated by Rieß et al. (2021),

we compute domain similarity using a sentence-level classifier, but in contrast with previous work, we based our classifier on a pre-trained language model. Given k domain corpora (one general domain corpus and n out-of-domain corpora), we trained a sentence classification model M based on BERT (Devlin et al., 2019). For a sentence x with a domain label L_x , a simple softmax is added to the top of the model M to predict the domain probability of sentence x :

$$P(x | h) = \text{softmax}(Wh),$$

where W is the parameter matrix of M and h is the hidden state of M . $P(x | h)$ is a probability set, which contains k probability scores indicating the similarity of sentence x to each domain. We finally select the probability of the general domain (from k probability scores) as the score of the sentence x and use this score as the curriculum to split the task in meta-learning (see more details in Section 3.3). A higher score indicates that the sentence is more similar to the general domain, so we will select it earlier.

3.3 Online Meta-Learning

After training the word-level domain mixing NMT model, we use it as a teacher model to initialize the meta-learning process. Algorithm 1 shows the complete algorithm.

Split Tasks. Zhan et al. (2021) propose a curriculum-based task splitting strategy, which uses divergence scores computed by a language model as the curriculum to split the corpus into small tasks. We follow a similar idea, but propose to use predictions from a domain classifier as the criterion for splitting the data. Concretely, we first train a domain classifier with BERT; the classifier scores sentences, indicating domain similarity between an in-domain sentence and a general domain sentence (see Section 3.2). The tasks are then split according to the scores; sentences more similar to the general domain sentences are selected in early tasks.

Balanced Sampling. Previous meta-learning approaches (Sharaf et al., 2020; Zhan et al., 2021) are based on token-size based sampling, which uses $8k$ or $16k$ token sizes split into many small tasks. However, the splitting process for the domain is not balanced, since some tasks did not contain all seen domains, especially in the early tasks. As we can see in Figure 2, the token-based splitting methods

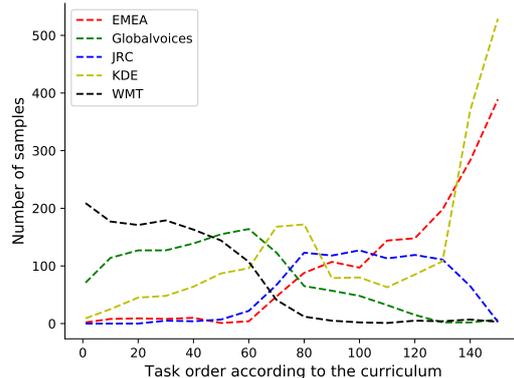


Figure 2: The statistic of samples in the task for the tokenization-based splitting strategy. More general domains are on the left and the more distinctive domains are on the right.

usually allocate more samples on domain-similar domains (*WMT*, *Globalvoices*) and allocate small samples on domain-distant domains (*EMEA*, *JRC*) in the sampling of early tasks. This can cause problems in our method since the model architecture is dynamically changing according to the number of domains (see more details in Section 3.1).

To address these issues, we sample the data uniformly from the domains to compensate for imbalanced domain distributions based on domain classifier scores.

Meta-Training. Following the balanced sampling, the process of meta-training is to update the current model parameter on $\mathcal{T}_{support}$ from θ to θ' , and then evaluate on \mathcal{T}_{query} . The model parameter θ' is updated to minimize the meta-learning loss through MAML.

Given a pre-trained model f_θ (initialized with parameters θ trained on word-level domain mixing) and the meta-train data $\mathcal{D}_{meta-train}$, for each task \mathcal{T} , we learn to use one gradient update to update the model parameters from θ to θ' as follows:

$$\theta' = \theta - \alpha \nabla_{\theta} L_{\mathcal{T}}(f_{\theta})$$

where α is the learning rate and L is the loss function. In our methods, we consider both the traditional sentence-level meta-learning loss $\mathcal{L}_{\mathcal{T}}(f_{\theta})$ and the word-level loss $\Gamma_{\mathcal{T}}(f_{\theta})$ (L^* of \mathcal{T}) calculated from the word-level domain mixing pre-trained model. More formally, the loss is updated as follows:

$$L_{\mathcal{T}}(f_{\theta}) = \mathcal{L}_{\mathcal{T}}(f_{\theta}) + \Gamma_{\mathcal{T}}(f_{\theta}).$$

Algorithm 1 RMLNMT (Robust Meta-Learning NMT Domain Adaptation)

Require: Domain classifier model cls ; Pretrained domain-mixing model θ ;

- 1: Score the sentence in $\mathcal{D}_{\text{meta-train}}$ using cls
- 2: **for** N epochs **do**
- 3: Split corpus into n tasks based on step 1
- 4: Balance sample through all tasks
- 5: **for** task $\mathcal{T}_i, i = 1 \dots n$ **do**
- 6: Evaluate loss $L_{\mathcal{T}}(f_{\theta})$
 $= \mathcal{L}_{\mathcal{T}_i}(f_{\theta}) + \Gamma_{\mathcal{T}_i}(f_{\theta})$ on support set
- 7: Update the gradient with parameters
 $\theta' = \theta - \alpha \nabla_{\theta} L_{\mathcal{T}}(f_{\theta})$
- 8: **end for**
- 9: Update the gradient with parameters
 $\theta = \theta - \beta \nabla_{\theta} L_{\mathcal{T}}(f_{\theta'})$ on query set
- 10: **end for**
- 11: **return** RMLNMT model parameter θ

Note that the meta-training phase is not adapted to a specific domain, so it can be used as a metric to evaluate the domain robustness of the model.

Meta-Adaptation. After the meta-training phase, the parameters are updated to adapt to each domain using the small *support set* of $\mathcal{D}_{\text{meta-test}}$ corpus to simulate the low-resource scenarios. Then performance is evaluated on the *query set* of $\mathcal{D}_{\text{meta-test}}$.

4 Experiments

Datasets. We experiment with English→German (*en2de*) and English→Chinese (*en2zh*) translation tasks. For the *en2de* task, we use the same corpora as Zhan et al. (2021). The data consists of corpora in nine domains (Bible, Books, ECB, EMEA, GlobalVoices, JRC, KDE, TED, WMT-News) publicly available on OPUS² (Tiedemann, 2012) and the COVID-19 corpus³. For *en2zh*, we use UM-Corpus (Tian et al., 2014) containing eight domains: Education, Microblog, Science, Subtitles, Laws, News, Spoken, Thesis. We use WMT14 (*en2de*) and WMT18 (*en2zh*) corpus published on the WMT website⁴ as our general domain corpora. We use WMT19 English monolingual corpora to train the LM model so that we can reproduce results from previous work.

²opus.nlpl.eu

³github.com/NLP2CT/Meta-Curriculum

⁴<http://www.statmt.org>

Data Preprocessing. For English and German, we preprocessed all data with the Moses tokenizer⁵ and use sentencepiece⁶ (Kudo and Richardson, 2018) to encode the corpus with a joint vocabulary, with size 40,000. After that, we filter the sentence longer than 175 tokens and deduplicate the corpus. For Chinese, we perform word segmentation using the Stanford Segmenter (Tseng et al., 2005). To have a fair comparison with previous methods (Sharaf et al., 2020; Zhan et al., 2021), we use the same setting, which randomly sub-sampled $\mathcal{D}_{\text{meta-train}}$ and $\mathcal{D}_{\text{meta-test}}$ for each domain with fixed token sizes in order to simulate domain adaptation tasks in low-resource scenarios. More details for data used in this paper can be found in Appendix A.1.

Baselines. We compare RMLNMT with the following baselines:

- **Vanilla.** A standard Transformer-based NMT system trained on the general domains (WMT14 for *en2de*, WMT18 for *en2zh*) and $\mathcal{D}_{\text{meta-train}}$ corpus in seen-domains. We use the $\mathcal{D}_{\text{meta-train}}$ corpus because meta-learning-based methods also use the $\mathcal{D}_{\text{meta-train}}$ corpus, this is a more fair and stronger baseline.
- **Plain fine-tuning.** Fine-tune the vanilla system on support set of $\mathcal{D}_{\text{meta-test}}$ for each individual domain.
- **Tag.** prepend a domain tag to each sentence to indicate what domain it belongs to (Kobus et al., 2017).
- **Meta-MT.** Standard meta-learning approach on domain adaptation task (Sharaf et al., 2020).
- **Meta-Curriculum (LM).** Meta-learning approach for domain adaptation using LM score as the curriculum to sample the task (Zhan et al., 2021).
- **Meta-based w/o FT.** This series of experiments uses the meta-learning system prior to adaptation to the specific domain. This can be used to evaluate the domain robustness of meta-based models (see more details in the meta-training part of Section 3.3).

⁵github.com/moses-smt/mosesdecoder

⁶github.com/google/sentencepiece

Models	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
1 Vanilla	24.34	12.08	12.61	29.96	27.89	37.27	24.19	39.84	27.75	27.38
2 Vanilla + tag	24.86	12.04	12.46	30.03	27.93	38.37	24.56	40.75	28.23	27.26
3 Meta-MT w/o FT	23.69	11.07	12.10	29.04	26.86	30.94	23.73	38.82	23.04	26.13
4 Meta-Curriculum (LM) w/o FT	23.70	11.16	12.24	28.22	27.21	33.49	24.27	39.21	27.60	25.83
5 RMLNMT w/o FT	25.48	11.48	13.11	31.42	28.05	47.00	26.35	51.13	32.80	28.37

Table 1: Domain Robustness: BLEU scores on the English \rightarrow German translation task. *w/o* denotes the meta-learning systems without fine-tuning, FT denotes fine-tuning. Best results are highlighted in bold.

Models	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
1 Plain FT	24.81	12.61	12.78	30.48	28.36	37.26	24.26	40.02	27.99	27.31
2 Plain FT + tag	25.31	12.57	12.83	30.57	28.39	39.54	24.91	41.51	29.14	27.58
3 Meta-MT + FT	25.83	14.20	13.39	30.36	28.57	34.69	24.64	39.15	27.47	26.38
4 Meta-Curriculum (LM) + FT	26.66	14.37	13.70	30.41	28.97	34.00	24.72	39.61	27.37	26.68
5 RMLNMT + FT	26.53	15.37	13.72	31.97	29.47	47.02	26.55	51.13	32.88	28.37

Table 2: Domain Adaptability: BLEU scores on the English \rightarrow German translation task.

Implementation. We use the Transformer model (Vaswani et al., 2017) as implemented in FairSeq⁷ (Ott et al., 2019). For our word-level domain-mixing modules, we dynamically adjust the network structure according to the number of domains since every domain has its multi-head layers. Hence, the number of model parameters in the attentive sub-layers of RMLNMT is k times the number in the standard transformer (k is the number of seen domains in the training data). Following Jiang et al. (2020), we enlarged the baseline models to have \sqrt{k} times larger embedding dimension, so the baseline has the same number of parameters. This should rule out that the improvements are due to increased parameter count rather than modeling improvements. For our meta-learning framework, we consider the general meta loss and word-adaptive loss together (as seen in Section 3.3). Following Zhan et al. (2021), the fine-tuning process in each models is strictly limited to 20 to simulate quick adaptation. Note that the meta-train stage only uses the seen domain corpus and the unseen domain corpus is only used in the meta-test stage. More details on hyper-parameters are listed in Appendix A.2.

Evaluation. For a fair comparison with previous work, we use the same data from the support set of $\mathcal{D}_{\text{meta-test}}$ to finetune the model and the same data from the query set of $\mathcal{D}_{\text{meta-test}}$ to evaluate the models. We measure case-sensitive detokenized BLEU with SacreBLEU⁸ (Post, 2018); beam search with a

beam of size five is used. Because of the recent criticism of BLEU score (Mathur et al., 2020), we also evaluate our models using chrF (Popović, 2015) and COMET⁹ (Rei et al., 2020); the results are listed in Appendix A.5.

Domain Robustness. Domain robustness shows the effectiveness of the model both in seen and unseen domains. Hence, we use the model without fine-tuning to evaluate the domain robustness.

Domain Adaptability. We evaluate the domain adaptability by testing that the model quickly adapts to new domains using just hundreds of in-domain parallel sentences. Therefore, we fine-tune the models on a small amount of domain-specific data.

Cross-Domain Robustness. To better show the cross-domain robustness of RMLNMT, we use the fine-tuned model of one specific domain to generate the translation for other domains. More formally, given k domains, we use the fine-tuned model M_J with the domain label of J to generate the translation of k domains.

5 Results

Table 1 and Table 3 show the domain robustness for English \rightarrow German and English \rightarrow Chinese respectively. Table 2 and Table 4 show the domain adaptability on both translation task.

Domain Robustness. As seen in Table 1 and Table 3, RMLNMT shows the best domain robust-

⁷github.com/facebookresearch/fairseq

⁸github.com/mjpost/sacrebleu

⁹github.com/Unbabel/COMET

Models	Unseen				Seen			
	Education	Microblog	Science	Subtitles	Laws	News	Spoken	Thesis
1 Vanilla	27.52	26.05	31.58	18.32	46.69	28.67	26.44	29.00
2 Vanilla + tag	27.36	26.11	31.53	18.25	47.13	28.75	26.71	29.19
3 Meta-MT w/o FT	28.76	26.41	32.41	17.38	43.74	27.31	25.98	28.11
4 Meta-Curriculum (LM) w/o FT	28.53	26.14	32.25	17.45	43.87	27.25	27.57	28.23
5 RMLNMT w/o FT	30.17	28.42	34.20	19.89	57.54	30.39	28.11	33.20

Table 3: Domain Robustness: BLEU scores on English \rightarrow Chinese translation tasks.

Models	Unseen				Seen			
	Education	Microblog	Science	Subtitles	Laws	News	Spoken	Thesis
1 Plain FT	27.05	26.31	32.09	17.77	47.64	28.28	25.73	28.47
2 Plain FT + tag	27.13	26.48	32.12	17.94	47.91	28.84	26.35	29.58
3 Meta-MT + FT	29.33	27.48	33.12	18.77	45.21	28.43	26.82	29.20
4 Meta-Curriculum (LM) + FT	28.91	27.20	33.19	18.93	45.46	28.17	27.84	29.47
5 RMLNMT + FT	30.91	28.52	34.51	20.13	57.58	30.42	28.03	32.25

Table 4: Domain Adaptability: BLEU scores on English \rightarrow Chinese translation tasks.

Methods	Avg
Meta-MT	-1.97
Meta-Curriculum (LM)	-0.96
Meta-Curriculum (cls)	-0.98
RMLNMT	2.64

Table 5: The average improvement over vanilla baseline.

ness compared with other models both in seen and unseen domains. In addition, the traditional meta-learning approach (Meta-MT, Meta-Curriculum) without fine-tuning is even worse than the standard transformer model in seen domains. This phenomenon is our motivation for improving the robustness of traditional meta-learning based approach. In other words, we cannot be sure whether the improvement of the meta-based method is due to the domain adaptability of meta-learning or the robustness of the teacher model. Note this setup differs from the previous work (Sharaf et al., 2020; Zhan et al., 2021) because we included the $\mathcal{D}_{\text{meta-train}}$ data to the vanilla system to insure all systems in the table use the same training data.¹⁰ Interestingly, the translation quality in the *WMT* domain is also improved which is different than (Zhan et al., 2021). They explain that their methods achieve maximum robustness on the *WMT* domain, while our results demonstrate that our model can further improve robustness even when trained on the same domain as the pre-trained model.

¹⁰We also confirmed with Zhan et al. (2021) via email that they did not deduplicate the corpus, which is another reason for the difference between our results and their results.

Domain Adaptability. From Tables 2 and 4, we observe that the traditional meta-learning approach shows high adaptability to unseen domains but fails on seen domains due to limited domain robustness. In contrast, RMLNMT shows its domain adaptability both in seen and unseen domains, and maintains the domain robustness simultaneously. Compared with RMLNMT, the traditional meta-learning approach show more improvement between the *w/o FT* model and *FT* model. For example, *Meta-MT* and *Meta-Curriculum (LM)* obtains 1.32 and 2.19 BLEU score improvement after finetuning on the *ECB* domain; improvement from *RMLNMT* only got 0.55. This phenomenon meets our expectations since *RMLNMT* without finetuning is already strong enough due to the domain robustness of word-level domain mixing. In other words, the improvement of the traditional meta-learning approach is to some extent due to the unrobustness of the model.

Cross-Domain Robustness. Table 5 reports the average difference of $k \times k$ BLEU scores; a larger positive value means a more robust model. We observed that the plain meta-learning based methods have a negative value, which means the performance gains in the specific domains come at the cost of performance decreases in other domains. In other words, the model is not domain robust enough. In contrast, RMLNMT has a positive difference with the vanilla system, showing that the model is robust. The specific BLEU scores are shown in Figure 3 of Appendix A.4.

The results of both domain robustness and do-

Classifier	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
CNN	24.12	13.57	12.74	30.31	28.14	46.12	25.17	50.52	31.15	26.34
BERT-many-labels	25.89	14.77	13.71	32.10	29.28	47.41	26.70	51.34	32.76	28.17
BERT-2-labels	26.10	14.85	13.58	31.99	29.17	46.80	26.46	51.56	32.83	28.37
mBERT-many-labels	26.10	14.73	13.69	31.93	29.11	47.02	26.33	51.13	32.69	27.91
mBERT-2-labels	26.53	15.37	13.71	31.97	29.47	47.02	26.55	51.13	32.88	28.37

Table 6: Different classifier: BLEU scores on the English → German translation task.

Sampling Strategy	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
Token-based sampling	25.30	11.38	12.70	31.61	28.01	47.51	26.50	51.31	32.88	28.03
Balance sampling	25.47	11.51	12.79	32.08	28.98	47.64	26.58	51.25	32.91	28.07

Table 7: Different sampling strategy: BLEU scores on the English → German translation task.

Finetune Strategy	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
FT-unseen	25.23	13.18	12.73	32.45	28.41	46.35	25.83	50.85	32.30	26.88
FT-seen	24.58	11.73	12.57	30.79	27.29	46.58	25.73	50.91	31.78	26.51
FT-all	15.00	7.77	9.06	21.33	16.98	24.69	14.63	27.59	12.77	15.75
FT-specific	26.53	15.37	13.71	31.97	29.47	47.02	26.33	51.13	32.83	28.37

Table 8: Different fine-tuning strategy: BLEU scores on the English → German translation task.

main adaptability are consistent for the chrF and COMET evaluation metrics (see more details in Tables 13 and 14 of Appendix A.5).

6 Analysis

In this section, we conduct additional experiments to better understand the strengths of RMLNMT. We analyze the contribution of different components in RMLNMT, through an ablation study.

Different classifiers. We evaluate the impact of different classifiers on translation performance. The main results are as shown in Table 6 (see more details in Appendix A.3). We observed that the performance of RMLNMT is not directly proportional to the accuracy of the classifier. In other words, slightly higher classification accuracy does not lead to better BLEU scores. This is because the accuracy of the classifier is close between BERT-based models and the primary role of the classifier is to construct the curriculum for splitting the tasks. When we use a significantly worse classifier, i.e., the CNN in our experiments, the overall performance of RMLNMT is worse than the BERT-based classifier.

Balanced sampling vs. Token-based sampling. Plain meta-learning uses a token-based sampling strategy to split sentences into small tasks. How-

ever, the token-based strategy could cause unbalanced domain distribution in some tasks, especially in the early stage of training due to domain mismatches (see the discussion of balanced sampling in Section 3.3). To address this issue, we proposed to balance the domain distribution after splitting the task. Table 7 shows that our methods can result in small improvements in performance. For example, in the *TED* domain, BLEU was 28.01 with token-based sampling, but with the balanced sampling strategy BLEU was 28.98. We keep the same number of tasks to have a fair comparison with previous methods.

Different fine-tuning strategies. As described in Section 3.1, the model for each domain has its own multi-head and feed-forward layers. During the fine-tuning stage of RMLNMT, we devise four strategies: i) **FT-unseen**: fine-tuning using all unseen domain corpora; ii) **FT-seen**: fine-tuning using all seen domain corpora; iii) **FT-all**: fine-tuning using all out-of-domain corpora (seen and unseen domains); iv) **FT-specific**: using the specific domain corpus to fine-tune the specific models. The results are shown in Table 8. *FT-specific* obtains robust results among all the strategies. Although other strategies outperform *FT-specific* in some domains, *FT-specific* is robust across all domains. Furthermore, *FT-specific* is the fairest comparison

because it uses only a specific domain corpus to fine-tune, which is the same as the baseline systems.

7 Related Work

Domain Adaptation for NMT. Current approaches can be categorized into two groups by granularity: From a sentence-level perspective, researchers either use data selection methods (Moore and Lewis, 2010; Axelrod et al., 2011) to select the training data that is similar to out-of-domain parallel corpora or train a classifier (Rieß et al., 2021) or utilize a language model (Wang et al., 2017; Zhan et al., 2021) to better weight the sentences. From a word-level perspective, researchers try to model domain distribution at the word level, since a word in a sentence can be related to more domains than just the sentence domain (Zeng et al., 2018; Yan et al., 2018; Hu et al., 2019; Sato et al., 2020; Jiang et al., 2020).

Curriculum Learning for NMT. Curriculum learning (Bengio et al., 2009) starts with easier tasks and then progressively gain experience to process more complex tasks, which has proved to be useful in NMT domain adaptation. Stojanovski and Fraser (2019) utilize curriculum learning to improve anaphora resolution in NMT systems. Zhang et al. (2019) and Zhan et al. (2021) use a language model to compute a similarity score between domains, from which a curriculum is devised for adapting NMT systems to specific domains from general domains.

Meta-Learning for NMT. Gu et al. (2018) apply model-agnostic meta-learning (MAML; Finn et al., 2017) to NMT. They show that MAML effectively improves low-resource NMT. Li et al. (2020), Sharaf et al. (2020) and Zhan et al. (2021) propose to formulate the problem of low-resource domain adaptation in NMT as a meta-learning problem: the model learns to quickly adapt to an unseen new domain from a general domain.

8 Conclusion

We presented RMLNMT, a robust meta-learning framework for low-resource NMT domain adaptation reaching both high domain adaptability and domain robustness (both in the seen domains and unseen domains). We found that domain robustness dominates the results compared to domain adaptability in meta-learning based approaches. The

results show that RMLNMT works best in setups that require high robustness in low-resource scenarios.

Acknowledgement

We thank Mengjie Zhao for the helpful comments. This work was supported by funding to Wen Lai’s PhD research from LMU-CSC (China Scholarship Council) Scholarship Program. This work has received funding from the European Research Council under the European Union’s Horizon 2020 research and innovation program (grant agreement #640550). This work was also supported by the DFG (grant FR 2829/4-1). The work at CUNI was supported by the European Commission via its Horizon 2020 research and innovation programme (870930).

References

- Amittai Axelrod, Xiaodong He, and Jianfeng Gao. 2011. [Domain adaptation via pseudo in-domain data selection](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 355–362, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *3rd International Conference on Learning Representations*.
- Trapit Bansal, Rishikesh Jha, Tsendsuren Munkhdalai, and Andrew McCallum. 2020. [Self-supervised meta-learning for few-shot natural language classification tasks](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 522–534, Online. Association for Computational Linguistics.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. [Curriculum learning](#). In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48.
- Praveen Dakwale and Christof Monz. 2017. [Finetuning for neural machine translation with limited degradation across in-and out-of-domain data](#). In *Proceedings of the XVI Machine Translation Summit*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

- Zi-Yi Dou, Keyi Yu, and Antonios Anastasopoulos. 2019. [Investigating meta-learning algorithms for low-resource natural language understanding tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1192–1197, Hong Kong, China. Association for Computational Linguistics.
- Yingjun Du, Nithin Holla, Xiantong Zhen, Cees Snoek, and Ekaterina Shutova. 2021. [Meta-learning with variational semantic memory for word sense disambiguation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5254–5268, Online. Association for Computational Linguistics.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. [Model-agnostic meta-learning for fast adaptation of deep networks](#). In *International Conference on Machine Learning*, pages 1126–1135. PMLR.
- Markus Freitag and Yaser Al-Onaizan. 2016. [Fast domain adaptation for neural machine translation](#). *arXiv preprint arXiv:1612.06897*.
- Robert M. French. 1999. [Catastrophic forgetting in connectionist networks](#). *Trends in Cognitive Sciences*, 3(4):128–135.
- Jiatao Gu, Yong Wang, Yun Chen, Victor O. K. Li, and Kyunghyun Cho. 2018. [Meta-learning for low-resource neural machine translation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3622–3631, Brussels, Belgium. Association for Computational Linguistics.
- Timothy M Hospedales, Antreas Antoniou, Paul Micaelli, and Amos J Storkey. 2021. [Meta-learning in neural networks: A survey](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Junjie Hu, Mengzhou Xia, Graham Neubig, and Jaime Carbonell. 2019. [Domain adaptation of neural machine translation by lexicon induction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2989–3001, Florence, Italy. Association for Computational Linguistics.
- Haoming Jiang, Chen Liang, Chong Wang, and Tuo Zhao. 2020. [Multi-domain neural machine translation with word-level adaptive layer-wise domain mixing](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1823–1834, Online. Association for Computational Linguistics.
- Yoon Kim. 2014. [Convolutional neural networks for sentence classification](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751, Doha, Qatar. Association for Computational Linguistics.
- Catherine Kobus, Josep Crego, and Jean Senellart. 2017. [Domain control for neural machine translation](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 372–378, Varna, Bulgaria. INCOMA Ltd.
- Taku Kudo. 2018. [Subword regularization: Improving neural network translation models with multiple subword candidates](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.
- Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Rumeng Li, Xun Wang, and Hong Yu. 2020. [Metamt, a meta learning method leveraging multiple domain data for low resource machine translation](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8245–8252.
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. [Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online. Association for Computational Linguistics.
- Robert C. Moore and William Lewis. 2010. [Intelligent selection of language model training data](#). In *Proceedings of the ACL 2010 Conference Short Papers*, pages 220–224, Uppsala, Sweden. Association for Computational Linguistics.
- Mathias Müller, Annette Rios Gonzales, and Rico Senrich. 2020. [Domain robustness in neural machine translation](#). In *Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 151–164.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. 2016. [Distillation as a](#)

- defense to adversarial perturbations against deep neural networks. In *2016 IEEE symposium on security and privacy (SP)*, pages 582–597. IEEE.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Kun Qian and Zhou Yu. 2019. [Domain adaptive dialog generation via meta learning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2639–2649, Florence, Italy. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. [COMET: A neural framework for MT evaluation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Simon Rieß, Matthias Huck, and Alex Fraser. 2021. [A comparison of sentence-weighting techniques for NMT](#). In *Proceedings of the 18th Biennial Machine Translation Summit (Volume 1: Research Track)*, pages 176–187, Virtual. Association for Machine Translation in the Americas.
- Shoetsu Sato, Jin Sakuma, Naoki Yoshinaga, Masashi Toyoda, and Masaru Kitsuregawa. 2020. [Vocabulary adaptation for domain adaptation in neural machine translation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4269–4279, Online. Association for Computational Linguistics.
- Amr Sharaf, Hany Hassan, and Hal Daumé III. 2020. [Meta-learning for few-shot NMT adaptation](#). In *Proceedings of the Fourth Workshop on Neural Generation and Translation*, pages 43–53, Online. Association for Computational Linguistics.
- Dario Stojanovski and Alexander Fraser. 2019. [Improving anaphora resolution in neural machine translation using curriculum learning](#). In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 140–150, Dublin, Ireland. European Association for Machine Translation.
- Brian Thompson, Jeremy Gwinnup, Huda Khayrallah, Kevin Duh, and Philipp Koehn. 2019. [Overcoming catastrophic forgetting during domain adaptation of neural machine translation](#). 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)*, pages 2062–2068, Minneapolis, Minnesota. Association for Computational Linguistics.
- Liang Tian, Derek F. Wong, Lidia S. Chao, Paulo Quaresma, Francisco Oliveira, Yi Lu, Shuo Li, Yiming Wang, and Longyue Wang. 2014. [UM-corpus: A large English-Chinese parallel corpus for statistical machine translation](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 1837–1842, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Jörg Tiedemann. 2012. [Parallel data, tools and interfaces in OPUS](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).
- Huihsin Tseng, Pichuan Chang, Galen Andrew, Daniel Jurafsky, and Christopher Manning. 2005. [A conditional random field word segmenter for sighthan bake-off 2005](#). In *Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing*.
- Zhaopeng Tu, Yang Liu, Lifeng Shang, Xiaohua Liu, and Hang Li. 2017. [Neural machine translation with reconstruction](#). In *Thirty-First AAAI Conference on Artificial Intelligence*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in neural information processing systems*, pages 5998–6008.
- Lingxiao Wang, Kevin Huang, Tengyu Ma, Quanquan Gu, and Jing Huang. 2021. [Variance-reduced first-order meta-learning for natural language processing tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2609–2615, Online. Association for Computational Linguistics.
- Rui Wang, Masao Utiyama, Lemao Liu, Kehai Chen, and Eiichiro Sumita. 2017. [Instance weighting for neural machine translation domain adaptation](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1482–1488, Copenhagen, Denmark. Association for Computational Linguistics.
- Shen Yan, Leonard Dahlmann, Pavel Petrushkov, Sanjika Hewavitharana, and Shahram Khadivi. 2018. [Word-based domain adaptation for neural machine translation](#). In *Proceedings of the 15th International Conference on Spoken Language Translation*, pages 31–38, Brussels. International Conference on Spoken Language Translation.
- Kyra Yee, Yann Dauphin, and Michael Auli. 2019. [Simple and effective noisy channel modeling for neural machine translation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*

(EMNLP-IJCNLP), pages 5696–5701, Hong Kong, China. Association for Computational Linguistics.

Wenpeng Yin. 2020. [Meta-learning for few-shot natural language processing: A survey](#). *arXiv preprint arXiv:2007.09604*.

Changlong Yu, Jialong Han, Haisong Zhang, and Wilfred Ng. 2020. [Hypernymy detection for low-resource languages via meta learning](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3656, Online. Association for Computational Linguistics.

Jiali Zeng, Jinsong Su, Huating Wen, Yang Liu, Jun Xie, Yongjing Yin, and Jianqiang Zhao. 2018. [Multi-domain neural machine translation with word-level domain context discrimination](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 447–457, Brussels, Belgium. Association for Computational Linguistics.

Runzhe Zhan, Xuebo Liu, Derek F. Wong, and Lidia S. Chao. 2021. [Meta-curriculum learning for domain adaptation in neural machine translation](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 14310–14318.

Xuan Zhang, Pamela Shapiro, Gaurav Kumar, Paul McNamee, Marine Carpuat, and Kevin Duh. 2019. [Curriculum learning for domain adaptation in neural machine translation](#). 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)*, pages 1903–1915, Minneapolis, Minnesota. Association for Computational Linguistics.

A Appendix

A.1 Datasets

For the OPUS corpus used in the English \rightarrow German translation task, we deduplicated the corpus, which is different from (Zhan et al., 2021) and is the main reason that we cannot reproduce the results in the original paper. The statistics of the original OPUS are shown in Table 9. The seen domains (EMEA, Globalvoices, JRC, KDE, WMT) contain a lot of duplicated sentences. The scores in the original paper are too high because the $\mathcal{D}_{\text{meta-train}}$ dataset overlaps with some sentences in $\mathcal{D}_{\text{meta-test}}$.

For the meta-learning phase, to have a fair comparison with previous methods, we use the same setting. We random split 160 tasks and 10 tasks respectively in $\mathcal{D}_{\text{meta-train}}$ and $\mathcal{D}_{\text{meta-test}}$ to simulate the low-resource scenarios. For each task, the token amount of support set and query set is a strict limit to $8K$ and $16K$. $\mathcal{D}_{\text{meta-dev}}$ corpus is limited to 5000 sentences for each domain. Table 10 and Table 11 shows the detailed statistics of the English \rightarrow German and English \rightarrow Chinese tasks.

Corpus	Original	Deduplicated
Covid	3,325	3,312
Bible	62,195	61,585
Books	51,467	51,106
ECB	113,116	113,081
TED	143,830	142,756
EMEA	1,103,807	360,833
Globalvoices	71,493	70,519
JRC	717,988	503,789
KDE	223,672	187,918
WMT	45,913	34,727

Table 9: Data statistic (sentences) of the original corpus for English \rightarrow German translation task

	$\mathcal{D}_{\text{meta-train}}$		$\mathcal{D}_{\text{meta-test}}$	
	Support	Query	Support	Query
Covid	/	/	309	612
Bible	/	/	280	548
Books	/	/	304	637
ECB	/	/	295	573
TED	/	/	390	772
EMEA	14856	29668	456	975
Globalvoices	11686	23319	368	699
JRC	7863	15769	254	519
KDE	24078	48284	756	1510
WMT	10939	21874	334	704

Table 10: Data statistic (sentences) of the meta-learning stage for English \rightarrow German translation task

A.2 Model Configuration

We use the Transformer Base architecture (Vaswani et al., 2017) as implemented in fairseq (Ott et al., 2019). We use the standard Transformer architecture with dimension 512, feed-forward layer 2048, 8 attention heads, 6 encoder layers and 6 decoder layers. For optimization, we use the Adam optimizer with a learning rate of $5 \cdot 10^{-5}$. To prevent overfitting, we applied a dropout of 0.3 on all layers. The number of warm-up steps was set to 4000. At the time of inference, a beam search of size 5 is used to balance the decoding time and accuracy of the search.

For the word-level domain-mixing model, we use the same setting as Jiang et al. (2020). The number of parameters of our model is dynamically adjusted with the domain numbers and k times higher than standard model architecture, since every domain has its multi-head attention layer and feed-forward layer. To have a fair comparison between baselines, we enlarged the baseline models to have \sqrt{k} times larger embedding dimension, so

	$\mathcal{D}_{\text{meta-train}}$		$\mathcal{D}_{\text{meta-test}}$	
	Support	Query	Support	Query
Education	/	/	395	785
Microblog	/	/	358	721
Science	/	/	392	852
Subtitles	/	/	612	1219
Laws	6379	13001	197	416
News	9004	18362	281	536
Spoken	18270	36569	571	1148
Thesis	8914	17883	298	547

Table 11: Data statistic (sentences) of the meta-learning stage for English→Chinese translation task

Classifier	Acc(%)
CNN	74.91%
BERT: many-labels	96.12%
BERT: 2-labels	95.35%
mBERT: many-labels	95.41%
mBERT: 2-labels	95.26%

Table 12: The accuracy of the different classifiers.

the baseline has the same number of parameters.

A.3 Different classifiers

With a general in-domain corpus and some out-of-domain corpora, we train five classifiers. We experiment with two different labeling schemes: `2-labels` where we distinguish only two classes: *out-of-domain* and *in-domain*; `many-labels` where sentences are labeled with the respective domain labels. Further, we experiment with two variants of the BERT model: first, we use monolingual English BERT on the source side only, and second, we use multilingual BERT (mBERT) to classify the parallel sentence pairs. For further comparison, we include also a CNN-based classifier (Kim, 2014). We present the accuracy of the English-German domain classifier in Table 12.

A.4 Cross-Domain Robustness

In Figure 3 we show the detailed results ($k \times k$ scores) of cross-domain robustness.

A.5 Evaluations

In addition to BLEU, we also use chrF (Popović, 2015) and COMET (Rei et al., 2020) as evaluation metrics. Table 13 and Table 14 show the results. Consistently with the BLEU score (Tables 1 and Table 2), we observed that RMLNMT is more effective than all previous methods.

Models	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
1 Vanilla	0.550	0.418	0.385	0.538	0.542	0.599	0.536	0.614	0.525	0.558
Plain FT	0.555	0.423	0.388	0.540	0.548	0.600	0.536	0.618	0.528	0.558
2 Vanilla + tag	0.555	0.418	0.384	0.540	0.544	0.657	0.545	0.627	0.531	0.558
Plain FT + tag	0.562	0.423	0.388	0.540	0.549	0.602	0.536	0.694	0.547	0.561
3 Meta-MT w/o FT	0.545	0.410	0.382	0.498	0.538	0.532	0.531	0.610	0.464	0.553
Meta-MT + FT	0.566	0.432	0.390	0.542	0.556	0.582	0.538	0.613	0.522	0.552
4 Meta-Curriculum (LM) w/o FT	0.548	0.412	0.384	0.523	0.543	0.560	0.536	0.611	0.521	0.554
Meta-Curriculum (LM) + FT	0.567	0.434	0.395	0.544	0.548	0.572	0.539	0.615	0.522	0.553
5 RMLNMT w/o FT	0.555	0.405	0.388	0.557	0.544	0.656	0.552	0.702	0.574	0.561
RMLNMT + FT	0.562	0.451	0.395	0.558	0.560	0.656	0.552	0.702	0.574	0.561

Table 13: chrF scores on the English → German translation task.

Models	Unseen					Seen				
	Covid	Bible	Books	ECB	TED	EMEA	Globalvoices	JRC	KDE	WMT
1 Vanilla	0.4967	-0.1250	-0.2225	0.3276	0.3400	0.3096	0.3199	0.5430	0.1836	0.4326
Plain FT	0.5066	-0.1105	-0.1985	0.3315	0.3553	0.3177	0.3276	0.5492	0.1813	0.4392
2 Vanilla + tag	0.4970	-0.1250	-0.2228	0.3277	0.3401	0.3176	0.3291	0.5495	0.1846	0.4311
Plain FT + tag	0.5078	-0.1105	-0.1981	0.3315	0.3553	0.3179	0.3341	0.5572	0.1973	0.4398
3 Meta-MT w/o FT	0.4850	-0.1454	-0.2228	0.0953	0.3506	0.0524	0.2985	0.5319	0.1304	0.4137
Meta-MT + FT	0.5175	-0.0650	-0.1878	0.3466	0.3824	0.2678	0.3189	0.5509	0.1316	0.4161
4 Meta-Curriculum (LM) w/o FT	0.4879	-0.1365	-0.2122	0.2568	0.3751	0.1968	0.3273	0.5246	0.0962	0.4206
Meta-Curriculum (LM) + FT	0.5347	-0.0604	-0.1773	0.3460	0.3729	0.2366	0.3141	0.5430	0.1467	0.4128
5 RMLNMT w/o FT	0.4943	-0.1956	-0.2179	0.3580	0.3394	0.4026	0.3769	0.6797	0.3014	0.4255
RMLNMT + FT	0.5302	-0.0543	-0.1610	0.3547	0.3867	0.4046	0.3771	0.6797	0.3015	0.4256

Table 14: COMET scores on the English → German translation task.

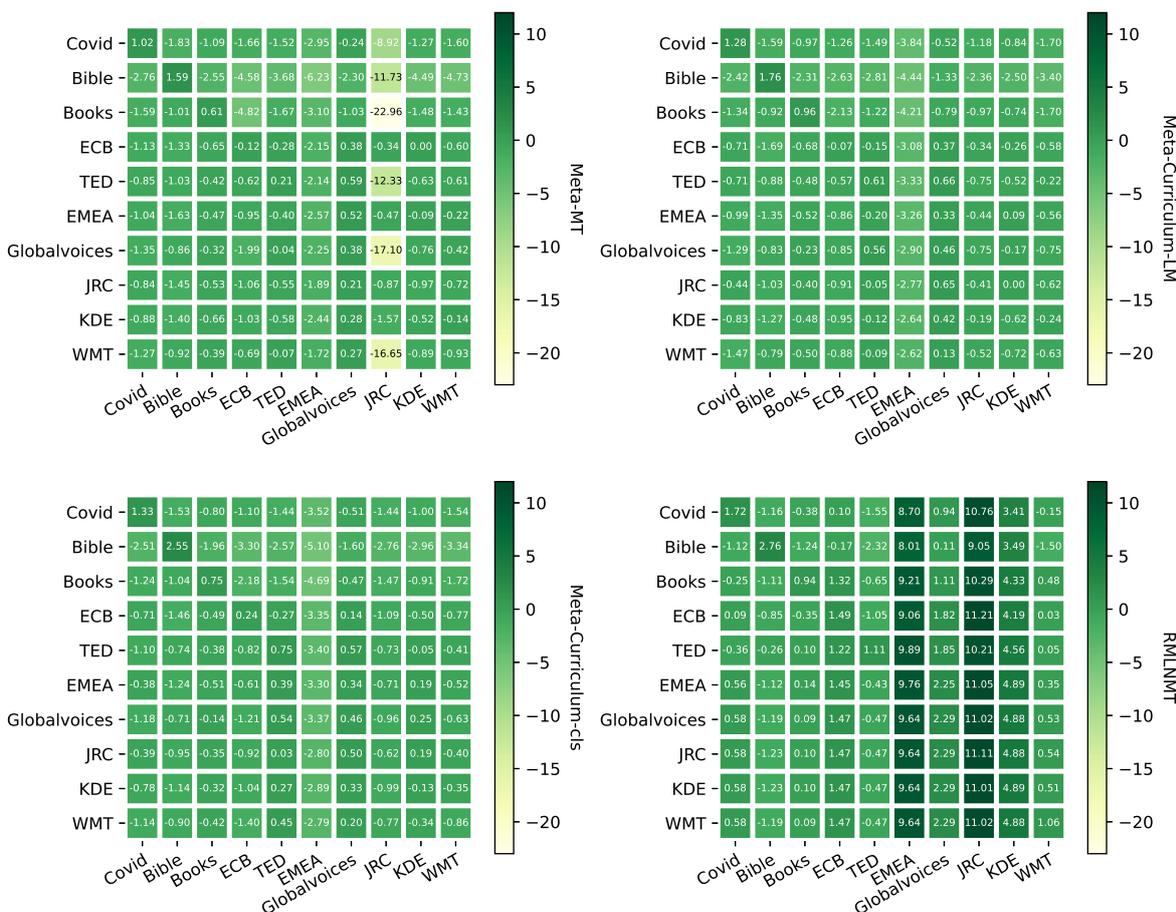


Figure 3: BLEU scores for one specific finetuned model on other domains for en2de translation.