

Towards a Fine-Grained Multi-Domain Neural Machine Translation Using Inter-Domain Relationships

Anonymous ACL submission

Abstract

While research on the domain adaptation task in neural machine translation has become popular recently, there exists no agreement on what constitutes a domain, and most previous studies only focus on coarse-grained domain adaptation and their methods cannot be generalized if the domain size is large. In this work, we argue the necessity to study a fine-grained domain adaptation problem. We build a new multilingual dataset from web sources that focus on fine-grained domains and inter-domain attributes and relationships. We also propose a simple but effective adaptation method to incorporate domain knowledge leveraging models in information networks.

1 Introduction

The success of most machine learning algorithms heavily relies on the assumption that feature spaces and underlying distributions of training data and testing data are similar. In reality, such assumption is often violated. This motivates the study of domain adaptation in many areas including computer vision (Wang and Deng, 2018) (Csurka, 2017), natural language processing (Ramponi and Plank, 2020) and recommendation systems (Pan, 2016).

In the neural machine translation (NMT) task, domain adaptation is also in high demand since previous researches have shown that a general translation system trained on open-domain corpus often performs poorly in specific domains (Koehn and Knowles, 2017). Since in-domain parallel corpora are often insufficient, leveraging both in-domain and out-of-domain resources becomes important. Previous researches exhibit several sub-tasks: a semi-supervised task that uses a small-sized parallel corpus of the target domain (Jia and Zhang, 2020), and a supervised meta-learning task that requires small-sized in-domain support data-set to fast adapt to the target domain (Li et al., 2020; Sharaf et al., 2020). We argue that there are two

problems regarding current domain adaptation studies.

What is a domain? There exists no common agreement on what constitutes a domain in NMT. Many researches attribute domain difference to different vocabulary (Blitzer et al., 2006) or system of a document (McClosky, 2010). In practice, people use different predetermined dataset to represent different domains (Plank, 2016; Ramponi and Plank, 2020). Such practices often neglect the heterogeneity inside the corpus, which Plank (2016) put forward a theoretical notion called the *variety space*. A corpus can be seen as a set of sub-domains drawn from the underlying variety space. It is necessary to consider the smaller granularity of domains.

How to utilize the prior knowledge of inter-domain relationships and attributes? Domain adaptation in NMT usually tries to utilize linguistic similarities shared between source and target domains (Britz et al., 2017; Gu et al., 2019). However, there are other attributes outside the corpus to utilize. Although parallel or monolingual training corpus in NMT for a newly-published document is hard to get, we can easily find attributes that are shared with other domains (e.g., authors, categories, tags) and their relation with previous domains (e.g. recommendations). They can be regarded as prior knowledge for domain adaptation methods. In computer vision, Gebru et al. (2017) utilizes domain-level attributes and adopts a multi-task method for fine-grained domain adaptation. However, in NMT, we did not find similar researches.

The goal of our work is to solve the above three problems. The contribution comes in two folds: We build a multi-lingual dataset and provide inter-domain attributes as well as relationships. We believe this dataset is the first public dataset that focuses on fine-grained domains and is a good resource to explore the open-set domain adaptation and domain-level knowledge integration.

Field Name	Description
id	unique index of a game
lang	language of this sample
name	name of the game
developer	developers of the game
publisher	publishers of the game
category*	game category
type	game type
tag*	most voted tags by users
text	game introduction
recommended	ids of similar games

Table 1: Field description of STEAM. Fields marked with * are multilingual. The text field across two languages for the same game is aligned in sentence-level

We also propose a new approach that utilizes prior knowledge of inter-domain relationships by training an information network (Zhu et al., 2020b), whose output is treated as additional feature embeddings. To the best of our knowledge, this is the first work that tries to incorporate prior inter-domain relationships in an open-set domain adaptation task. Please find more information about the related work and baseline implementation in the appendix.

2 STEAM: A Multilingual Parallel Dataset for Domain Adaptation

Few datasets are designed specifically for domain adaptation and most datasets do not consider heterogeneity between samples. This corpus is a reaction to the growing importance of NMT in specific domains. While FDMT (Zhu et al., 2020a) make fine-grained domain corpus, it does not provide training data and only has English-Chinese pair. Most researches in multi-domain adaptation only consider a small subset of listed domains and cannot be generalized to the open-set domain adaptation settings. Domains are very distant from each other. No dataset describes inter-domain relationships and attributes.

This dataset is created based on Steam¹, a video game platform that contains thousands of game introductions in many languages. We produced parallel sentences based on HTML positions and LASER-based alignment (Artetxe and Schwenk, 2019). We also find relations between games based on information of game producers and user recommendations.

¹<https://store.steampowered.com/>

2.1 Data Collection & Analysis

We use the python framework Scrapy² to crawl text information from game introductions. The steps are in the appendix. 8 languages are chosen which are Chinese(ch), German(ge), English(en), Spanish(es), Russian(ru), Japanese(ja), Korean(ko) and French(fr). The field description is shown in 1. The data processing part is divided into two steps, monolingual and bilingual sentence matching. The details can be found in the appendix. The goal is to make the mapping of the same statement in different languages end up in the same neighborhood.

There are 10k games with 90k aligned parallel sentences after the cleaning. Table 3 in the appendix illustrates the detailed information. Codes and examples of bilingual parallel sentences can be found on our Google repository.³

3 Problem Formulation and Proposed methods

3.1 Task Definition

D_i is the i^{th} domain, which is a tuple of four elements (S_i, T_i, R_i, A_i) . S_i and T_i refer to paralleled source and target data in i^{th} domain. R_i refers to a set of other domains related to the i^{th} domain. A_i refers to domain-level attributes shared to all samples in this domain. The goal is to build a model that can generate best \hat{T}_i based on (S_i, R_i, A_i) . Our domain graph is defined as $G = (D, E)$, where E is the set of relationships between domains. We want to generate a feature embedding for each domain by modelling all domains in an undirected graph to avoid the scalability issue occurred.

3.1.1 Inter-Domain Relation Extraction

In this dataset, each domain is regarded as a node. The purpose of the inter-domain is to extract the relationship among nodes. There are several methods to learn nodes representations for inter-domain knowledge. Kipf and Welling (2017) generalizes neural models to work with structured graphs for semi-supervised learning. Tang et al. (2015) proposes LINE to generate node embeddings in a network by preserving the first second order proximities among nodes. We the use graph representation learning method named deep GRaph Contrastive rEpresentation learning (GRACE)(Zhu et al., 2020b). Two graph views are generated at

²<https://scrapy.org/>

³<https://drive.google.com/drive/folders/1wlj5sNETFL5hatC5sZfrHyAYf3ACK2nd?usp=sharing>

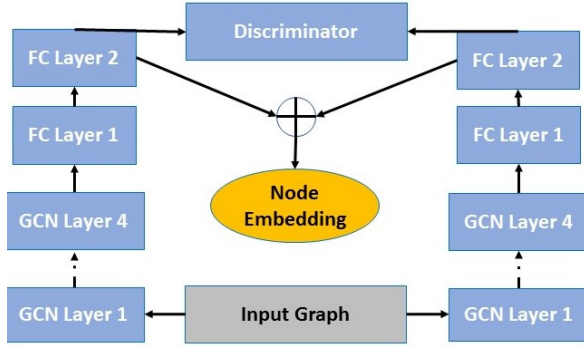


Figure 1: Illustration of the modified GRACE model

each iteration so that each node in the original graph has two embeddings. A discriminator is used to distinguish the two embeddings. We choose the method because it is stated that the framework is suitable for large-scale graphs and which fits our dataset. Also, the connections in our dataset are generated by the recommendation rules of steam and they are not clearly stated. It is rational to assume there are complicated relationships among the connected nodes. We extend the original framework by using a four-layer GCN network instead of two in the generation of node embeddings. Moreover, we increase the number of nodes in the MLP layers to 512 which fits the input of our NMT model. The final embedding of each node is obtained by adding the two hidden layers together and dividing them by two. The modified GRACE framework is shown in 1. The code can be found in Google Drive link⁴.

3.1.2 Handling Isolated Nodes

In the domain adaptation settings, there can be some domains not related to other domains. We call these domains isolated nodes. The proportion of isolated domains in our STEAM dataset is shown in the appendix. We cannot simply feed those isolated nodes into GRACE that takes as input a list of edges. We augment every node with a self-loop to have every node encoded and take an average of all the isolated node embedding vectors to get a final representation for all the isolated nodes.

3.1.3 Domain-Aware Encoder

We use the Transformer (Vaswani et al., 2017) as our backbone NMT model. We are motivated by Dou et al. (2019) to assign a domain-aware feature embedding on the model encoder side as an additional input representation to disentangle encoder

⁴<https://drive.google.com/drive/folders/1LCYP3WH489G-KfgnrgsKy8-HluXR26M?usp=sharing>

representation and learn domain knowledge. Although (Dou et al., 2019) trains different domains embeddings for each encoder layer, we found it not scalable in our setting. We combine three different ways to implement the domain-aware encoder, which is shown in figure2. Given the feature representation of in domain i as d_i , the j^{th} token of the sentence as x_j , the first method is to add the domain features only after the input embeddings. The formula is:

$$h_j^{(0)} = 0.5 * (W_c x_j + p_j + W_d d_i) \quad (1)$$

where $h_j^{(0)}$ is the input embeddings of x_j , W_c is the contextual embeddings of tokens, and p_j is the positional embeddings. W_d is a trainable linear projection that transforms the graph-based domain representations into domain embeddings for the NMT model. Note that domain is labeled at the sentence level, so tokens within a sentence share the same domain embeddings.

The second method is to add the projected domain features right after the Multi-Head Attention module of each layer. The formula is:

$$\tilde{H}^{(l)} = 0.5 * (\text{Multihead}(Q, K, V) + W_d^{(l)} d_i) \quad (2)$$

where $\tilde{H}^{(l)}$ is the intermediate hidden representation before feed-forward network and residual connection in each layer and $\text{Multihead}(Q, K, V)$ is the projected output of the concatenated attention heads defined in Transformer’s original paper.

The third method is to add domain embeddings at the end of each transformer encoder layer, which is similar to the method used in (Dou et al., 2019). The formula is

$$H^{(l)} = 0.5 * (\text{layer}_i(H^{(l-1)}) + W_d^{(l)} d_i) \quad (3)$$

The main difference between our proposed methods and methods introduced in Dou et al. (2019) and Michel and Neubig (2018) (noted as previous methods) is that they use one-hot encoded domain labels as the input of domain embedding/adaptor layers, while we use graph-based domain feature representation as the input. We claim two advantages of using graph-based domain feature representation.

Reduced number of parameters The method introduced in Dou et al. (2019) requires at least $r|D|$ more parameters and the method in Michel and Neubig (2018) requires at least $r(|D| + |V|)$

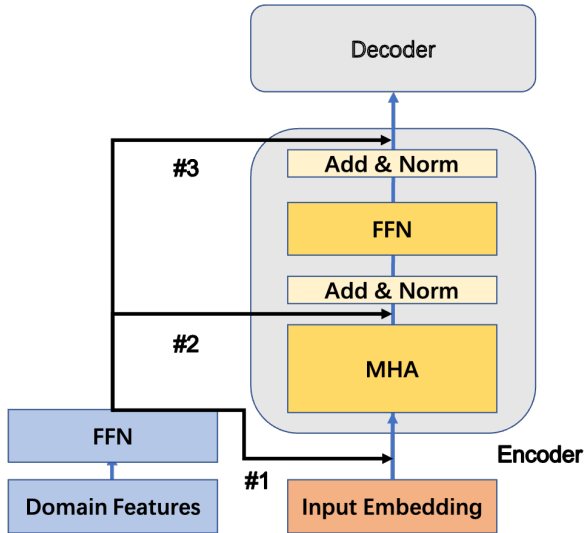


Figure 2: Architecture of the Domain-Aware Encoder. #X denotes the X^{th} method mentioned in 3.1.3

more parameters, where r is the hidden dimension of the transformer encoder, $|D|$ and $|V|$ are the size of domains and vocabularies which are linearly correlated with the size of domains. Our method only requires rk more parameters where k is the dimension of domain representation and is more practical when the size of domains is very large. In our dataset, the domain size for the en-zh pair is 9392, and we only use 5% more parameters compared with previous methods.

Inference on unseen domains Previous methods require the model to use either monolingual or bilingual corpus of all domains. Our methods can be used in an open-set domain situation where the corpus of the testing domains is not available in the training step based on domain similarities extracted from the graph-based domain feature representation.

4 Experiments

4.1 Datasets and Mixed fine-tuning

We choose ch-en, fr-en, es-en and ru-en as tested language pairs because they are four distinctively different language families. Domains in the test set are excluded in the training set to guarantee an open-set domain adaptation setting. We use mixed fine-tuning (Chu et al., 2017) as our training strategy. We collect UN Corpus⁵ as the out-of-domain data. We first pretrain a vanilla transformer model based on this out-of-domain data for 5 epochs. Then we sample 2 million sentences pairs

⁵<https://conferences.unite.un.org/UNCorpus/>

Method	zh-en	es-en	fr-en	ru-en
Baseline	17.22	28.31	23.18	14.28
Noise Embedding	15.80	27.41	12.86	7.58
DA Encoder	17.85	29.52	26.08	15.52

Table 2: BLEU in the test set for the baseline and proposed models. Bold text highlights the best results

from UN Corpus and over-sample another 2 million sentence pairs from the STEAM corpus for mixed fine-tuning. For our baseline model, pre-trained NMT models are directly loaded and fine-tuned in the mixed data, while for our proposed models, parameters related to the domain features are randomly initialized in the pretrain step. We also generate noise embeddings for comparison. The mean and variance of the noise embeddings are the same graph embeddings which leads to a similar distribution. The implementation details are in the appendix and our code can be found in our Google link⁶.

4.2 Main Results

Table 2 shows that our proposed models generally improve the overall translation performance in the test set. Regarding that our proposed method only requires very few extra parameters, this improvement is satisfying. Also, with noise embeddings, the results are worse than the baseline.

In Table 4 in the appendix, we can see although the BLEU improvement is relatively small, our proposed methods did improve the quality of translated sentence by correctly translating some key words.

5 Conclusion

This paper presents a new research topic on fine-grained multi-domain adaptation in NMT. We contribute a new dataset that focuses on the fine-grained domains and inter-domain relationships and proposes a novel method to utilize inter-domain relationships.

References

Mikel Artetxe and Holger Schwenk. 2019. [Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond](#). *Transactions of the Association for Computational Linguistics*, 7:597–610.

⁶<https://drive.google.com/drive/folders/1gRSv9B2ZMqzTKdDD9ge9tmDJkEdLTYWz?usp=sharing>

309	John Blitzer, Ryan McDonald, and Fernando Pereira.	Rumeng Li, Xun Wang, and Hong Yu. 2020. Metamt,	363
310	2006. Domain adaptation with structural correspon-	a meta learning method leveraging multiple domain	364
311	dence learning. In <i>Proceedings of the 2006 con-</i>	data for low resource machine translation. In <i>Pro-</i>	365
312	<i>ference on empirical methods in natural language</i>	<i>ceedings of the AAAI Conference on Artificial Intelli-</i>	366
313	<i>processing</i> , pages 120–128.	<i>gence</i> , volume 34, pages 8245–8252.	367
314	Denny Britz, Quoc Le, and Reid Pryzant. 2017. Effec-	David McClosky. 2010. <i>Any domain parsing: auto-</i>	368
315	tive domain mixing for neural machine translation.	<i>matic domain adaptation for natural language pars-</i>	369
316	In <i>Proceedings of the Second Conference on Machine</i>	<i>ing</i> . Ph.D. thesis, Brown University.	370
317	<i>Translation</i> , pages 118–126, Copenhagen, Denmark.		
318	Association for Computational Linguistics.	Paul Michel and Graham Neubig. 2018. Extreme adap-	371
319	Chenhui Chu, Raj Dabre, and Sadao Kurohashi. 2017.	tation for personalized neural machine translation.	372
320	An empirical comparison of domain adaptation meth-	In <i>Proceedings of the 56th Annual Meeting of the As-</i>	373
321	ods for neural machine translation. In <i>Proceedings</i>	<i>sociation for Computational Linguistics (Volume 2:</i>	374
322	<i>of the 55th Annual Meeting of the Association for</i>	<i>Short Papers)</i> , pages 312–318, Melbourne, Australia.	375
323	<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	Association for Computational Linguistics.	376
324	pages 385–391, Vancouver, Canada. Association for		
325	Computational Linguistics.	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan,	377
326	Gabriela Csurka. 2017. Domain adaptation for vi-	Sam Gross, Nathan Ng, David Grangier, and Michael	378
327	visual applications: A comprehensive survey. <i>arXiv</i>	Auli. 2019. fairseq: A fast, extensible toolkit for	379
328	<i>preprint arXiv:1702.05374.</i>	sequence modeling. In <i>Proceedings of the 2019 Con-</i>	380
329	Zi-Yi Dou, Junjie Hu, Antonios Anastasopoulos, and	<i>ference of the North American Chapter of the Associa-</i>	381
330	Graham Neubig. 2019. Unsupervised domain adap-	<i>tion for Computational Linguistics (Demonstrations)</i> ,	382
331	tation for neural machine translation with domain-	pages 48–53, Minneapolis, Minnesota. Association	383
332	aware feature embeddings. In <i>Proceedings of the</i>	for Computational Linguistics.	384
333	<i>2019 Conference on Empirical Methods in Natu-</i>		
334	<i>ral Language Processing and the 9th International</i>	Weike Pan. 2016. A survey of transfer learning for	385
335	<i>Joint Conference on Natural Language Processing</i>	collaborative recommendation with auxiliary data.	386
336	<i>(EMNLP-IJCNLP)</i> , pages 1417–1422.	<i>Neurocomputing</i> , 177:447–453.	387
337	T. Gebru, J. Hoffman, and L. Fei-Fei. 2017. Fine-	Barbara Plank. 2016. What to do about non-standard	388
338	grained recognition in the wild: A multi-task domain	(or non-canonical) language in nlp.	389
339	adaptation approach. In <i>2017 IEEE International</i>		
340	<i>Conference on Computer Vision (ICCV)</i> , pages 1358–	Alan Ramponi and Barbara Plank. 2020. Neural unsu-	390
341	1367.	pervised domain adaptation in NLP—A survey. In	391
342	Shuhao Gu, Yang Feng, and Qun Liu. 2019. Improving	<i>Proceedings of the 28th International Conference</i>	392
343	domain adaptation translation with domain invariant	<i>on Computational Linguistics</i> , pages 6838–6855,	393
344	and specific information. In <i>Proceedings of the 2019</i>	Barcelona, Spain (Online). International Committee	394
345	<i>Conference of the North American Chapter of the</i>	on Computational Linguistics.	395
346	<i>Association for Computational Linguistics: Human</i>		
347	<i>Language Technologies, Volume 1 (Long and Short</i>	Amr Sharaf, Hany Hassan, and Hal Daumé III. 2020.	396
348	<i>Papers)</i> , pages 3081–3091, Minneapolis, Minnesota.	Meta-learning for few-shot NMT adaptation. In <i>Pro-</i>	397
349	Association for Computational Linguistics.	<i>ceedings of the Fourth Workshop on Neural Genera-</i>	398
350	Chen Jia and Yue Zhang. 2020. Multi-cell composi-	<i>tion and Translation</i> , pages 43–53, Online. Associa-	399
351	tional LSTM for NER domain adaptation. In <i>Pro-</i>	tion for Computational Linguistics.	400
352	<i>ceedings of the 58th Annual Meeting of the Asso-</i>		
353	<i>ciation for Computational Linguistics</i> , pages 5906–	Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun	401
354	5917, Online. Association for Computational Lin-	Yan, and Qiaozhu Mei. 2015. Line: Large-scale	402
355	guistics.	information network embedding. In <i>Proceedings of</i>	403
356	Thomas N. Kipf and Max Welling. 2017. Semi-	<i>the 24th international conference on world wide web</i> ,	404
357	supervised classification with graph convolutional	pages 1067–1077.	405
358	networks. In <i>Proceedings of the 5th International</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	406
359	<i>Conference on Learning Representations, ICLR ’17.</i>	Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz	407
360	Philipp Koehn and Rebecca Knowles. 2017. Six chal-	Kaiser, and Illia Polosukhin. 2017. Attention is all	408
361	lenges for neural machine translation. <i>arXiv preprint</i>	<i>you need. arXiv preprint arXiv:1706.03762.</i>	409
362	<i>arXiv:1706.03872.</i>	Mei Wang and Weihong Deng. 2018. Deep visual	410
		domain adaptation: A survey. <i>Neurocomputing</i> ,	411
		312:135–153.	412
		Wenhao Zhu, Shujian Huang, Tong Pu, Xu Zhang,	413
		Jian Yu, Wei Chen, Yanfeng Wang, and Jiajun Chen.	414
		2020a. Fdmnt: A benchmark dataset for fine-grained	415
		domain adaptation in machine translation.	416

417 Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu
418 Wu, and Liang Wang. 2020b. Deep graph con-
419 trastive representation learning. *arXiv preprint*
420 *arXiv:2006.04131*.

421 **A Example Appendix**

422 **A.1 Data Crawling and Cleaning**

423 Setting request cookies such as birth time, age
424 checking, header, initial searching URL, and lan-
425 guage code.

426 Parsing information in the game search page
427 of each language including the URL, game name,
428 game ID, and information language.

429 Getting detailed information from the game page
430 including review text, game tags, game categories,
431 developer, publisher, and recommendations. The
432 recommendation is to recommend similar games
433 based on the current game.

434 Monolingual sentence cleaning focuses on filter-
435 ing abnormal sentences, words, and sentences that
436 are not in the labeled languages, using a language
437 identification tool called langid⁷. We use two steps
438 to build the bilingual parallel sentences and

439 First we find the passage that can be segmented
440 into small paragraphs based on different positions
441 in the website. Then we use Stanza⁸ to segment
442 sentences in the paragraph pair. For each paragraph
443 pair, if the numbers of segments for two paragraphs
444 are the same, each pair of segments is treated as
445 a pair of parallel sentences. Otherwise, we imple-
446 ment LASER⁹ to compute the sentence embedding
447 and use cosine similarity to pick out similar sen-
448 tence pairs. LASER is a toolkit to represent sen-
449 tences by vectors which are generated with respect
450 to both the input language and the NLP task.

451 **A.2 Implementation details**

452 We use fairseq (Ott et al., 2019) as our code base,
453 and use $1e - 4$ as the learning rate with the inverse
454 learning rate decay. All models share the same
455 vocabulary list, which is generated using sentence-
456 piece¹⁰ on the 4m mixed data mentioned above.

⁷<https://github.com/saffsd/langid.py>

⁸<https://stanfordnlp.github.io/stanza/>

⁹<https://github.com/facebookresearch/LASER>

¹⁰<https://github.com/google/sentencepiece>

	Sent. Num	Game Num	Avg. Length	Cat/Tag Num	Recommended
zh-en	82819	9392	29(CH)/15(EN)	29/418	12
ru-en	91719	10069	13(RU)/14(EN)	29/419	13
fr-en	93357	10550	17(FR)/15(EN)	30/425	12
es-en	92265	9264	16(ES)/15(EN)	29/419	12
ja-en	58445	10678	40(JA)/15(EN)	28/420	14
ko-en	69539	7982	35(KO)/14(EN)	28/423	12
de-en	93203	9504	17(DE)/15(EN)	29/420	13

Table 3: Statistics of our bilingual parallel dataset.¹¹

Reference	<p>You are afraid. You are forced to hide. Those who served you till recently have suddenly turned against you. Your only option is to hide, surviving with no hope, not understanding what’s going on around you. You are an outcast, an outlaw because you are human.</p>
Baseline	<p>You feel terrified and forced to crawl on. Now, they have branded your friends and family. Your only option is to survive with despair. You have no idea what you’re hiding because you are human, you’re the outcasts and outlaws.</p>
Our Model(DA Encoder)	<p>You’re scared and forced to hide. The robots that served you were suddenly betrayed. Your only option is hiding and survive with despair. You don’t know where they are or why they do this. Because you are a human being, so you are the bearer and the outcasts.</p>

Table 4: Comparison between the system output and the baseline.

Language	Total Games	Total Isolated Games
en	73306	26530
zh	15858	6459
ru	16654	6925
fr	18757	9583
es	17525	8342

Table 5: Statistics of isolated games for each language.