

Zero-Shot Paraphrase Generation with Multilingual Language Models

Yinpeng Guo¹, Yi Liao¹, Xin Jiang¹, Qing Zhang², Yibo Zhang², Qun Liu¹

¹Huawei Noah’s Ark Lab

²Intelligence Engineering Department, Huawei Consumer Business Group

{guo.yinpeng, liao.yi, jiang.xin, zhangqing49,
yibo.cheung, qun.liu}@huawei.com

Abstract

Leveraging multilingual parallel texts to automatically generate paraphrases has drawn much attention as size of high-quality paraphrase corpus is limited. Round-trip translation, also known as the *pivoting* method, is a typical approach to this end. However, we notice that the pivoting process involves multiple machine translation models and is likely to incur semantic drift during the two-step translations. In this paper, inspired by the Transformer-based language models, we propose a simple and unified paraphrasing model, which is purely trained on multilingual parallel data and can conduct *zero-shot* paraphrase generation in one step. Compared with the pivoting approach, paraphrases generated by our model is more semantically similar to the input sentence. Moreover, since our model shares the same architecture as GPT (Radford and Sutskever, 2018), we are able to pre-train the model on large-scale unparallel corpus, which further improves the fluency of the output sentences. In addition, we introduce the mechanism of denoising auto-encoder (DAE) to improve diversity and robustness of the model. Experimental results show that our model surpasses the pivoting method in terms of relevance, diversity, fluency and efficiency.

1 Introduction

Paraphrasing is to express the same meaning using different expressions. Paraphrase generation plays an important role in various natural language processing (NLP) tasks such as response diversification in dialogue system, query reformulation in information retrieval, and data augmentation in machine translation. Recently, models based on Seq2Seq learning (Ilya Sutskever, 2014) have achieved the state-of-the-art results on paraphrase generation. Most of these models (Prakash et al., 2016; Ziqiang Cao, 2017; Ankush Gupta, 2018;

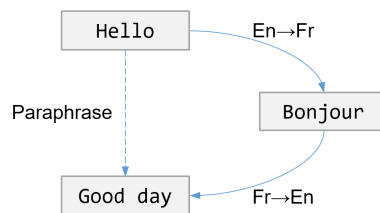


Figure 1: Paraphrase generation via round-trip translation.

Zichao Li, 2018, 2019) focus on training the paraphrasing models based on a paraphrase corpus, which contains a number of pairs of paraphrases. However, high-quality paraphrases are usually difficult to acquire in practice, which becomes the major limitation of these methods. Therefore, we focus on *zero-shot paraphrase generation* approach in this paper, which aims to generate paraphrases without requiring a paraphrase corpus.

A natural choice is to leverage the bilingual or multilingual parallel data used in machine translation, which are of great quantity and quality. The basic assumption is that if two sentences in one language (e.g., English) have the same translation in another language (e.g., French), they are assumed to have the same meaning, i.e., they are paraphrases of each other. Therefore, one typical solution for paraphrasing in one language is to *pivot* over a translation in another language. Specifically, it is implemented as the *round-trip* translation, where the input sentence is translated into a foreign sentence, then back-translated into a sentence in the same language as input (Jonathan Mallinson and Lapata, 2017). The process is shown in Figure 1. Apparently, two machine translation systems (English→French and French←English) are needed to conduct the generation of a paraphrase.

Although the pivoting approach works in general, there are several intrinsic defects. First,

the round-trip system can hardly explore all the paths of paraphrasing, since it is pivoted through the finite intermedia outputs of a translation system. More formally, let Z denote the meaning representation of a sentence X , and finding paraphrases of X can be treated as sampling another sentence Y conditioning on the representation Z . Ideally, paraphrases should be generated by following $P(Y|X) = \int_Z P(Y|Z)P(Z|X)dZ$, which is marginalized over all possible values of Z . However, in the round-trip translation, only one or several Z s are sampled from the machine translation system $P(Z|X)$, which can lead to an inaccurate approximation of the whole distribution and is prone to the problem of *semantic drift* due to the sampling variances. Second, the results are determined by the pre-existing translation systems, and it is difficult to optimize the pipeline end-to-end. Last, the system is not efficient especially at the inference stage, because it needs two rounds of translation decoding.

To address these issues, we propose a single-step zero-shot paraphrase generation model, which can be trained on machine translation corpora in an end-to-end fashion. Unlike the pivoting approach, our proposed model does not involve explicit translation between multiple languages. Instead, it directly learns the paraphrasing distribution $P(Y|X)$ from the parallel data sampled from $P(Z|X)$ and $P(Y|Z)$. Specifically, we build a Transformer-based (Ashish Vaswani, 2017) language model, which is trained on the concatenated bilingual parallel sentences with language indicators. At inference stage, given a input sentence in a particular language, the model is guided to generate sentences in the same language, which are deemed as paraphrases of the input. Our model is simple and compact, and can empirically reduce the risk of semantic drift to a large extent. Moreover, we can initialize our model with generative pre-training (GPT) (Radford and Sutskever, 2018) on monolingual data, which can benefit the generation in low-resource languages. Finally, we borrow the idea of denoising auto-encoder (DAE) to further enhance robustness in paraphrase generation.

We conduct experiments on zero-shot paraphrase generation task, and find that the proposed model significantly outperforms the pivoting approach in terms of both automatic and human evaluations. Meanwhile, the training and inference cost are largely reduced compared to the pivot-based meth-

ods which involves multiple systems.

2 Methodology

2.1 Transformer-based Language Model

Transformer-based language model (TLM) is a neural language model constructed with a stack of Transformer *decoder* layers (Ashish Vaswani, 2017). Given a sequence of tokens, TLM is trained with maximizing the likelihood:

$$L(X) = \sum_{i=1}^n \log P(x_i|x_{1,\dots,i-1}; \theta) \quad (1)$$

where $X = [x_1, x_2, \dots, x_n]$ is a sentence in a language (e.g., English), and θ denotes the parameters of the model. Each Transformer layer is composed of multi-head self-attention, layer normalization and a feed-forward network. We refer reader to the original paper for details of each component. Formally, the decoding probability is given by

$$\begin{aligned} [e_1, \dots, e_{i-1}] &= [W_e x_1 + p_1, \dots, W_e x_{i-1} + p_{i-1}], \\ [h_1, \dots, h_{i-1}] &= \text{Transformer}([e_1, \dots, e_{i-1}]), \\ P(x_i|x_{1,\dots,i-1}; \theta) &= \text{Softmax}(W_o h_{i-1}), \end{aligned} \quad (2)$$

where x_i denotes the token embedding, p_i denote the positional embedding and h_i denotes the output states of the i -th token, and W_e and W_o are the input and output embedding matrices.

Although TLM is normally employed to model monolingual sequences, there is no barrier to utilize TLM to model sequences in multiple languages. In this paper, inspired by Lample and Conneau (2019), we concatenate pairs of sentences from bilingual parallel corpora (e.g., English→French) as training instances to the model. Let X and Y denote the parallel sentences in two different languages, the training objective becomes

$$\begin{aligned} L(X, Y) &= \sum_{i=1}^n \log P(x_i|x_{1,\dots,i-1}; \theta) \\ &+ \sum_{j=1}^m \log P(y_j|x_{1,\dots,n}, y_{1,\dots,j-1}; \theta). \end{aligned} \quad (3)$$

This bilingual language model can be regarded as the decoder-only model compared to the traditional encoder-decoder model. It has been proved to work effectively on monolingual text-to-text generation tasks such as summarization (Peter J. Liu, 2018). The advantages of such architecture include less

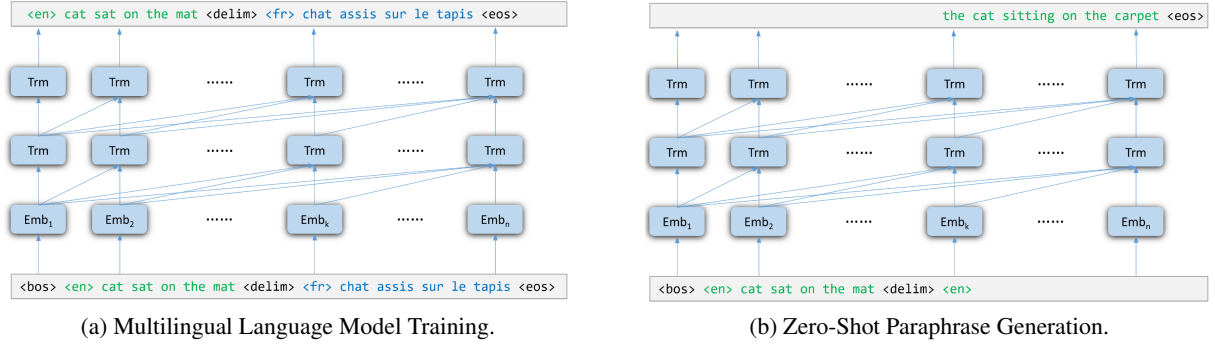


Figure 2: Paraphrase generation via multilingual language model training.

model parameters, easier optimization and potential better performance for longer sequences. Furthermore, it naturally integrates with language models pre-training on monolingual corpus.

For each input sequence of concatenated sentences, we add special tokens $\langle bos \rangle$ and $\langle eos \rangle$ at the beginning and the end, and $\langle delim \rangle$ in between the sentences. Moreover, at the beginning of each sentence, we add a special token as its language identifier, for instance, $\langle en \rangle$ for English, $\langle fr \rangle$ for French. One example of English \rightarrow French training sequence is “ $\langle bos \rangle \langle en \rangle$ cat sat on the mat $\langle delim \rangle \langle fr \rangle$ chat assis sur le tapis $\langle eos \rangle$ ”.

At inference stage, the model predicts the next word as the conventional auto-regressive model:

$$\hat{y}_j \sim P(y_j | x_{1,\dots,n}, y_{1,\dots,j-1}; \theta) \quad (4)$$

2.2 Zero-shot Paraphrase Generation

We train the bilingual language model on multiple bilingual corpora, for example, English \leftrightarrow French and German \leftrightarrow Chinese. Once the language model has been trained, we can conduct zero-shot paraphrase generation based on the model. Specifically, given an input sentence that is fed into the language model, we set the output language identifier the same as input, and then simply conduct decoding to generate paraphrases of the input sentence.

Figure 2 illustrates the training and decoding process of our model. In the training stage, the model is trained to sequentially generate the input sentence and its translation in a specific language. Training is conducted in the way of teacher-forcing. In the decoding stage, after an English sentence “ $\langle bos \rangle \langle en \rangle$ cat sat on the mat $\langle delim \rangle$ ” is fed to the model, we intentionally set the output language identifier as “ $\langle en \rangle$ ”, in order to guide the model to continue to generate English words. At the same

time, since the model has been trained on translation corpus, it implicitly learns to keep the semantic meaning of the output sentence the same as the input. Accordingly, the model will probably generate the paraphrases of the input sentence, such as “the cat sitting on the carpet $\langle eos \rangle$ ”.

It should be noted our model can obviously be trained on parallel paraphrase data without any modification. But in this paper, we will mainly focus on the research and evaluation in the zero-shot learning setting.

In the preliminary experiments of zero-shot paraphrasing, we find the model does not perform consistently well and sometimes fails to generate the words in the correct language as indicated by the language identifier. Similar phenomenon has been observed in the research of zero-shot neural machine translation (Sestorain et al., 2018; Arivazhagan et al., 2019; Jiatao Gu, 2019), which is referred as the *degeneracy* problem by Jiatao Gu (2019). To address these problems in zero-shot paraphrase generation, we propose several techniques to improve the quality and diversity of the model as follows.

2.2.1 Language Embeddings

The language identifier prior to the sentence does not always guarantee the language of the sequences generated by the model. In order to keep the language consistency, we introduce language embeddings, where each language is assigned a specific vector representation. Supposing that the language embedding for the i -th token in a sentence is a_i , we concatenate the language embedding with the Transformer output states and feed it to the softmax layer for predicting each token:

$$P(y_j | x_{1,\dots,n}, y_{1,\dots,j-1}; \theta) = \text{Softmax}(W_o[h_j, a_j]) \quad (5)$$

We empirically demonstrate that the language embedding added to each tokens can effectively guide the model to generate sentences in the required language. Note that we still let the model to learn the output distribution for each language rather than simply restricting the vocabularies of output space. This offers flexibility to handle coding switching cases commonly seen in real-world data, e.g., English words could also appear in French sentences.

2.2.2 Pre-Training on Monolingual Corpora

Language model pre-training has shown its effectiveness in language generation tasks such as machine translation, text summarization and generative question answering (Radford et al., 2019; Dong et al., 2019; Song et al., 2019). It is particularly helpful to the low/zero-resource tasks since the knowledge learned from large-scale monolingual corpus can be transferred to downstream tasks via the pre-training-then-fine-tuning approach. Since our model for paraphrase generation shares the same architecture as the language model, we are able to pre-train the model on massive monolingual data.

Pre-training on monolingual data is conducted in the same way as training on parallel data, except that each training example contains only one sentence with the beginning/end of sequence tokens and the language identifier. The language embeddings are also employed. The pre-training objective is the same as Equation (1).

In our experiments, we first pre-train the model on monolingual corpora of multiple languages respectively, and then fine-tune the model on parallel corpora.

2.2.3 Denoising Auto-Encoder

We adopt the idea of denoising auto-encoder (DAE) to further improve the robustness of our paraphrasing model. DAE is originally proposed to learn intermediate representations that are robust to partial corruption of the inputs in training auto-encoders (Pascal Vincent, 2008). Specifically, the initial input X is first partially corrupted as \tilde{X} , which can be treated as sampling from a noise distribution $\tilde{X} \sim q(\tilde{X}|X)$. Then, an auto-encoder is trained to recover the original X from the noisy input \tilde{X} by minimizing the reconstruction error. In the applications of text generation (Freitag and Roy, 2018) and machine translation (Yunsu Kim, 2018), DAE has shown to be able to learn representations that are more robust to input noises and also

generalize to unseen examples.

Inspired by (Yunsu Kim, 2018), we directly inject three different types of noises into input sentence that are commonly encountered in real applications.

1) *Deletion*: We randomly delete 1% tokens from source sentences, for example, “*cat sat on the mat*” \mapsto “*cat on the mat*.”

2) *Insertion*: We insert a random token into source sentences in 1% random positions, for example, “*cat sat on the mat*” \mapsto “*cat sat on red the mat*.”

3) *Reordering*: We randomly swap 1% tokens in source sentences, and keep the distance between tokens being swapped within 5. “*cat sat on the mat*” \mapsto “*mat sat on the cat*.”

By introducing such noises into the input sentences while keeping the target sentences clean in training, our model can be more stable in generating paraphrases and generalisable to unseen sentences in the training corpus. The training objective with DAE becomes

$$\begin{aligned} L(X, Y) &= L(X) + L(Y|\tilde{X})q(\tilde{X}|X) \\ &= \sum_{i=1}^n \log P(x_i|x_{1,\dots,i-1}; \theta) \\ &\quad + \sum_{j=1}^m \log P(y_j|\tilde{x}_{1,\dots,n}, y_{1,\dots,j-1}; \theta). \end{aligned} \quad (6)$$

Once the model is trained, we generate paraphrases of a given sentence based on $P(Y|X; \theta)$.

3 Experiments

3.1 Datasets

We adopt the mixture of two multilingual translation corpus as our training data: MultiUN (Andreas Eisele, 2010) and OpenSubtitles (Pierre Li-son, 2016). MultiUN consists of 463,406 official documents in six languages, containing around 300M words for each language. OpenSubtitles is a corpus consisting of movie and TV subtitles, which contains 2.6B sentences over 60 languages. We select four shared languages of the two corpora: English, Spanish, Russian and Chinese. Statistics of the training corpus are shown in Table 1. Sentences are tokenized by Wordpiece as in BERT. A multilingual vocabulary of 50K tokens is used. For validation and testing, we randomly sample 10000 sentences respectively from each language pair. The rest data are used for training. For mono-

Table 1: Statistics of training data (#sentences).

	En \leftrightarrow Es	En \leftrightarrow Ru	En \leftrightarrow Zh	Es \leftrightarrow Ru	Es \leftrightarrow Zh	Ru \leftrightarrow Zh
OpenSubtitles	11.7M	11.7M	11.2M	10.5M	8.5M	9.6M
MultiUN	11.4M	11.7M	9.6M	10.6M	9.8M	9.6M
Total	23.1M	23.4M	20.8M	21.1M	18.3M	19.2M

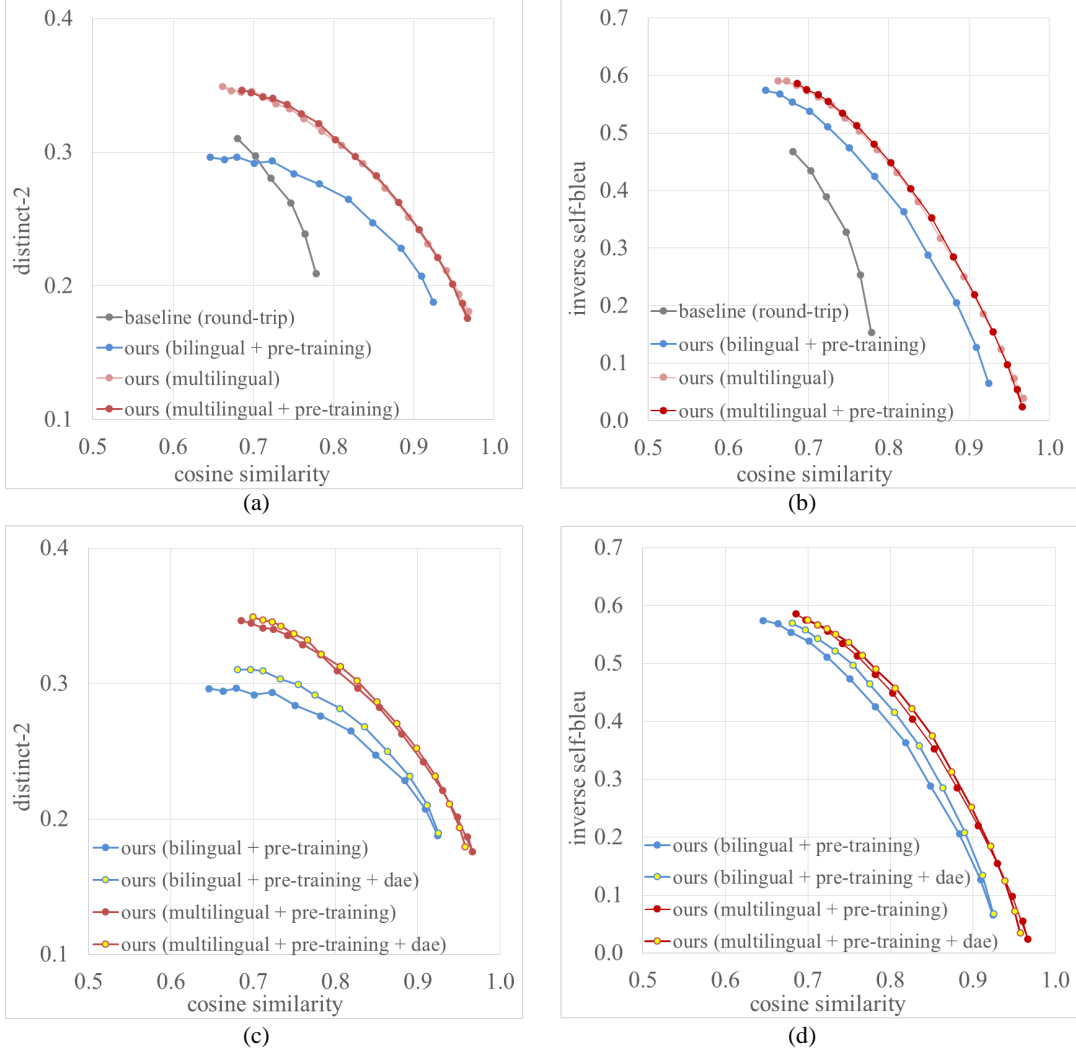


Figure 3: Automatic evaluation: (a)(c) Distinct-2 versus Relevance; (b)(d) Inverse Self-BLEU versus Relevance.

lingual pre-training, we use English Wikipedia¹ corpus, which contains 2,500M words.

3.2 Experimental Settings

We implement our model in Tensorflow (Abadi et al., 2016). The size of our Transformer model is identical to BERT-base (Jacob Devlin, 2019). The model is constituted by 12 layers of Transformer blocks. Number of dimension of token embedding, position embedding and transformer hidden

state are 768, while that of states in position-wise feed-forward networks are 3072. The number of attention heads is 12. Models are trained using Adam optimization (Diederik P. Kingma) with a learning rate up to $1e-4$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $L2$ weight decay of 0.01. We use top-k truncated random sampling strategy for inference that only sample from k candidate words with highest probabilities.

Throughout our experiments, we train and evaluate two models for paraphrase generation: the bilingual model and the multilingual model.

¹<https://dumps.wikimedia.org/enwiki/latest/>

The bilingual models are trained only with English \leftrightarrow Chinese, while the multilingual models are trained with all the data between the four languages. The round-trip translation baseline is based on the Transformer-based neural translation model.

3.3 Automatic Evaluation

We evaluate the relevance between input and generated paraphrase as well as the diversity among multiple generated paraphrases from the same input. For relevance, we use the cosine similarity between the sentential representations (Chia-Wei Liu, 2016). Specifically, we use the Glove-840B embeddings (Jeffrey Pennington, 2014) for word representation and Vector Extrema (Chia-Wei Liu, 2016) for sentential representation. For generation diversity, we employ two evaluation metrics: Distinct-2² and *inverse* Self-BLEU (defined as: $1 - \text{Self-BLEU}$) (Yaoming Zhu, 2018). Larger values of Distinct-2 and *inverse* Self-BLEU indicate higher diversity of the generation.

For each model, we draw curves in Figure 3 with the aforementioned metrics as coordinates, and each data-point is obtained at a specific sampling temperature. Since a good paraphrasing model should generate both relevant and diverse paraphrases, the model with curve lying towards the up-right corner is regarded as with good performance.

3.3.1 Comparison with Baseline

First we compare our models with the conventional pivoting method, i.e., round-trip translation. As shown in Figure 3 (a)(b), either the bilingual or the multilingual model is better than the baseline in terms of relevance and diversity in most cases. In other words, with the same generation diversity (measured by both Distinct-2 and Self-BLEU), our models can generate paraphrase with more semantic similarity to the input sentence.

Note that in Figure 3 (a), there is a cross point between the curve of the bilingual model and the baseline curve when relevance is around 0.71. We particularly investigate generated paraphrases around this point and find that the baseline actually achieves better relevance when Distinct-2 is at a high level (>0.3). It means our bilingual model is semantically drifting faster than the baseline model as the Distinct-2 diversity increases. The round-trip

translation performs two-round of supervised translations, while the zero-shot paraphrasing performs single-round unsupervised ‘translation’ (paraphrasing). We suspect that the unsupervised paraphrasing can be more sensitive to the decoding strategy. It also implies the latent, language-agnostic representation may be not well learned in our bilingual model. While on the other hand, our multilingual model alleviate this insufficiency. We further verify and analyze it as follows.

3.3.2 Multilingual Models

As mentioned above, our bilingual model can be unstable in some cases due to the lack of a well-learned language-agnostic semantic representation. A natural method is to introduce multilingual corpus, which consists of various translation directions. Training over multilingual corpus forces the model to decouple the language type and semantic representation.

Empirical results shows that our multilingual model performs significantly better than the bilingual model. The red and blue curves in Figure 3 (a)(b) demonstrates a great improvement of our multilingual model over the bilingual model. In addition, the multilingual model also significantly outperforms the baseline in the setting with the reasonable relevance scores.

3.3.3 Denoising Auto-Encoder

To verify the effectiveness of DAE in our model, various experiments with different hyperparameters were conducted. We find that DAE works the best when uniformly perturbing input sentences with probability 0.01, using only *Deletion* and *Reordering* operations. We investigate DAE over both bilingual and multilingual models as plotted in Figure 3 (c)(d). Curves with the yellow circles represent models with DAE training.

Results in the Figure 3 (c)(d) demonstrate positive effects of DAE in either bilingual or multilingual models. It is worth to note that, while DAE have marginal impact on multilingual model, it improves bilingual model significantly. This is an evidence indicating that DAE can improve the model in learning a more robust representation.

More specifically, since *Deletion* forces model to focus on sentence-level semantics rather than word-level meaning while *Reordering* forces model to focus more on meaning rather than their positions, it would be more difficult for a model to learn shortcuts (e.g. copy words). In other words, DAE

²<https://github.com/neural-dialogue-metrics/Distinct-N>

improves models’ capability in extracting deep semantic representation, which has a similar effect to introducing multilingual data.

3.3.4 Monolingual Pre-Training

Table 2: Log-probabilities of the generated sentences. \checkmark and \times symbols denote learning with or without pre-training respectively, **bold** font denotes greater values.

Model	Sampling	Pre-Training	Log-Prob
Multilingual	greedy, temp=1	\checkmark	-0.1427
		\times	-0.1428
	top-3, temp=1	\checkmark	-0.1425
		\times	-0.1448
	top-3, temp=1.5	\checkmark	-0.1420
		\times	-0.1425
Bilingual	greedy, temp=1	\checkmark	-0.1472
		\times	-0.1484
	top-3, temp=1	\checkmark	-0.1487
		\times	-0.1502
	top-3, temp=1.5	\checkmark	-0.1461
		\times	-0.1506

As shown in Figure 3 (a)(b), the model with language model pre-training almost performs equally to its contemporary without pre-training. However, evaluations on fluency uncover the value of pre-training. We evaluate a group of models over our test set in terms of fluency, using a n-grams language model³ trained on 14k public domain books.

As depicted in Table 2, models with language model pre-training stably achieves greater log-probabilities than the model without pre-training. Namely, language model pre-training brings better fluency.

3.4 Human Evaluation

200 sentences are sampled from our test set for human evaluation. The human evaluation guidance generally follows that of (Zichao Li, 2018) but with a compressed scoring range from [1, 5] to [1, 4]. We recruit five human annotators to evaluate models in semantic relevance and fluency. A test example consists of one input sentence, one generated sentence from baseline model and one generated sentence from our model. We randomly permute a pair of generated sentences to reduce annotators’ bias on a certain model. Each example is evaluated by two annotators.

As shown in Table 3, our method outperforms the baseline in both relevance and fluency significantly. We further calculate agreement (Cohen’s kappa) between two annotators.

³<http://www.openslr.org/11/>

Table 3: Human evaluation results.

Model	Relevance	Fluency	Agreement
Round-trip	2.72	3.61	0.36
Multilingual (ours)	3.43	3.75	0.35

Both round-trip translation and our method performs well as to fluency. But the huge gap of relevance between the two systems draw much attention of us. We investigate the test set in details and find that round-trip approach indeed generate more noise as shown in case studies.

3.5 Case Studies

We further study some generated cases from different models. All results in Table 4 are generated over our test set using randomly sampling. For both baseline and multilingual model, we tune their sampling temperatures to control the Distinct-2 and the inverse Self-BLEU at 0.31 and 0.47 respectively.

In the case studies, we find that our method usually generates sentences with better relevance to source inputs, while the round-trip translation method can sometimes run into serious semantic drift. In the second case, our model demonstrates a good feature that it maintains the meaning and even a proper noun *guide* unchanged while modifies the source sentence by both changing and reordering words. This feature may be introduced by DAE perturbation strategies which improves model’s robustness and diversity simultaneously. These results evidence that our methods outperforms the baseline in both relevance and diversity.

4 Related Work

Generating paraphrases based on deep neural networks, especially Seq2Seq models, has become the mainstream approach. A majority of neural paraphrasing models tried to improve generation quality and diversity with high-quality paraphrase corpora. Prakash et al. (2016) starts a deep learning line of paraphrase generation through introducing stacked residual LSTM network. A word constraint model proposed by Ziqiang Cao (2017) improves both generation quality and diversity. Ankush Gupta (2018) adopts variational auto-encoder to further improve generation diversity. Zichao Li (2018) utilize neural reinforcement learning and adversarial training to promote generation quality. Zichao Li (2019) decompose paraphrase generation into phrase-level and sentence-level.

Table 4: Case studies. For each input source, we randomly sample three paraphrases for comparison.

Source	I guess I kinda felt insignificant.
Round-trip	<i>I think I just don't feel right about that.</i> <i>I guess I'm a little uncomfortable.</i> <i>I think I'm a little bit of a problem right now.</i>
Multilingual (ours)	<i>I guess I was feeling a bit insignificant.</i> <i>I guess I felt some sorts of insignificant.</i> <i>I guess I kind of felt insignificant.</i>
Source	This site will make better use of the guide and will increase its distribution.
Round-trip	<i>The site would make better use of the guidelines and would be expanded.</i> <i>The site will make the best use of guides and expand them.</i> <i>This site would have made use more of the guidelines and could be expanded to its distribution.</i>
Multilingual (ours)	<i>This web site will make better use of the guide and will increase its dissemination.</i> <i>This site will better utilize the guide, and will improve its distribution.</i> <i>The web site is going to make the guide's use more efficient and its distribution will grow.</i>
Source	That's how eric got the passcodes.
Round-trip	<i>Then eric has a code.</i> <i>Then eric has the codes.</i> <i>Then erik'll have the codes.</i>
Multilingual (ours)	<i>That's the way eric got the password codes.</i> <i>That's how eric got passwords.</i> <i>That's where eric gets the passcodes.</i>

Several works tried to generate paraphrases from monolingual non-parallel or translation corpora. Lilin Zhang (2016) exploits Markov Network model to extract paraphrase tables from monolingual corpus. Quirk, Brockett, and Dolan (2004), Wubben, van den Bosch, and Krahmer (2010) and Wubben, van den Bosch, and Krahmer (2014) create paraphrase corpus through clustering and aligning paraphrases from crawled articles or headlines. With parallel translation corpora, pivoting approaches such round-trip translation (Jonathan Mallinson and Lapata, 2017) and back-translation (John Wieting, 2018) are explored.

However, to the best knowledge of us, none of these paraphrase generation models has been trained directly from parallel translation corpora as a single-round end-to-end model.

5 Conclusions

In this work, we have proposed a Transformer-based model for zero-shot paraphrase generation, which can leverage huge amount of off-the-shelf translation corpora. Moreover, we improve generation fluency of our model with language model pre-training. Empirical results from both automatic and human evaluation demonstrate that our model surpasses the conventional pivoting approaches in terms of relevance, diversity, fluency and efficiency. Nevertheless, there are some interesting directions to be explored. For instance, how to obtain a better latent semantic representation with multi-modal data and how to further improve the generation di-

versity without sacrificing relevance. We plan to strike these challenging yet valuable problems in the future.

References

- Abadi, M.; Barham, P.; Chen, J.; Chen, Z.; Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Irving, G.; Isard, M.; Kudlur, M.; Levenberg, J.; Monga, R.; Moore, S.; Murray, D. G.; Steiner, B.; Tucker, P.; Vasudevan, V.; Warden, P.; Wicke, M.; Yu, Y.; and Zheng, X. 2016. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, 265–283.
- Andreas Eisele, Y. C. 2010. Multiun: A multilingual corpus from united nation documents. In *LREC*.
- Ankush Gupta, Arvind Agarwal, P. S. P. R. 2018. A deep generative framework for paraphrase generation. In *AAAI*.
- Arivazhagan, N.; Bapna, A.; Firat, O.; Aharoni, R.; Johnson, M.; and Macherey, W. 2019. The missing ingredient in zero-shot neural machine translation. *arXiv preprint arXiv:1903.07091*.
- Ashish Vaswani, Noam Shazeer, N. P. J. U. L. J. A. N. G. L. K. I. P. 2017. Attention is all you need. In *NIPS*.
- Chia-Wei Liu, Ryan Lowe, I. V. S. M. N. L. C. J. P. 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *EMNLP*.
- Diederik P. Kingma, J. B. Adam: A method for stochastic optimization. In *ICLR*.

- Dong, L.; Yang, N.; Wang, W.; Wei, F.; Liu, X.; Wang, Y.; Gao, J.; Zhou, M.; and Hon, H.-W. 2019. Unified language model pre-training for natural language understanding and generation. *arXiv preprint arXiv:1905.03197*.
- Freitag, M., and Roy, S. 2018. Unsupervised natural language generation with denoising autoencoders. In *EMNLP*.
- Ilya Sutskever, Oriol Vinyals, Q. V. L. 2014. Sequence to sequence learning with neural networks. In *NIPS*.
- Jacob Devlin, Ming-Wei Chang, K. L. K. T. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- Jeffrey Pennington, Richard Socher, C. D. M. 2014. Glove: Global vectors for word representation. In *EMNLP*.
- Jiatao Gu, Yong Wang, K. C. V. O. L. 2019. Improved zero-shot neural machine translation via ignoring spurious correlations. In *ACL*.
- John Wieting, K. G. 2018. Paranzmt-50m - pushing the limits of paraphrastic sentence embeddings with millions of machine translations. In *ACL*.
- Jonathan Mallinson, R. S., and Lapata, M. 2017. Paraphrasing revisited with neural machine translation. In *EACL*.
- Lample, G., and Conneau, A. 2019. Cross-lingual language model pretraining. *arXiv preprint arXiv:1901.07291*.
- Lilin Zhang, Zhen Weng, W. X. J. W. Z. C. Y. T. M. L. M. W. 2016. Extract domain-specific paraphrase from monolingual corpus for automatic evaluation of machine translation. In *MT*.
- Pascal Vincent, Hugo Larochelle, Y. B. P.-A. M. 2008. Extracting and composing robust features with denoising autoencoders. In *ICML*.
- Peter J. Liu, Mohammad Saleh, E. P.-B. G. R. S. L. K. N. S. 2018. Generating wikipedia by summarizing long sequences. In *ICLR*.
- Pierre Lison, J. T. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles. In *LREC*.
- Prakash, A.; Hasan, S. A.; Lee, K.; Datla, V.; Qadir, A.; Liu, J.; and Farri, O. 2016. Neural paraphrase generation with stacked residual LSTM networks. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 2923–2934. Osaka, Japan: The COLING 2016 Organizing Committee.
- Quirk, C.; Brockett, C.; and Dolan, W. 2004. Monolingual machine translation for paraphrase generation. In *EMNLP*, 142–149. Barcelona, Spain: Association for Computational Linguistics.
- Radford, A., and Sutskever, I. 2018. Improving language understanding by generative pre-training. In *arxiv*.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language models are unsupervised multitask learners. *OpenAI Blog* 1(8).
- Sestorain, L.; Ciaramita, M.; Buck, C.; and Hofmann, T. 2018. Zero-shot dual machine translation. *arXiv preprint arXiv:1805.10338*.
- Song, K.; Tan, X.; Qin, T.; Lu, J.; and Liu, T.-Y. 2019. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*.
- Wubben, S.; van den Bosch, A.; and Krahmer, E. 2010. Paraphrase generation as monolingual translation: Data and evaluation. In *INLG, INLG '10*, 203–207. Stroudsburg, PA, USA: Association for Computational Linguistics.
- Wubben, S.; van den Bosch, A.; and Krahmer, E. 2014. Creating and using large monolingual parallel corpora for sentential paraphrase generation. In *LREC*, 4292–4299. Reykjavik, Iceland: European Languages Resources Association (ELRA).
- Yaoming Zhu, Sidi Lu, L. Z. J. G. W. Z. J. W. Y. Y. 2018. Txygen: A benchmarking platform for text generation models. In *SIGIR*.
- Yunsu Kim, Jiahui Geng, H. N. 2018. Improving unsupervised word-by-word translation with language model and denoising autoencoder. In *EMNLP*.
- Zichao Li, Xin Jiang, L. S. H. L. 2018. Paraphrase generation with deep reinforcement learning. In *EMNLP*.
- Zichao Li, Xin Jiang, L. S. Q. L. 2019. Decomposable neural paraphrase generation. In *ACL*.
- Ziqiang Cao, Chuwei Luo, W. L. S. L. 2017. Joint copying and restricted generation for paraphrase. In *AAAI*.