

# Improving the Diversity of Unsupervised Paraphrasing with Embedding Outputs

**Monisha Jegadeesan\***

Indian Institute of Technology Madras  
monishaj.65@gmail

**John Wieting\***

Carnegie Mellon University  
jwieting@cs.cmu.edu

**Sachin Kumar**

Carnegie Mellon University  
sachink@cs.cmu.edu

**Yulia Tsvetkov**

University of Washington  
yuliats@cs.washington.edu

## Abstract

We present a novel technique for zero-shot paraphrase generation. The key contribution is an end-to-end multilingual paraphrasing model that is trained using translated parallel corpora to generate paraphrases into “meaning spaces” – replacing the final softmax layer with word embeddings. This architectural modification, plus a training procedure that incorporates an autoencoding objective, enables effective parameter sharing across languages for more fluent monolingual rewriting, and facilitates fluency and diversity in generation. Our continuous-output paraphrase generation models outperform zero-shot paraphrasing baselines, when evaluated on two languages using a battery of computational metrics as well as in human assessment.<sup>1</sup>

## 1 Introduction

Paraphrasing aims to rewrite text while preserving its meaning and achieving a different surface realization. It is an eminently practical task, useful in educational applications (Inui et al., 2003; Petersen and Ostendorf, 2007; Pavlick and Callison-Burch, 2016; Xu et al., 2016), information retrieval (Duboue and Chu-Carroll, 2006; Harabagiu and Hickl, 2006; Fader et al., 2014), in dialogue systems (Yan et al., 2016), as well as for data augmentation in a plethora of other tasks (Berant and Liang, 2014; Romano et al., 2006; Fadaee et al., 2017; Jin et al., 2018; Hou et al., 2018).

Generating diverse and coherent paraphrases is a difficult task. Unlike in machine translation, where naturally occurring parallel data in the form of translated news, books and talks are available in abundance on the web, naturally occurring paraphrase corpora are scarce. Most common

approaches to paraphrasing are based on translation, in the form of bilingual pivoting (Mallinson et al., 2017a,b) or back-translation (Wieting and Gimpel, 2018; Hu et al., 2019a,b). This stems from the hypothesis that if two sentences in a language (e.g. English) have the same translation in another, (e.g. French) they must be paraphrases of each other. While these pipeline approaches bypass the problem of missing data, they propagate errors. Further, all neural paraphrasing models (e.g., Prakash et al., 2016; Gupta et al., 2018; Wang et al., 2019) predict discrete tokens through a final softmax layer. We hypothesize that softmax-based architectures restrict the diversity of outputs, biasing the models to copy words and phrases from the input, which has an effect opposite to the intended one in paraphrasing.

In this work, we introduce PARAVMF – a simple and effective method of training paraphrasing models by generating into embedding spaces (§2). Since parallel paraphrasing data is not available even in otherwise high-resource languages like French, we focus on an unsupervised approach. Using bilingual parallel corpora, we adapt multilingual machine translation (Johnson et al., 2017) to monolingual translation. We propose to train this model with translation and autoencoding objectives. The latter helps simplify the training setup by using only one language pair, whereas prior work required multiple language pairs and more data to stabilize training (Tiedemann and Scherrer, 2019; Buck et al., 2018; Guo et al., 2019; Thompson and Post, 2020). To encourage diversity, we propose to replace the final softmax layer in the decoder with a layer that learns to predict word vectors (Kumar and Tsvetkov, 2019). We show that predicting into word meaning representations increases diversity in paraphrasing by generating semantically similar words and phrases which are often neighbors in the embedding space.

We evaluate our proposed model on paraphras-

\*Currently at Google LLC

<sup>1</sup>The code is available at [https://github.com/monisha-jega/paraphrasing\\_embedding\\_outputs](https://github.com/monisha-jega/paraphrasing_embedding_outputs)

ing English and French sentences (§3). In several setups, standard automatic metrics and human judgment experiments show that our zero-shot paraphrasing model with embedding outputs generates more diverse and fluent paraphrases, compared to state-of-the-art methods (§4).

## 2 The PARAVMF Model

Let the language to paraphrase in be  $L_1$ . Our goal is to learn a mapping  $f(\mathbf{x}; \theta)$  parameterized by  $\theta$ .  $f$  takes a text  $\mathbf{x} = (x_1, x_2, \dots, x_m)$  containing  $m$  words as input, which can be a sentence or a segment in  $L_1$ . It then generates  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  of length  $n$  in the same language such that  $\mathbf{x}$  and  $\mathbf{y}$  are paraphrases. That is,  $\mathbf{y}$  represents the same meaning as  $\mathbf{x}$  using different phrasing. We assume that no direct supervision data is available, but there exists a bilingual parallel corpus between  $L_1$  and another language  $L_2$ . We are also given pre-trained embeddings (Bojanowski et al., 2017) for words in both  $L_1$  and  $L_2$ . The dimension of both the embedding spaces is  $d$ .

We use a standard transformer-based encoder-decoder model (Vaswani et al., 2017) as the underlying architecture for  $f$ . As visualized in the system diagram presented in the Appendix,  $f$  is jointly trained to perform three tasks with a shared encoder and decoder: (1) translation from  $L_1$  to  $L_2$ , (2) translation from  $L_2$  to  $L_1$  and (3) reconstructing the input text in  $L_1$  (autoencoding).<sup>2</sup>

Towards our primary goal of meaning preservation, the translation objectives help the encoder map the inputs in both the languages to a common semantic space, whereas the decoder learns to generate language-specific outputs. On the other hand, with the autoencoding objective, we expose the model to examples where the input and output are in the same language, biasing the model to adhere to the start token supplied to it and decode monolingually. Using this training algorithm, we find in our experiments (§4), that the resulting paraphrases albeit meaning-preserving still lack in diversity. We identify two reasons for this issue. First, the model overfits to the autoencoding objective and just learns to copy the input sentences. We address this issue by using only a small random sample of the total training sentences for training

with this objective.<sup>3</sup>

Second, we find that cross-entropy loss used to train the model results in peaky distributions at each decoding step where the target words get most of the probability mass. This distribution being another signal of overfitting also reduces diversity (Meister et al., 2020). We find in our preliminary experiments, that prior work to address this issue by augmenting diversity inducing objectives to the training loss (Vijayakumar et al., 2018) often comes at a cost of reducing meaning preservation. In this work, we propose using a different training loss which naturally promotes output diversity. We follow Kumar and Tsvetkov (2019), and instead of treating each word  $w$  in the vocabulary as a discrete unit, we represent it using a unit-normalized pre-trained vector  $\mathbf{e}$  learned using monolingual corpora (Bojanowski et al., 2017). At each decoding step, instead of predicting a probability distribution over the vocabulary using a softmax layer, we predict a  $d$ -dimensional continuous-valued vector  $\hat{\mathbf{e}}$ . We train our proposed model by minimizing von Mises-Fisher (vMF) loss—a probabilistic variant of cosine distance—between the predicted vector and the pre-trained vector. At each step of decoding, the output word is generated by finding the closest neighbor (using cosine similarity) of the predicted output vector  $\hat{\mathbf{e}}$  in the pre-trained embedding table. Since this loss does not directly optimize for a specific token but for a vector subspace which contains many similar meaning words, we observe that it has a higher tendency to generate diverse outputs than softmax-based models, both at the lexical and syntactic level as we show in our experiments.

Overall, the contribution of this work is twofold: (1) a translation and autoencoding based training objective to enable paraphrasing while preserving meaning without any parallel paraphrasing data, and (2) optimizing for vector subspaces instead of token probabilities to induce diversity of outputs.

## 3 Experiments

**Datasets** We evaluate paraphrasing in two languages: English and French. IWSLT’16 En↔Fr corpus (Cettolo et al., 2016) with  $\sim 220\text{K}$  sentence pairs is used for training with translation objective, and 4450 sentences, randomly sampled  $\sim 1\%$  of the training data in  $L_1$  (either En or Fr), for au-

<sup>2</sup>To bias the model against always decoding in the other language, unlike in Johnson et al. (2017); Tiedemann and Scherrer (2019), we provide a language-specific start token in the encoder input, in addition to the decoder input.

<sup>3</sup>We empirically determine this sample size to be  $\sim 1\%$  of the total number of training examples.

toencoding. We use the  $L_1$  side of the IWSLT’16 dev set for early stopping with the autoencoding objective. We use IWSLT’16 test set for automatic evaluation consisting of 2331 samples in En and Fr each. For human evaluation we subsample 200 sentences from this set. We tokenize and true-case all the data using Moses preprocessing scripts (Koehn et al., 2007). We conduct additional experiments with a larger En–Fr corpus constructed using a 2M sentence-pair subset of the combination of the WMT’10 Gigaword (Tiedemann, 2012) and the OpenSubtitles corpora (Lison and Tiedemann, 2016).

**Implementation** We modify the standard seq2seq transformer model in OpenNMT (Klein et al., 2017) to generate word embeddings (Kumar and Tsvetkov, 2019), and train it with the vMF loss with respect to target vectors. We initialize and fix the input embeddings of the encoder and decoder with off-the-shelf (sub-word based) fasttext embeddings (Bojanowski et al., 2017) for both En and Fr and align the embeddings to encourage cross-lingual sharing (Artetxe et al., 2018). With a vocabulary size of 50K for each language, the combined vocabulary size of the encoder and the decoder is 100K. Both encoder and decoder consist of 6 layers with 4 attention heads. The model is optimized using Adam (Kingma and Ba, 2015), with batch size 4K, and 0.3 dropout. The hidden dimension size is 1024, the dimension of the embedding layers is 512. We add a linear layer to transform 300-dimensional pre-trained embeddings to 512-dimensional input vectors to the model. After decoding, we postprocess the generated output to replace words from  $L_2$  by a look-up in the dictionary induced from the aligned embedding spaces.

**Baselines** Although unsupervised methods of paraphrasing with only monolingual data have been explored in recent works (Gupta et al., 2018; Yang et al., 2019; Roy and Grangier, 2019; Patro et al., 2018; Park et al., 2019) they have not been shown to outperform translation based baselines (West et al., 2020). Hence we compare our proposed approach with translation-based baselines only. First, we compare with bilingual pivoting baselines (Mallinson et al., 2017a,b) which pipeline two separate translation models,  $L_1 \rightarrow L_2$ , and  $L_2 \rightarrow L_1$ . We use two bilingual pivoting baselines, one based on continuous-output model (**BP-vMF**; the output vectors of the first model are first converted to discrete tokens before being fed to the next) and

another based on softmax-based model (**BP-CE**).

To evaluate the impact of embedding outputs, we also compare our proposed model **PARAVMF** to softmax-based baseline **PARACE**, leaving other model components unchanged. **PARACE** is a modified bilingual version of the multilingual method proposed in Guo et al. (2019), the current state-of-the-art in zero-shot paraphrasing.

**Evaluation setup** There are many ways to paraphrase a sentence, but no manually crafted multi-reference paraphrase datasets exist, that could be used as test sets (and there are no datasets in languages other than English). We thus evaluate the generated paraphrases on semantic similarity and lexical diversity compared to the *input text*. Following prior work, we use the  $n$ -gram based metric **METEOR** (Banerjee and Lavie, 2005). Despite accounting for synonyms, it is not well-suited to evaluate paraphrases, since it typically assigns lower scores to novel phrasings, due to incomplete synonym dictionaries. We thus also include **BERTScore** (Zhang et al., 2020), computing cosine similarity between the contextual embeddings of two sentences. Naturally, just copying the inputs can also lead to high scores in these metrics. To evaluate lexical diversity, we follow Hu et al. (2019b) and include **IoU** – Intersection over Union (also called Jaccard Index) and Word Error Rate (**WER**). To measure structural diversity we use (constituency) Parse Tree Edit distance (**PTED**).<sup>4</sup> Note that model outputs that do not preserve meaning in paraphrasing (and generate totally different sentences) will also obtain high diversity scores, but these are not indicative of quality paraphrasing but will falsely contribute to high diversity scores if averaged across the entire test set. We thus measure the diversity only on subsets of the test set for which the strongest baseline (PARACE) and our model generate meaning-preserving paraphrases measured using BERTScore thresholds. We report the diversity scores for three such thresholds: 0.95, 0.9, 0.85, selected empirically such that the sample size is sufficiently large.

## 4 Results

**Automatic evaluation** We observe in table 1 that **PARAVMF** outperforms all baselines in meaning-

<sup>4</sup>Before computing the PTED, we prune the tree to a max height of 3, and discard all the terminal nodes. We employ Stanford CoreNLP (Manning et al., 2014) for parsing and APTEd algorithm for edit distance (Pawlik and Augsten, 2015).

Model	ENGLISH		FRENCH	
	BS $\uparrow$	MET $\uparrow$	BS $\uparrow$	MET $\uparrow$
BP-CE	75.0	75.0	69.4	67.5
BP-VMF	72.1	72.2	65.5	64.2
PARACE	83.5	87.4	82.3	81.6
PARAVMF	<b>88.6</b>	<b>91.6</b>	<b>87.2</b>	<b>86.4</b>

Table 1: Meaning-preservation in generated paraphrases. BS: BertScore, MET: METEOR

preservation. Both pivoting based baselines perform poorly on average. This is a consequence of error propagation exacerbated in BP-VMF<sup>5</sup>. As a result, a very small fraction of generated sentences show meaning preservation (as measured by achieving a BERTScore greater than 0.85). Hence, we only compare the diversity in the two best meaning-preserving models, **PARACE** and **PARAVMF**. As shown in table 2, across all thresholds the latter model achieves higher lexical and syntactic diversity in the outputs. Ablation results in the Appendix show that both the autoencoding objective and the final embedding layer contribute to the improved quality of paraphrases. An additional benefit of our proposed model is that by replacing the softmax layer with word embeddings, PARAVMF is trained 3x faster than the PARACE baseline.

We further conduct a **manual evaluation** which quantifies the rate at which annotators find paraphrases fluent, consistent with input meaning, and novel in phrasing. In an A/B testing setup, we compare our proposed approach with the strongest baseline PARACE.<sup>6</sup> 200 sentences sampled from the IWSLT English test were scored by two annotators independently, which yielded the inter-annotator agreement of 0.37 (fair agreement). Out of the sentences on which both annotators agree (142 out of 200), we find that PARAVMF model outperforms the PARACE model in 73% of votes. We show more details and some examples of PARAVMF and PARACE system outputs in the Appendix.

Finally, we also evaluate that our results hold on

<sup>5</sup>This is expected as VMF has been shown to slightly underperform CE for translation in prior work (Kumar and Tsvetkov, 2019). Our training procedure with an autoencoding objective alleviates this issue in PARAVMF.

<sup>6</sup>Each judge is presented with a set of questions, each consisting of an input sentence and paraphrases generated by the two models as options, and is asked to choose the sentence that is fluent, meaning-preserving and offers a novel phrasing of the input. They are asked to choose neither if both sentences are dis-fluent and/or not able to preserve content. The options are shuffled.

a **larger dataset in different domain**. We retrain PARAVMF and PARACE on 2M En-Fr corpus described in §3.<sup>7</sup> The results of automatic evaluation are presented in the Appendix. We conduct human evaluation on a sample of 200 sentences from this test set following the same A/B testing procedure as described above, with each sample rated by three annotators, resulting in a pairwise-average kappa agreement index of 0.21.<sup>8</sup> 42.9% PARAVMF outputs were selected as better paraphrases, compared to 24.5% outputs from PARACE, supporting our main results on the IWSLT dataset.

## 5 Related Work

Bilingual pivoting is a common technique used with bilingual data (Barzilay and McKeown, 2001; Ganitkevitch et al., 2013; Pavlick et al., 2015; Mallinson et al., 2017a). PARANMT (Wieting and Gimpel, 2018) is a large psuedo-parallel paraphrase corpus constructed through back-translation (Wieting et al., 2017). Iyyer et al. (2018) augment it with syntactic constraints for controlled paraphrasing; PARABANK (Hu et al., 2019a) improves upon PARANMT via lexical constraining of decoding; and PARABANK 2 (Hu et al., 2019b) improves the diversity of paraphrases in PARABANK through a clustering-based approach. Note that these works are focused on English. Here, we propose a language-independent approach relying only on abundant bilingual data. Our approach is most similar to Guo et al. (2019) who use bilingual and multilingual translation for zero-shot paraphrasing. They, however, observe that bilingual models are insufficient for paraphrasing and are often unable to produce the output in the correct language. We incorporate an autoencoding objective which simplifies and stabilizes training, and embedding-based outputs improving the diversity in paraphrasing.

## 6 Conclusion

We present PARAVMF, an end-to-end model for generating paraphrases, trained solely with bilingual data, without any paraphrase supervision. We propose to generate paraphrases into meaning

<sup>7</sup>We use 4K English sentences subsampled ( $\sim 0.1\%$  of the training data) from the same corpus for autoencoding. To further discourage copying, we use denoised autoencoding (Lample et al., 2018).

<sup>8</sup>We discarded around 53 samples with no clear majority among the annotator ratings and report the results on the remaining samples, further ignoring cases where the paraphrases from both the models were rated to be of similar quality.



BERTScore threshold	Model	# (out of 2K)	IoU↓	ENGLISH WER↑	PTED↑	# (out of 2K)	IoU↓	FRENCH WER↑	PTED↑
0.85	PARACE	710	94.3	4	<b>0.5</b>	710	94.3	3.9	<b>0.55</b>
	PARAVMF		<b>92.4</b>	<b>4.1</b>	0.42		<b>92.7</b>	<b>4.1</b>	0.42
0.9	PARACE	539	96.2	2.6	<b>0.34</b>	580	96.1	2.6	<b>0.34</b>
	PARAVMF		<b>94.5</b>	<b>2.9</b>	0.29		<b>94.5</b>	<b>2.9</b>	0.29
0.95	PARACE	300	98.8	0.8	0.15	380	98.7	0.8	0.15
	PARAVMF		<b>97.7</b>	<b>1.2</b>	<b>0.16</b>		<b>97.7</b>	<b>1.2</b>	<b>0.16</b>

Table 2: Diversity of meaning-preserving paraphrases compared to the test set. PARAVMF outperforms a strong baseline PARACE for both English and French, across all metrics for thresholds 0.85 and 0.9, and in IoU and WER for threshold of 0.95.

spaces as opposed to discrete tokens. This leads to significant improvements in quality and diversity of paraphrasing over strong baselines.

## Acknowledgments

This material is based upon work supported by the National Science Foundation (NSF) under Grants No. IIS2040926 and IIS2007960. The views and opinions of authors expressed herein do not necessarily state or reflect those of the NSF.

## References

- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. [A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 789–798, Melbourne, Australia.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Regina Barzilay and Kathleen R. McKeown. 2001. [Extracting paraphrases from a parallel corpus](#). In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 50–57.
- Jonathan Berant and Percy Liang. 2014. [Semantic parsing via paraphrasing](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1415–1425.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Christian Buck, Jannis Bulian, Massimiliano Ciaramita, Wojciech Gajewski, Andrea Gesmundo, Neil Houlsby, and Wei Wang. 2018. [Ask the right questions: Active question reformulation with reinforcement learning](#). In *International Conference on Learning Representations*.
- Mauro Cettolo, Niehues Jan, Stüker Sebastian, Luisa Bentivogli, Roldano Cattoni, and Marcello Federico. 2016. The iwslt 2016 evaluation campaign. In *International Workshop on Spoken Language Translation*.
- Pablo Duboue and Jennifer Chu-Carroll. 2006. [Answering the question you wish they had asked: The impact of paraphrasing for question answering](#). In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, pages 33–36. Association for Computational Linguistics.
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. [Data augmentation for low-resource neural machine translation](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 567–573.
- Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. 2014. [Open question answering over curated and extracted knowledge bases](#). In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 1156–1165. Association for Computing Machinery.
- Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2013. [PPDB: The paraphrase database](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 758–764. Association for Computational Linguistics.
- Yinpeng Guo, Yi Liao, Xin Jiang, Qing Zhang, Yibo Zhang, and Qun Liu. 2019. Zero-shot paraphrase generation with multilingual language models. *arXiv preprint arXiv:1911.03597*.
- Ankush Gupta, Arvind Agarwal, Prawaan Singh, and Piyush Rai. 2018. A deep generative framework for paraphrase generation. In *The Thirty-Second AAAI Conference on Artificial Intelligence*.

- Sanda Harabagiu and Andrew Hickl. 2006. [Methods for using textual entailment in open-domain question answering](#). In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 905–912. Association for Computational Linguistics.
- Yutai Hou, Yijia Liu, Wanxiang Che, and Ting Liu. 2018. [Sequence-to-sequence data augmentation for dialogue language understanding](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1234–1245. Association for Computational Linguistics.
- J Edward Hu, Rachel Rudinger, Matt Post, and Benjamin Van Durme. 2019a. PARABANK: Monolingual bitext generation and sentential paraphrasing via lexically-constrained neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6521–6528.
- J. Edward Hu, Abhinav Singh, Nils Holzenberger, Matt Post, and Benjamin Van Durme. 2019b. [Large-scale, diverse, paraphrastic bitexts via sampling and clustering](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 44–54. Association for Computational Linguistics.
- Kentaro Inui, Atsushi Fujita, Tetsuro Takahashi, Ryu Iida, and Tomoya Iwakura. 2003. [Text simplification for reading assistance: A project note](#). In *Proceedings of the Second International Workshop on Paraphrasing - Volume 16*, page 9–16. Association for Computational Linguistics.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. [Adversarial example generation with syntactically controlled paraphrase networks](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1875–1885.
- Lifeng Jin, David King, Amad Hussein, Michael White, and Douglas Danforth. 2018. [Using paraphrasing and memory-augmented models to combat data sparsity in question interpretation with a virtual patient dialogue system](#). In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 13–23. Association for Computational Linguistics.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. [Google’s multilingual neural machine translation system: Enabling zero-shot translation](#). volume 5, pages 339–351.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*.
- Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. [OpenNMT: Open-source toolkit for neural machine translation](#). In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72. Association for Computational Linguistics.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions*, page 177–180. Association for Computational Linguistics.
- Sachin Kumar and Yulia Tsvetkov. 2019. [Von Mises-Fisher loss for training sequence to sequence models with continuous outputs](#). In *International Conference on Learning Representations*.
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018. [Unsupervised machine translation using monolingual corpora only](#). In *International Conference on Learning Representations*.
- Pierre Lison and Jörg Tiedemann. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles.
- Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2017a. [Paraphrasing revisited with neural machine translation](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 881–893.
- Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2017b. [Paraphrasing revisited with neural machine translation](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 881–893. Association for Computational Linguistics.
- Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. [The Stanford CoreNLP natural language processing toolkit](#). In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60.
- Clara Meister, Elizabeth Salesky, and Ryan Cotterell. 2020. [Generalized entropy regularization or: There’s nothing special about label smoothing](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6870–6886, Online. Association for Computational Linguistics.
- Sunghyun Park, Seung-won Hwang, Fuxiang Chen, Jaegul Choo, Jung-Woo Ha, Sunghun Kim, and

- Jinyeong Yim. 2019. Paraphrase diversification using counterfactual debiasing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6883–6891.
- Badri Narayana Patro, Vinod Kumar Kurmi, Sandeep Kumar, and Vinay Namboodiri. 2018. [Learning semantic sentence embeddings using sequential pairwise discriminator](#). pages 2715–2729.
- Ellie Pavlick and Chris Callison-Burch. 2016. [Simple PPDB: A paraphrase database for simplification](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 143–148.
- Ellie Pavlick, Pushpendre Rastogi, Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. 2015. [PPDB 2.0: Better paraphrase rankings, fine-grained entailment relations, word embeddings, and style classification](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 425–430. Association for Computational Linguistics.
- Mateusz Pawlik and Nikolaus Augsten. 2015. [Efficient computation of the tree edit distance](#). *ACM Trans. Database Syst.*, 40(1).
- Sarah E Petersen and Mari Ostendorf. 2007. Text simplification for language learners: A corpus analysis. In *Workshop on Speech and Language Technology in Education*.
- Aaditya Prakash, Sadid A. Hasan, Kathy Lee, Vivek V. Datla, Ashequl Qadir, Joey Liu, and Oladimeji Farri. 2016. [Neural paraphrase generation with stacked residual LSTM networks](#). *CoRR*, abs/1610.03098.
- Lorenza Romano, Milen Kouylekov, Idan Szpektor, Ido Dagan, and Alberto Lavelli. 2006. [Investigating a generic paraphrase-based approach for relation extraction](#). In *11th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Aurko Roy and David Grangier. 2019. Unsupervised paraphrasing without translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Brian Thompson and Matt Post. 2020. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. *arXiv preprint arXiv:2008.04935*.
- Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC’12)*, Istanbul, Turkey. European Language Resources Association (ELRA).
- Jörg Tiedemann and Yves Scherrer. 2019. [Measuring semantic abstraction of multilingual NMT with paraphrase recognition and generation tasks](#). In *Proceedings of the 3rd Workshop on Evaluating Vector Space Representations for NLP*, pages 35–42.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, undekasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, page 6000–6010.
- Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R. Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2018. [Diverse beam search: Decoding diverse solutions from neural sequence models](#).
- Su Wang, Rahul Gupta, Nancy Chang, and Jason Baldridge. 2019. A task in a suit and a tie: Paraphrase generation with semantic augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7176–7183.
- Peter West, Ximing Lu, Ari Holtzman, Chandra Bhagavatula, Jena Hwang, and Yejin Choi. 2020. Reflective decoding: Unsupervised paraphrasing and abductive reasoning. *arXiv preprint arXiv:2010.08566*.
- John Wieting and Kevin Gimpel. 2018. [PARAMT-50M: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations](#). pages 451–462.
- John Wieting, Jonathan Mallinson, and Kevin Gimpel. 2017. [Learning paraphrastic sentence embeddings from back-translated bitext](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 274–285. Association for Computational Linguistics.
- Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. [Optimizing statistical machine translation for text simplification](#). *Transactions of the Association for Computational Linguistics*, 4:401–415.
- Zhao Yan, Nan Duan, Junwei Bao, Peng Chen, Ming Zhou, Zhoujun Li, and Jianshe Zhou. 2016. [Doc-Chat: An information retrieval approach for chatbot engines using unstructured documents](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 516–525.
- Qian Yang, Zhouyuan Huo, Dinghan Shen, Yong Cheng, Wenlin Wang, Guoyin Wang, and Lawrence Carin. 2019. [An end-to-end generative architecture for paraphrase generation](#). 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 3132–3142. Association for Computational Linguistics.

Model	Semantic Similarity	
	BERTScore	METEOR
PARANMT	61.6	62.1
BP (vMF)	44.6	57.4
BP (CE)	45.0	60.4
PARACE	65.9	81.7
<b>PARAvMF</b>	<b>68.9</b>	<b>83.9</b>

Table 3: Evaluation of paraphrase generation on the PARANMT test set.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

## A System diagram

The PARAvMF system is represented diagrammatically in Figure 1.

## B Example outputs

Sample outputs of the PARAvMF and PARACE models are shown in table 5.

## C Training on a Larger Translation Dataset

To measure the impact of the size of parallel translation data used for training, we conduct an experiment with a larger French-English corpus constructed using a 2M sentence-pair subset of the combination of the WMT’10 Gigaword (Tiedemann, 2012) and the OpenSubtitles corpora (Lison and Tiedemann, 2016). The semantic similarity scores and the diversity results are presented in table 4. The results of human evaluation are presented in the main paper.

## D Evaluation on PARANMT-50M Test Set

We evaluate the PARAvMF model (trained on English-French two-way translation data and English autoencoding data from the IWSLT’16 dataset) on test data sampled from PARANMT-50M (Wieting and Gimpel, 2018), to demonstrate its paraphrasing ability on out-of-domain input, in addition to enabling direct comparison with back-translated data, as shown in table 3. However, it is to be noted that the comparison is not a fair one, since PARAvMF is trained on just 220K data samples, whereas PARANMT is back-translated using

a translation model that was trained on a bilingual dataset with a size of around 70M.

## E Ablation

We proposed three changes in a multilingual MT setup to use bilingual data for paraphrasing, (1) predicting continuous outputs and training with vMF loss, (2) language-specific start tokens in the encoder, and (3) an autoencoding objective. In the results section of the main paper, by comparing our method to PARACE, we already established the importance of using vMF compared to cross-entropy. As shown in table 7, ablating either of the other remaining two components leads to considerable performance drop. This is because the ablated models generate outputs in  $L_2$  since they are never exposed to monolingual examples during training. Additional, in our preliminary experiments, we also observe that increasing the size of autoencoding data too much beyond  $\sim 1\%$  of the size of parallel translation data leads to a performance drop because the model just starts to learn to copy the input as-is rather than rephrasing.



Model	ENGLISH	
	BERTScore↑	METEOR↑
PARACE	62.2	73.6
PARAvMF	<b>71.6</b>	<b>79.6</b>

(a) Semantic similarity between the test set and generated paraphrases

BERTScore threshold	Model	# (out of 2K)	ENGLISH		
			IoU↓	WER↑	PTED↑
0.85	PARACE	559	85.5	11.9	<b>1.43</b>
	PARAvMF		<b>82.5</b>	<b>12.4</b>	1.42
0.9	PARACE	327	91.2	7.0	<b>0.81</b>
	PARAvMF		<b>87.9</b>	<b>8.3</b>	0.64
0.95	PARACE	196	95.9	3.3	0.28
	PARAvMF		<b>93.9</b>	<b>3.9</b>	<b>0.29</b>

(b) Diversity of meaning-preserving paraphrases compared to the test set

Table 4: Evaluation of paraphrase generation with PARAvMF trained on  $2M$  English-French sentence pairs. It outperforms a strong cross-entropy based baseline (PARACE) on semantic similarity and majority of diversity metrics.

Input	It 's expensive , it takes a long time , and it 's very complicated .
PARACE	It 's expensive takes a time , and it 's very complicated .
PARAvMF	It 's costly , It takes a long time , and it 's very difficult .
Input	These are things to talk about and think about now , with your family and your loved ones .
PARACE	These are things to talk about and think about now , with your family and your loved ones .
PARAvMF	These are things to speak of and think of now , with your family and the ones you love.
Input	So what opened my eyes ?
PARACE	So what opened my eyes ?
PARAvMF	So what is it that opened my eyes up ?
Input	And this work has been wonderful . It 's been great .
PARACE	And this work has been wonderful . It 's been great .
PARAvMF	This work has been wonderful and great .
Input	I wasn 't doing anything that was out of the ordinary at all .
PARACE	I wasn 't doing anything that was out of the regular regular at all .
PARAvMF	I was doing nothing that was not ordinary .
Input	It will make tons of people watch , because people want this experience .
PARACE	It will make tons of people watch , because people want this .
PARAvMF	Tonnes of people will look because they want this experience .

Table 5: Comparison of selected sample outputs for the IWSLT Test Set between PARAvMF model and the baselines. PARAvMF not only exhibits content preservation, but also demonstrates fluency as well as lexical and syntactic diversity.

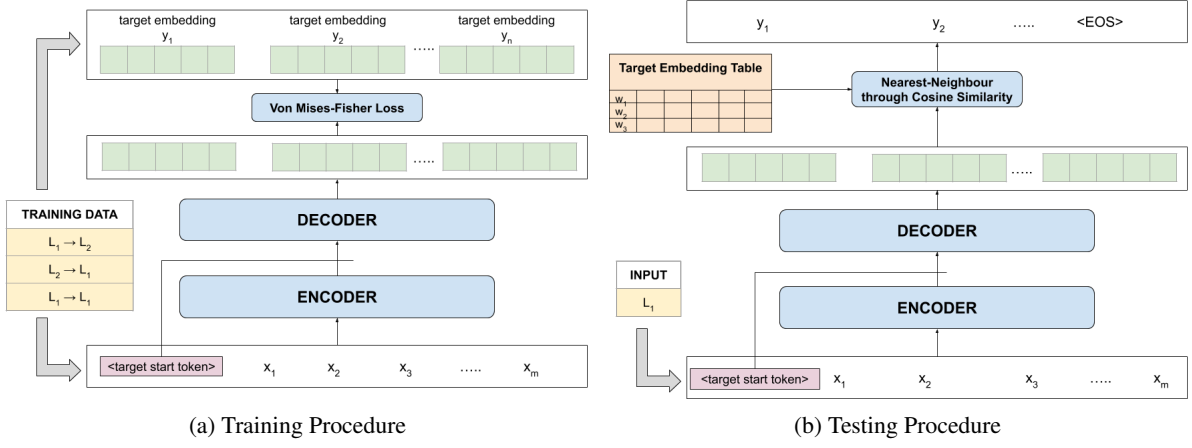


Figure 1: The PARAvMF Model: The decoder generates continuous-valued vectors at each step. It is trained by minimizing von Mises-Fisher loss between the output vectors and the pre-trained embeddings of the target words. Start tokens signalling the target language are supplied to both the encoder and the decoder. The training data consists of translation samples,  $L_1 \leftrightarrow L_2$  and autoencoding samples,  $L_1 \rightarrow L_1$ . During testing, the word in the target vocabulary whose embedding is closest to the generated output in terms of cosine similarity is output.

Model	Votes (%)
PARACE	39 (27.3%)
PARAvMF	<b>104 (72.7%)</b>

Table 6: PARAvMF outperforms the baseline in manual A/B testing (English).

Model	BLEU	BS	MET.
PARAvMF	64.0	88.6	91.6
- encoder start token	0.86	46.0	12.0
- autoencoding	0.85	46.0	12.1

Table 7: Performance of PARAvMF without the proposed enhancements - removing either leads to a drastic performance drop