# Simultaneous paraphrasing and translation by fine-tuning Transformer models

**Rakesh Chada**
Amazon.com, Inc.,
`rakchada@amazon.com`

## Abstract

This paper describes the third place submission to the shared task on simultaneous translation and paraphrasing for language education at the 4th workshop on Neural Generation and Translation (WNGT) for ACL 2020. The final system leverages pre-trained translation models and uses a Transformer architecture combined with an oversampling strategy to achieve a competitive performance. This system significantly outperforms the baseline on Hungarian (27% absolute improvement in Weighted Macro F1 score) and Portuguese (33% absolute improvement) languages.

## 1 Introduction

This paper describes the third place submission to the shared task Mayhew et al. (2020) on simultaneous translation and paraphrasing for language education at the 4th workshop on Neural Generation and Translation (WNGT) for ACL 2020. The shared task involves generating multiple translations for a given source text in English and a target language. The five target languages in the task are Hungarian (hu), Portuguese (pt), Japanese (ja), Korean (ko) and Vietnamese (vi). We competed in the Hungarian and Portuguese tracks. A goal of the shared task, hosted by Duolingo, is to enable development of automated grading processes and curation systems for language learners' responses. A high-coverage and precise multi-output translation and paraphrasing system would vastly help such automated efforts. For the task, participants were provided with hand-crafted and field-tested sets of several possible translations for each English sentence. Each of these translations were also ranked and weighted according to actual learner response frequency and these weights were provided as additional features. Along with these, translations from AWS were provided as a baseline and additional data. The challenges associated with the shared task are two-fold: i) Translating from English to target languages and ii) Producing multiple valid translations (paraphrases) while balancing precision with the coverage. We conduct several experiments to address these two challenges and develop a simple system that leverages pre-trained transformer Vaswani et al. (2017) models and a wide beam search strategy. Furthermore, we leverage the provided translation scores and experiment with multiple training distribution strategies to develop a simple oversampling strategy that produces improvements over the vanilla method of using one translation one time.

## 2 Related work

Paraphrasing and machine translation are well-studied research areas in general but there's not much research specifically in the context of multi-output translation systems, especially for low resource languages. Tan et al. (2019) train a Transformer-based Neural Machine Translation model for Hungarian-English and Portugese-English translation. However, their goal was to assess the benefits of multilingual modeling by clustering languages and is different from that of a multi-output translation system. For English-Portuguese, Aires et al. (2016) build a phrase-based machine translation system to translate biomedical texts. For multilingual parahrasing, Ganitkevitch and Callison-Burch (2014) release a database consisting of paraphrases for several languages, including Hungarian and Portuguese, at lexical, phrasal and syntactic level. Guo et al. (2019) build a zero-shot multilingual paraphrase generation model to show mixed results. However, their end goal was to generate paraphrases in the same language (English) as opposed to our shared task which requires generating paraphrases in a different language.

| Target Language | Train | | | | | | Dev | Test |
|---|---|---|---|---|---|---|---|---|
| | Prompts | Pairs | MSL | MTL | 99p SL | 99p TL | Prompts | Pairs |
| Hungarian (hu) | 4000 | 251442 | 21 | 21 | 11 | 14 | 500 | 500 |
| Portuguese (pt) | 4000 | 526466 | 33 | 21 | 25 | 15 | 500 | 500 |

Table 1: Dataset statistics. MSL=Maximum Source Length. MTL=Maximum Target Length. 99p SL=99th percentile Source Length. 99p TL=99th percentile Target Length.

Ippolito et al. (2019) study diverse decoding methods on conditional language models and show promising results on movie dialogue corpus and image captioning tasks.

## 3 Task

We describe dataset statistics and evaluation metrics in this section.

### 3.1 Data

There are two phases of the competition - Dev and Test. Table 1 shows data statistics for all phases. There were 4000 train prompts provided, in English, for both Hungarian and Portuguese languages. However, each of these prompts were accompanied with multiple translations leading to 251,442 English-Hungarian (en-hu) pairs and 526,466 English-Portuguese (en-pt) pairs. There were 500 prompts in both dev and test phases. After tokenization, for en-hu, most of the source sentences were shorter than 11 tokens and target sentences were shorter than 14 tokens. For en-pt, most of the source sentences were shorter than 25 tokens and target sentences were shorter than 15 tokens.

### 3.2 Evaluation Metrics

The main scoring metric for the competition is the weighted macro F1 score. This is a measure of how well the system returns all human-curated translations weighted by the likelihood that an English learner would respond with each translation. For each prompt p, weighted macro F1 is calculated as the harmonic mean of precision and weighted recall (note that the precision is unweighted). To calculated weighted recall for each example, we first calculate Weighted True Positives (WTP) and Weighted False Negatives (WFN) as:

$$WTP_p = \sum_{t \in TP_p} weight(t)$$

$$WFN_p = \sum_{t \in FN_p} weight(t)$$

Then, weighted recall (WR) is calculated as:

$$WR_p = \frac{WTP_p}{WTP_p + WFN_p}$$

The weighted Macro F1 (WF) over all prompts P is then calculated by averaging over all prompts in the corpus as:

$$WF = \sum_{p \in P} \frac{WF_p}{|P|}$$

## 4 System Design

We now describe the final submitted system design in detail. We have experimented with several other variants and describe these in a later section 5.

### 4.1 Data sampling

For the final system, we chose to use weighted sampling of the data where the weights correspond to the provided learner response frequency. Specifically, we multiply the frequency of the translation (a number between 0 and 1) with a heuristic value of 50 and duplicate the source-translation pair that many number of times. In effect, this would create repeated samples of certain pairs whose frequency is greater than 0.02 while eliminating pairs whose frequency is less than 0.02. With this sampling, we end up with 40,500 en-hu pairs and 42,000 en-pt pairs. We separate 15% of the provided prompts as a validation set. The performance on this validation set is used to pick the best model.

### 4.2 Preprocessing

For text pre-processing, we use sentencepiece tokenization Kudo and Richardson (2018) for en-hu and byte-pair encoding Sennrich et al. (2016) for en-pt data. We use pre-trained tokenization models provided in OPUS-MT.

### 4.3 Model Architecture

The final submitted model architecture, shown in Figure 1, uses the standard Transformer sequence-to-sequence model. This has 6 encoder and 6 decoder layers and an 8-headed attention mechanism

Figure 1: Architecture of the final system

in both encoder and decoder. We initialize the model with the pre-trained representations obtained from the OPUS-MT data. This model is then fine-tuned on the task data. We tie the encoder, decoder and output embedding weights and use a shared vocab size of 60,522. For position-wise feed-forward layers, the Swish activation function Ramachandran et al. (2018) is used. The whole model is fine-tuned, through an early stopping mechanism, on the dataset constructed as detailed in 4.1 .

For fine-tuning, we use the standard cross-entropy loss objective on the target sequence along with a label smoothing loss Szegedy et al. (2016).

For decoding, we use beam search with a beam size of 10 and select top 10 hypotheses for en-hu track. For en-pt track, we use a beam size of 28 and select top 28 hypotheses. We implement the model in Marian NMT Junczys-Dowmunt et al. (2018).

### 4.4 Postprocessing

The beam search outputs scores for each individual token. These scores represent the log likelihood of that token in the output sentence. As a post-processing step, we remove all translation predictions where the maximum of these token-level scores is less than -3.5. This value was determined by studying the impact of the maximum score thresholding on validation set performance.

### 4.5 Hyperparameters

We use the following hyperparameters. Batch size is set to 500. Dropout is set to 0.1. Label smoothing is set to 0.1. We use Adam optimizer with learning rate of 3e-4, $\beta_1$=0.9, $\beta_2$=0.98 and epsilon = 1e-9. We decay the learning rate by an inverse square root mechanism for 16000 steps. The gradient clip norm is set to 5. And patience for early stopping is set to 5.

## 5 Ablations

### 5.1 Ablations

We have performed several ablation studies on the en-hu task. The results of all these studies are listed in Table 3. We list the experiment methodologies

below.

**No fine-tuning**: Here, we applied the pre-trained translation model directly on the task without any fine-tuning. The decoding was done using beam search beam size of 12 and by selecting top 12 hypotheses (determined based on validation performance).

**No oversampling**: Here, we use all provided translation pairs without any filtering based on the learner response frequency. We fine-tune the pre-trained model on this dataset and decode using beam search with a beam size of 15 and selecting top 15 hypotheses.

**No post-processing**: This is the same as the final submitted model without the post-processing (maximum score thresholding).

## 6 Other Modeling Variants

We experimented with different modeling alternatives for the shared task. We describe them in this section. The results of these variations are listed in Table 4.

### 6.1 Multi-output sequence formulation

Here, we re-formulate the task as a multi-output prediction task by taking the top 5 translation pairs (based on the learner response frequency) and concatenating them into a single target sequence. The pre-trained model is then fine-tuned on this dataset.

**Nucleus sampling**: Here, we use the above multi-output sequence model and add Nucleus sampling Holtzman et al. (2019) while decoding with p value set to 0.95.

### 6.2 Back Translation

Here, we start with a pre-trained hu-en translation model. We then construct a hu-en dataset from the provided en-hu translation pairs. The pre-trained model is fine-tuned on this dataset. We apply this fine-tuned hu-en model on the provided reference AWS translations of the target hu sentences. With a beam size of 15 and top-5 hypotheses selection,

| Model | Validation | | | Dev | | | Test | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | WR | WF | P | WR | WF | P | WR | WF |
| Fairseq Baseline (en-hu) | - | - | - | 19.35 | 12.47 | 13.02 | 18.3 | 11.8 | 12.17 |
| AWS Baseline (en-hu) | - | - | - | **84.6** | 19.9 | 29.85 | **86.8** | 18.9 | 28.1 |
| Fine-tuned Transformer (en-hu) | 75.14 | 50.34 | 56.72 | 75.2 | **55.2** | **59.8** | 75.5 | **49.2** | **55.08** |
| Fairseq Baseline (en-pt) | - | - | - | 29.86 | 13.3 | 15.14 | 28.2 | 11.7 | 13.57 |
| AWS Baseline (en-pt) | - | - | - | **86.8** | 14.09 | 21.15 | **87.8** | 13.9 | 21.3 |
| Fine-tuned Transformer (en-pt) | 72.14 | 49.22 | 54.25 | 69.96 | **52.55** | **55.03** | 72.06 | **50.11** | **54.39** |

Table 2: Final submission results. **Bold** indicates best performance. P=Precision. WR=Weighted Recall. WF=Weighted Macro F1.

| Model | Validation | | | | | | |
|---|---|---|---|---|---|---|---|
| | P | R | WR | MiF | MaF | WMiF | WMaF |
| No fine-tuning | 52.41 | 6.32 | 41.18 | 11.28 | 19.21 | 46.12 | 41.31 |
| No oversampling | 58.40 | 13.26 | 46.70 | **21.62** | **32.34** | 51.90 | 45.98 |
| No post-processing | 74.04 | 9.28 | 49.71 | 16.49 | 28.81 | 59.49 | 54.93 |

Table 3: Performance of various en-hu ablations on validation dataset. **Bold** indicates best performance. R=Recall. MiF=Micro F1. MaF=Macro F1. WMiF=Weighted Micro F1. WMaF=Weighted Macro F1.

| Model | Validation | | | | | | |
|---|---|---|---|---|---|---|---|
| | P | R | WR | MiF | MaF | WMiF | WMaF |
| Multi-output sequence | **74.44** | 7.33 | 44.35 | 13.35 | 23.58 | 55.59 | 52.29 |
| Nucleus sampling | 72.98 | 7.70 | 45.13 | 13.93 | 24.27 | 55.77 | 52.67 |
| Back Translation | 70.98 | 7.42 | 44.45 | 13.43 | 23.63 | 54.67 | 52.08 |
| Model-based Prediction Filtering | 72.71 | **10.60** | **51.90** | 18.51 | 31.01 | **60.56** | **56.10** |

Table 4: Performance of modeling variants on en-hu validation dataset. **Bold** indicates best performance.

we generate 5 English paraphrases for each given English prompt. Now, the en-hu fine-tuned model from the "Multi-output sequence formulation" ablation is made to predict separately for each of the generated English paraphrases and all the outputs are combined into the final prediction.

### 6.3 Model-based Prediction Filtering

Here, we start with the final submission model and build a binary XGBoost classifier on top of it to filter predictions (accept vs reject). The features of the XGBoost model are the token-level scores, as described in Section 4.4, that are obtained from the final submission model. As different sequences have different lengths, we build a fixed size feature

vector by truncating or padding all sequences to a length of 11. This is the 99 percentile source length listed in Table 1. The binary labels for training are obtained by comparing output translation with the provided gold translations. We do a randomized search on "max_depth", "colsample_bytree", "colsample_bylevel" and "n_estimators" hyperparameters of the XGBoost model to find the best set of values. We then perform a 5-fold cross-validation to identify the best model. The F1 score of this model on the "accept" class is 0.81 and on the "reject" class is 0.48. The overall accuracy is about 72%.

## 7 Results & Discussion

Table 2 shows results of the final submission, for en-hu and en-pt tracks, along with a comparison to the baseline. As per the main evaluation metric (Weighted Macro F1 score), our model outperforms the strong AWS baseline by a significant margin on both en-hu and en-pt tracks. For en-hu, the improvement is about 30 absolute points on the dev dataset and 27 points on the test dataset. For en-pt, the improvement is about 34 absolute points on the dev dataset and 33 absolute points on the test dataset. This model ranked 1st on the dev leaderboard and 2nd on the test leaderboard for en-hu track. It ranked 2nd on the dev leaderboard and 3rd on the test leaderboard for en-pt track.

Table 3 shows the results for several ablations for en-hu model listed in section 5. And Table 4 shows results for several modeling variants listed in section 6. There are several interesting observations to be made from these ablations and variants. First, there's a clear improvement of about 15.4 points in Weighted Macro F1 from fine-tuning the pre-trained model on the provided dataset. The simple post-processing strategy of score thresholding yielded a gain of about 1.79 absolute points. Similarly, there's also a big improvement of about 10.7 absolute points from the oversampling strategy we used (as opposed to no oversampling). However, this gap seemed to have been closed by a big margin (about 7 absolute points) through the multi-output sequence formulation and slightly more by adding Nucleus sampling on top of it. A separate approach that uses back translation seemed to also have yielded similar gains upon the "No oversampling" approach. The model-based prediction filtering yielded an improvement of about 4 abso-

lute points. Interestingly, all of these variants still ended up inferior (by varying levels) to the simple oversampling + fine-tuning + post-processing strategy that was used for the final submission.

## 8 Summary

We describe the system for our submission to the shared task on simultaneous translation and paraphrasing for language education at the 4th workshop on Neural Generation and Translation (WNGT) for ACL 2020. The final submitted system leverages pre-trained translation models, with Transformer architecture, and an oversampling strategy to achieve competitive performance. For future, it'd be interesting to see if initializing the model with latest state-of-the-art sequence-to-sequence pre-trained models such as BART Lewis et al. (2019) and T5 Raffel et al. (2019) and fine-tuning could help boost performance. It would also be a promising direction to explore the benefit of using cross-lingual models such as XLM-Roberta Conneau et al. (2019). One way to use them would be to initialize the encoder part of the architecture with pre-trained representations. Given the shared representations, it might be interesting to see if concatenating several language pairs' train datasets and training a joint model produces additional benefits.

## Acknowledgments

## References

José Aires, Gabriel Lopes, and Luís Gomes. 2016. English-Portuguese biomedical translation task using a genuine phrase-based statistical machine translation approach. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 456–462, Berlin, Germany. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale.

Juri Ganitkevitch and Chris Callison-Burch. 2014. The multilingual paraphrase database. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, pages 4276–

4283, Reykjavik, Iceland. European Languages Resources Association (ELRA).

Yinpeng Guo, Yi Liao, Xin Jiang, Qing Zhang, Yibo Zhang, and Qun Liu. 2019. Zero-shot paraphrase generation with multilingual language models.

Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *ArXiv*, abs/1904.09751.

Daphne Ippolito, Reno Kriz, Joao Sedoc, Maria Kustikova, and Chris Callison-Burch. 2019. Comparison of diverse decoding methods from conditional language models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3752–3762, Florence, Italy. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension.

Stephen Mayhew, Klinton Bicknell, Chris Brust, Bill McDowell, Will Monroe, and Burr Settles. 2020. Simultaneous translation and paraphrase for language education. In *Proceedings of the ACL Workshop on Neural Generation and Translation (WNGT)*. ACL.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer.

Prajit Ramachandran, Barret Zoph, and Quoc V. Le. 2018. Searching for activation functions.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826.

Xu Tan, Jiale Chen, Di He, Yingce Xia, Tao Qin, and Tie-Yan Liu. 2019. Multilingual neural machine translation with language clustering.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.