# HW-TSC Translation Systems for the WMT22 Chat Translation Task

**Jinlong Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, Xiaoyu Chen,
Zhengzhe Yu, Zhiqiang Rao, Shaojun Li, Zhanglin Wu, Yuhao Xie,
Yuanchang Luo, Ting Zhu, Yanqing Zhao, Lizhi Lei, Hao Yang, Ying Qin**

Huawei Translation Service Center, Beijing, China

{yangjinlong7,lizongyao,weidaimeng,shanghengchao,chenxiaoyu35,
yuzhengzhe,raozhiqiang,lishaojun18,wuzhanglin2,xieyuhao2,luoyuanchang,
zhuting20,zhaoyanqing,leilizhi,yanghao30,qinying}@huawei.com

## Abstract

This paper describes the submissions of Huawei Translation Services Center(HW-TSC) to WMT22 chat translation shared task on English↔German (en-de) bidirection with results of zero-shot and few-shot tracks. We use the deep transformer architecture with a larger parameter size. Our submissions to the WMT21 News Translation task are used as the baselines. We adopt strategies such as back translation, forward translation, domain transfer, data selection, and noisy forward translation in task, and achieve competitive results on the development set. We also test the effectiveness of document translation on chat tasks. Due to the lack of chat data, the results on the development set show that it is not as effective as sentence-level translation models.

## 1 Introduction

Neural Machine Translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2015; Gehring et al., 2017) has achieved good translation results in most scenarios, but few researches have been done in the field of chat translation, mainly because of insufficient chat data.

WMT20 holds the chat translation shared task (Farajian et al., 2020) for the first time. The data set mainly includes pre-sales conversations between customers and agents (meal booking, air ticket reservation, etc.). This year, the data set focuses on post-sales conversations between customers and agents. Although the translation content is all about chat, the domains are slightly different. The results show that the data from previous years can not effectively improve the quality of the model for this year's task.

We participate in the en-de bidirectional translation task. The en-de bidirectional models we submitted to the WMT21 news task (Wei et al., 2021) are used as the baseline models and the architecture is deep transformer (Vaswani et al., 2017;

Dou et al., 2018). Commonly-used optimization strategies are used, such as domain transfer, data selection, back translation, self-training, noisy self-training, finetuning and model averaging.

Considering that the chat task is a context-aware translation task, we conduct a series of document-level (Wang et al., 2017) experiments using WMT document data, but it does not work well on development sets. The analysis shows that the document data deviates greatly from the chat domain, and the data therefore cannot effectively improve chat translation quality. According to the results, the best models are obtained by selecting in-domain data from out domain data by the development sets.

This paper is structured as follows: Section 2 describes our data volume and data pre-processing method. The model structure and method we used are presented in Section 3. Section 4 details our experiment setting. We present the results in Section 5, and finally we conclude our work in Section 6.

## 2 Data

### 2.1 Data Size

We use WMT21 news en-de bidirection models as our baselines (Wei et al., 2021). Bilingual data comes for WMT20 chat task Farajian et al. (2020), and monolingual data is from Byrne et al. (2019)

We select data of three related domains, including conversation, subtitle, and shopping, from our in-house English corpus for domain transfer. In addition, the document-level data from WMT22 general task [1] is used to train the document-level translation model. In addition, 40M in-house general bilingual data is used.

For details about the data size, see Table 1 and Table 2.

---

[1] Data is available from https://www.statmt.org/wmt22/

| | general | chat 20 | chat 22 | doc |
|---|---|---|---|---|
| en-de | 40M | 17847 | 2109 | 400K |

Table 1: Sentences size of bilingual data used for training

| | Domain-related | chat 20 | chat 22 | doc |
|---|---|---|---|---|
| en | 5M | 1M | 6389 | 20M |
| de | - | - | 7011 | 20M |

Table 2: Sentences size of monolingual data used for training

## 2.2 Data pre-processing

Considering that the data sizes of WMT20 and WMT22 chat tasks are limited, we do not cleanse the chat data. We use the following data cleansing methods for other data:

- Remove duplicate sentences (Khayrallah and Koehn, 2018; Ott et al., 2018)

- Filter out sentences with more than 150 words

- Filter out sentences with length ratios greater than 1.5

- Apply langid (Joulin et al., 2016, 2017) to filter out sentences in other languages

- Use fast-align (Dyer et al., 2013) to filter out sentence pairs that are poorly aligned.

Besides, we adopt joint SentencePiece Model(SPM) (Kudo and Richardson, 2018; Kudo, 2018) for word segmentation with a vocabulary of 32K.

## 3 System Overview

### 3.1 Model

Transformer has been widely used for neural machine translation in recent years, which has achieved good performance even with the most primitive architecture. Therefore, the baseline models for WMT21 news en-de task use the Transformer-Big architecture. Deep transformer is an improvement of Transformer, which increases the number of encoder layers and uses pre-layer-normalization to further improve model performance. Therefore, in this task, we adopt the following model architecture:

- Deep 25-6 large Model: This model features 25-layer encoder, 6-layer decoder, 1024 dimensions of word vector, 4096 domensions

of FFN, 16-head self-attention, and pre-layer-normalization.

### 3.2 Document-level NMT

Document-level machine translation (Ma et al., 2021) conditions on surrounding sentences to produce coherent translations. There has been a lot of work on custom model architectures to integrate document context into translation models.

There are many document translation strategies, such as Doc2Sent, Window2Window, Doc2Doc (Junczys-Dowmunt, 2019), DocBT (Junczys-Dowmunt, 2019), DocRepair (Voita et al., 2019), NoisyChannelDoc (Yu et al., 2019) and G-Transformer (Bao et al., 2021). Among the methods mentioned above, Doc2Doc and DocBT are preferred by us because the data processing procedures are simple and the model requires no modification.

To train our document-level model, the bilingual document data is spliced into a long sequence based on paragraph information and the sentences are separated by numbered <SEPX> symbols. For document-level monolinguals, we first generate synthetic bilingual data by back translation and use the same strategy to construct doc2doc data. We then use the document data to fine-tune the sentence-level translation model to ensure the model capable of translating long sequences.

We use two methods for inference. The first one translates single sentences just like a standard translation model. The other method combines a sentence with its context to construct a long-sequence input. After decoding, the model splits the result into single sentences and sacreBLEU[2] (Post, 2018) is calculated on the single sentences.

### 3.3 Data Selection

Data selection (van der Wees et al., 2017) is a data augmentation method that we use to select in-domain data from out-of-domain data.

For monolingual data selection, we train a Fast-Text (Joulin et al., 2016) classifier using a small number of English monolinguals in subtitle, conversation, and shopping domains, and then select in-domain English monolinguals from the common corpus.

For bilingual data selection, as mentioned by Wang et al. (2019, 2018) , we use the in-domain data to fine-tune the out-domain model, and then

---

[2] https://github.com/mjpost/sacrebleu

| System | 20 en→de test | 20 de→en test | 22 en→de dev | 22 de→en dev |
|---|---|---|---|---|
| baseline | 45.1 | 46.7 | 50.3 | 58.7 |
| + Data Selection | 49.2(+4.1) | 48.1(+1.4) | 62.5(+12.2) | 65.5(+6.8) |
| + Noisy FT | 49.0(+3.9) | 49.9(+3.2) | **64.3**(+13.1) | 65.5(+6.8) |
| + Model Average | 49.8(+4.7) | 49.2(+2.5) | 63.2(+12.9) | **65.7**(+7.0) |

Table 3: sacreBLEU score on chat20 test set and chat22 dev set

use the model before and after the fine-tuning to calculate the decoding probability score of the out-domain bilingual data. The data with a higher score on the fine-tuned model is selected as the in-domain bilingual data. The specific scoring is carried out according to the formula 1.

$$score = \frac{\log P(y|x;\theta_{in}) - \log P(y|x;\theta_{out})}{|y|} \quad (1)$$

Where $\theta_{out}$ represents the model trained with out-domain data, and $\theta_{in}$ represents the model after fine-tuning with a small amount of in-domain bilingual data, and |y| represents the length of the target sentence.

### 3.4 Forward Translation

Forward Translation (FT) (Wu et al., 2019), also known as Self-Training (Imamura and Sumita, 2018) , usually refers to using a forward NMT model to translate source-side monolingual data to target-side text so as to generate synthetic bilinguals. The data is then used to train the forward translation model. Generally, beam search (Freitag and Al-Onaizan, 2017) is used for forward translation. In our experiment, the beam size is set to 4.

Noisy self-training (He et al., 2020) adds noise to the source-side of the pseudo parallel corpus generated by forward translation. Experiments show that this method is effective in low resource tasks. Noisy self-training is therefore used in the last step when a small amount of in-domain monolinguals is used.

### 3.5 Back Translation

Back-translation(BT) (Edunov et al., 2018) has been recognized as the most effective data augmentation strategy for enhancing NMT model performance. Contrary to forward translation, it translates target-side monolinguals into source-side to generate synthetic parallel corpus. Among the many back translation methods, sampling (Graça et al., 2019), noise (Edunov et al., 2018) and tagged

back translation (Caswell et al.) work better. In our experiment, sampling back-translation is chosen.

### 3.6 Fine-tuning

Fine-tuning (Dakwale and Monz, 2017) is a way to achieve domain transfer. In our translation task, we adopt a three-stage fine-tuning strategy. Firstly, we use synthetic corpus from similar domains to fine-tune the out-of-domain NMT model, and then use bilingual data selected from general domain according to the development set to improve the model performance. After that, we use the synthetic data generated from the in-domain monolingual data to fine-tune the in-domain model for more fine-grained domain transfer.

### 3.7 Model Averaging

Model averaging (Dormann et al., 2018) is a commonly used technique to improve translation quality. Generally, models (5 in our experiment) that perform best on the development set are selected for parameter averaging, result to significantly improvement.

## 4 Experiment Setting

During the training phase, we use Pytorch-based Fairseq[3] (Ott et al., 2019) open-source framework as our benchmark system. Each model is trained using 8 GPUs with a batch size of 2048. The update frequency is 4 and the learning rate is 5e-4. The label smoothing rate is set to 0.1, the warm-up steps to 4000, and the dropout to 0.3. Adam optimizer (Kingma and Ba, 2015) with $\beta1$=0.9 and $\beta2$=0.98 is also used. In the evaluation phase, we use Marian[4] (Junczys-Dowmunt et al., 2018) for decoding and then calculate the sacreBLEU scores on the WMT22 chat translation task dev sets to measure the performance of each model.

| System | 22 en→de sent | 22 en→de doc | 22 de→en sent | 22 de→en doc |
|---|---|---|---|---|
| baseline | 50.3 | - | 58.7 | - |
| + Data Selection | 62.5 | - | 65.5 | - |
| + Bilingual Doc | 36.7 | 37.1 | 45.2 | 44.5 |
| + Bilingual & Doc bt | 54.6 | 55.6 | 56.8 | 56.3 |

Table 4: The results of different strategies in the document-level model. Bilingual Doc means using WMT news bilingual doc data to finetune previous model. Bilingual & Doc bt means using WMT news bilingual doc data and pseudo corpus generated from WMT news monolingual doc data to finetune previous model.

| System | sacrebleu↑ | total↑ | er↑ | es↑ | sie↑ |
|---|---|---|---|---|---|
| Baseline | 25.8 | 0.37 | 0.12 | 0.82 | 0.16 |
| + Data Selection | **26.3** | 0.36 | 0.13 | 0.81 | 0.14 |
| +Bilingual Doc | 23.0 | **0.41** | 0.19 | 0.69 | 0.34 |
| +Bilingual Doc & DOC BT | 24.7 | 0.36 | 0.11 | 0.82 | 0.15 |

Table 5: The higher the accuracy of pronoun translation, the better the model combines contextual information.

## 5 Result and Analysis

Table 3 shows the main results on the development sets. Bilingual data selection gains significant improvement on dev sets. Although bilingual data is selected based on dev sets, the selected data consists of 13M sentences. Therefore, there is no overfitting risk. This strategy also improves model performance on chat20 test sets.

Since the dev set is already used to select data, we no longer use the dev set to fine-tune model. We then continue to train our model using noisy self-training strategy on monolingual in-domain data. The result shows that there is an increase in BLEU on the en→de track. After model averaging, performance on de→en track improves, but performance on en→de deteriorates.

Finally, we select the result of data selection after model averaging as the primary submission, noisy self-training after model averaging as the contrastive2. Note that, the two submissions before are few-shot. The result of the baseline model as the submission of the zero-shot track, is contrastive1.

### 5.1 Document-level NMT

According to the test results shown in Table 4, using document-level data to optimize models does not work well, mostly because this data is from the news domain. From our subsequent experiments, we also find that chat tasks have high requirements on data domain.

From rows 4 and 5 in Table 4, the model using

bilingual document data has worse results than the model using DocBT data. We assume that there are two reasons for this phenomenon. One is that the size of bilingual documents is limited, and the other is that the DocBT data generated using the data selection model is closer to the chat domain than the original bilinguals.

To verify the effectiveness of our document-level translation model, we evaluate our model on (Müller et al., 2018) test set, which is a pronoun translation accuracy task.

As shown in Table 5, the pronoun translation accuracy of bilingual document-level model was significantly better than that of other models. But the BLEU is the lowest due to the minimum amount of data. From subsequent domian transfer experiments, we can also find, chat tasks are extremely sensitive to the domain of the data, but we cannot find enough chat data to train the document-level translation model. Therefore, we cannot draw a conclusion that document-level translation is useless for chat translation tasks. Further researches can be carried out when sufficient chat data is available.

### 5.2 Domain Transfer

Since no chat training data is provided except for the development set, we continue to train the baseline model using development set and monolingual data from previous chat tasks. As shown in rows 3 and 4 in the Table 6, models training with chat20 development set performs well on the chat20 test set. However, little improvement is observed on chat22 dev set. As mentioned above, the data dis-

| System | 20 en→de test | 20 de→en test | 22 en→de dev | 22 de→en dev |
|---|---|---|---|---|
| baseline | 45.1 | 46.7 | 50.3 | 58.7 |
| + 20 dev fine-tune | 60.5 | 63.5 | 52.4(+1.9) | 53.5(-5.1) |
| + 20 mono en FT/BT | 59.6 | 64.1 | 45.7(-4.6) | 31.6(-27.1) |
| + Subtitle en FT/BT | 57.1 | 53.5 | 43.7(-6.6) | 28.5(-30.2) |
| + Conversation en FT/BT | 56.7 | 51.4 | 49.5(-0.5) | 43.8(-14.9) |
| + Shopping en FT/BT | 56.2 | 55.9 | 50.3(-) | 43.3(-15.4) |
| + Data Selection | 49.2 | 48.1 | 62.5(+12.2) | 65.5(+6.8) |

Table 6: The results of different strategies for the sentence-level model. FT/BT means that forward translation in en→de direction and back translation in de→en direction.

tribution for these two tasks is not consistent.

Monolingual data of similar domains, such as subtitle, conversation, and shopping, is then used for FT or BT enhancement. From rows 5, 6 and 7 in the Table 6, the results are worse than using chat20 data. Although the monolingual data is of higher quality, its domain and style are far away from chat data. So it brings no improvement.

### 5.3 Bilingual Data Selection

Through the above experiments, we find that this year's chat task has unique features and is very sensitive to domain differences. Using the idea proposed by Wang et al. (2019, 2018), we select 13M data from 40M general bilingual data to optimize our baseline model.

As can be seen from row 8 in the Table 6, this strategy is effective and improves the translation quality in both directions. Besides, we find that the data selected using the chat22 development set also improves model performance on the chat20 task, indicating that this strategy is a general method. We will test its applicability in the future.

## 6 Conclusion

This paper presents the submissions of HW-TSC to the WMT 2022 Chat Translation Shared Task. For both direction in customer-agent translation task, we perform experiments with a series of pre-processing and training strategies. The results show that bilingual data selection achieves the best results. In the future, we will continue to explore the applicability of bilingual data selection mentioned in this paper.

Besides, the performance of document-level translation model is limited given the amount of data. It has not achieved the expected results on this task, and we will continue to explore the impact of context for the chat task.

## References

Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.

Guangsheng Bao, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. 2021. G-transformer for document-level machine translation. In *ACL*.

Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Daniel Duckworth, Semih Yavuz, Ben Goodrich, Amit Dubey, Andy Cedilnik, and Kyu-Young Kim. 2019. Taskmaster-1: Toward a realistic and diverse dialog dataset. In *EMNLP*.

Isaac Caswell, Ciprian Chelba, and David Grangier. Tagged back-translation. *WMT 2019*, page 53.

P Dakwale and C Monz. 2017. Fine-tuning for neural machine translation with limited degradation across in-and out-of-domain data.

Carsten F. Dormann, Justin M. Calabrese, Gurutzeta Guillera-Arroita, Eleni Matechou, Volker Bahn, Kamil Bartoń, Colin M. Beale, Simone Ciuti, Jane Elith, Katharina Gerstner, Jérôme Guelat, Petr Keil, José J. Lahoz-Monfort, Laura J. Pollock, Björn Reineking, David R. Roberts, Boris Schröder, Wilfried Thuiller, David I. Warton, Brendan A. Wintle, Simon N. Wood, Rafael O. Wüest, and Florian Hartig. 2018. Model averaging in ecology: a review of bayesian, information-theoretic, and tactical approaches for predictive inference. *Ecological Monographs*, 88(4):485–504.

Zi-Yi Dou, Zhaopeng Tu, Xing Wang, Shuming Shi, and Tong Zhang. 2018. Exploiting deep representations for neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4253–4262.

Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 489–500.

M. Amin Farajian, António V. Lopes, André F. T. Martins, Sameen Maruf, and Gholamreza Haffari. 2020. Findings of the WMT 2020 shared task on chat translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 65–75, Online. Association for Computational Linguistics.

Markus Freitag and Yaser Al-Onaizan. 2017. Beam search strategies for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 56–60.

Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1243–1252.

Miguel Graça, Yunsu Kim, Julian Schamper, Shahram Khadivi, and Hermann Ney. 2019. Generalizing back-translation in neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 45–52.

Junxian He, Jiatao Gu, Jiajun Shen, and Marc'Aurelio Ranzato. 2020. Revisiting self-training for neural sequence generation. In *Proceedings of ICLR*.

Kenji Imamura and Eiichiro Sumita. 2018. Nict self-training approach to neural machine translation at nmt-2018. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 110–115.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Herve Jegou, and Tomas Mikolov. 2016. Fasttext. zip: Compressing text classification models.

Armand Joulin, Édouard Grave, Piotr Bojanowski, and Tomáš Mikolov. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431.

Marcin Junczys-Dowmunt. 2019. Microsoft translator at wmt 2019: Towards large-scale document-level neural machine translation. In *WMT*.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, et al. 2018. Marian: Fast neural machine translation in c++. In *ACL (4)*.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83.

Diederik P Kingma and Jimmy Lei Ba. 2015. Adam: A method for stochastic optimization.

Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *ACL (1)*.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *EMNLP (Demonstration)*.

Zhiyi Ma, Sergey Edunov, and Michael Auli. 2021. A comparison of approaches to document-level machine translation. *ArXiv*, abs/2101.11040.

Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A Large-Scale Test Set for the Evaluation of Context-Aware Pronoun Translation in Neural Machine Translation. In *WMT 2018*, Brussels, Belgium. Association for Computational Linguistics.

Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In *ICML*.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Matt Post. 2018. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, pages 3104–3112.

Marlies van der Wees, Arianna Bisazza, and Christof Monz. 2017. Dynamic data selection for neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1400–1410.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Elena Voita, Rico Sennrich, and Ivan Titov. 2019. Context-aware monolingual repair for neural machine translation. In *EMNLP*.

Longyue Wang, Zhaopeng Tu, Andy Way, and Qun Liu. 2017. Exploiting cross-sentence context for neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2826–2831, Copenhagen, Denmark. Association for Computational Linguistics.

Wei Wang, Isaac Caswell, and Ciprian Chelba. 2019. Dynamically composing domain-data selection with clean-data selection by "co-curricular learning" for neural machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1282–1292, Florence, Italy. Association for Computational Linguistics.

Wei Wang, Taro Watanabe, Macduff Hughes, Tetsuji Nakagawa, and Ciprian Chelba. 2018. Denoising neural machine translation training with trusted data and online data selection. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 133–143, Brussels, Belgium. Association for Computational Linguistics.

Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, Hao Yang, and Ying Qin. 2021. HW-TSC's participation in the WMT 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 225–231, Online. Association for Computational Linguistics.

Lijun Wu, Yiren Wang, Yingce Xia, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2019. Exploiting monolingual data at scale for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4207–4216.

Lei Yu, Laurent Sartran, Wojciech Stokowiec, Wang Ling, Lingpeng Kong, Phil Blunsom, and Chris Dyer. 2019. Putting machine translation in context with the noisy channel model. *ArXiv*, abs/1910.00553.