

When do Contrastive Word Alignments Improve Many-to-many Neural Machine Translation?

Anonymous ACL submission

Abstract

Word alignment has proven to benefit many-to-many neural machine translation (NMT). However, high-quality ground-truth bilingual dictionaries were used for pre-editing in previous methods, which are unavailable for most language pairs. Meanwhile, the contrastive objective can implicitly utilize automatically learned word alignment, which has not been explored in many-to-many NMT. This work proposes a word-level contrastive objective to leverage word alignments for many-to-many NMT. Empirical results show that this leads to 0.8 BLEU gains for several language pairs. Analyses reveal that in many-to-many NMT, the encoder's retrieval performance highly correlates with the translation quality, which explains when the proposed method impacts translation. This motivates future exploration for many-to-many NMT focusing on improving the encoder retrieval performance.

1 Introduction

Many-to-many neural machine translation (NMT) (Firat et al., 2016; Johnson et al., 2017; Aharoni et al., 2019; Sen et al., 2019; Arivazhagan et al., 2019; Lin et al., 2020; Pan et al., 2021b) jointly trains a translation system for multiple language pairs and obtain significant gains consistently across many translation directions. Previous work (Lin et al., 2020) shows that word alignment information helps improve pre-training for many-to-many NMT. However, cleaned high-quality ground-truth bilingual dictionaries are used to pre-edit the source sentences, which are unavailable for most language pairs.

Recently, contrastive objectives (Clark et al., 2020; Gunel et al., 2021; Giorgi et al., 2021; Wei et al., 2021) have been shown to be superior at leveraging alignment knowledge in various NLP tasks by contrasting the representations of positive and negative samples in a discriminative manner. This objective, which implicitly utilizes word alignment

learned by any toolkit refraining the constraints of using manually constructed dictionaries, has not been explored in the context of leveraging word alignment for many-to-many NMT.

An existing contrastive method (Pan et al., 2021b) (mRASP2) for many-to-many NMT relies on sentence-level alignments. Given that the incorporation of word alignments has led to improvements in previous work, we believe that fine-grained contrastive objectives focusing on word alignments should help improve translation. Therefore, this paper proposes word-level contrastive learning for many-to-many NMT using the word alignment extracted by automatic aligners. We conduct experiments on three many-to-many NMT systems covering general and spoken language domains. Results show that our proposed method achieves significant BLEU gains in the general domain compared to previous word alignment based methods and the sentence-level contrastive method.

We then analyze how the word-level contrastive objective affects NMT training. Inspired by previous work (Artetxe and Schwenk, 2019) training sentence retrieval model using many-to-many NMT, we speculate that our contrastive objectives affect the sentence retrieval performance and subsequently impact the translation quality. Further investigation reveals that in many-to-many NMT, the sentence retrieval precision of the multilingual encoder for a language pair strongly correlates with its translation quality (BLEU), which provides insight about when contrastive alignment improves translation. This revelation emphasizes the importance of improving the retrieval performance of the encoder for many-to-many NMT.

2 Word-level Contrastive Learning for Many-to-many NMT

Inspired by the contrastive learning framework (Chen et al., 2020) and the sentence-level contrastive learning objective of mRASP2, we pro-

pose a word-level contrastive learning objective to explicitly guide the training of the multilingual encoder to obtain well-aligned cross-lingual representations. Specifically, we use word alignments, obtained using automatic word aligners, to supervise the training of the multilingual encoder by a contrastive objective alongside the NMT objective.

Alignment Extraction Two main approaches for automatically extracting aligned words from a sentence pair are: using a bilingual dictionary and using unsupervised word aligners. The former extracts fewer but precise alignments, whereas the latter extracts more but noisy alignments. We extract word-level alignments by both methods and explore how they impact NMT training. For the former approach, we use word2word (Choe et al., 2020) to construct bilingual lexicons and then extract word pairs from parallel sentences. The extracted word pairs are combined to form a phrase if words are consecutive in the source and target sentence. For the latter approach, we use FastAlign (Dyer et al., 2013) and use only 1-to-1 mappings for training.

Word-level Contrastive Learning With the extracted alignments, we propose a word-level contrastive learning objective for the multilingual encoder by the motivation that the aligned words within a sentence pair should have a similar contextual representation. We expect the supervision of the contrastive objective on the corresponding contextual word representation leads to a robust multilingual encoder. Assume that the tokenized source and target parallel sentences in i -th batch are $\mathcal{D}_i = \{src_{ij}, tgt_{ij}\}_{j=1}^B$, and the extracted alignments from all the sentence pairs in each batch are $\mathcal{A}_i = \{s_{ik}, t_{ik}\}_{k=1}^N$, where B and N denote the batch-size and the number of alignments, respectively. Note that s_{ik} and t_{ik} may contain several tokens after the word combination for word2word or subword tokenization for NMT. Then the word-level contrastive loss in a batch is:

$$\mathcal{L}_{align}^{(i)} = - \sum_{k=1}^N \left(\log \frac{\exp(\text{sim}(s_{ik}, t_{ik})/\mathcal{T})}{\sum_{m=1}^N \exp(\text{sim}(s_{ik}, t_{im})/\mathcal{T})} + \log \frac{\exp(\text{sim}(s_{ik}, t_{ik})/\mathcal{T})}{\sum_{m=1}^N \exp(\text{sim}(s_{im}, t_{ik})/\mathcal{T})} \right) \quad (1)$$

where \mathcal{T} denotes a similarity scaling temperature. The similarity between two words is measured by:

$$\text{sim}(\text{word}_x, \text{word}_y) = \cos(g(\bar{\mathbf{x}}), g(\bar{\mathbf{y}})) \quad (2)$$

where $g(\mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$ and $\bar{\mathbf{x}}$ denotes the average of contextual hidden states of the correspond-

La. pair	Source	Size	N (w2w)	N (FA)
en-et	WMT18	1.9M	5,762,977	38,454,477
en-it	IWSLT17	231k	603,032	3,000,011
en-ja	IWSLT17	223k	684,583	2,797,882
en-kk	WMT19	124k	124,511	279,429
en-my	ALT	18k	75,383	377,392
en-nl	IWSLT17	237k	564,697	2,836,873
en-ro	WMT16	612k	3,271,848	13,092,240
en-tr	WMT17	207k	770,873	2,885,102
en-vi	IWSLT15	133k	354,167	2,120,755

Table 1: **Data Source and number of the extracted word pairs.** La. pair, N (w2w) and N (FA) denote the language pair, the number of the word pairs extracted by word2word and FastAlign, respectively. Refer to Appendix B for details of the dataset splits.

ing subword positions on top of the multilingual encoder. Following (Chen et al., 2020), we use an MLP between contrastive loss and the contextual representation for NMT loss. ReLU activation is used for σ , \mathbf{W}_1 is $d \times d$ and \mathbf{W}_2 is $d \times d'$, where d is the encoder’s hidden dimension and $d' < d$.

Finally, to jointly train with the NMT loss, we use the following equation to combine our proposed word-level contrastive loss for a batch:

$$\mathcal{L}^{(i)} = \frac{1}{B} (\mathcal{L}_{NMT}^{(i)} + w \frac{N_T}{2N} \mathcal{L}_{align}^{(i)}) \quad (3)$$

where N_T is the number of the tokens within a batch, $\frac{N_T}{2N}$ is a multiplier that scales the contrastive loss to be consistent with NMT loss, and w is a weight to balance the joint training.

3 Experimental Settings

Datasets and Preprocessing We selected ten languages, including English (en), Estonian (et), Italian (it), Japanese (ja), Kazakh (kk), Burmese (my), Dutch (nl), Romanian (ro), Turkish (tr), Vietnamese (vi) from different language families to train the NMT systems. We used the parallel datasets from different domains for the selected nine language pairs, including IWSLT, WMT, and ALT. We followed mBART (Liu et al., 2020) for tokenization. Details are given in Appendix A. For each parallel dataset, we implemented two approaches as stated in Section 2 to extract word pairs for the contrastive training objective. Data source and the number of the extracted word pairs are shown in Table 1. To ensure high alignment quality, we used large-scale out-of-domain (see Appendix B) parallel corpora with FastAlign.

Methods	222_en-ja	626_I	626_II
MLSC	13.90	23.76	13.55
+align	13.90	23.67	13.39
+w2w (ours)	13.85	23.44	13.69
+FA (ours)	13.30	23.68	13.48
mBART FT	18.90	29.11	20.64
+align	18.55	28.87	20.42
+w2w (ours)	18.80	29.08	20.89
+FA (ours)	18.65	29.01	20.87

Table 2: **Overall average BLEU of all the systems.** 626_I and 626_II denote “626_en-it-ja-nl-tr-vi” and “626_en-tr-ro-et-my-kk,” respectively. Results better than MLSC or mBART FT are marked **bold**. Refer to Appendix D for the detailed scores of all the systems.

Many-to-many NMT systems We established three many-to-many NMT systems as follows:

222_en-ja: Bidirectional en-ja NMT model using en-ja parallel corpus.

626_en-it-ja-nl-tr-vi: 6-to-6 multilingual NMT model using spoken language domain corpora for en-it, en-ja, en-nl, en-tr and en-vi.

626_en-tr-ro-et-my-kk: 6-to-6 multilingual NMT model using general domain corpora for en-tr, en-ro, en-et, en-my and en-kk.

Baselines and Ours For each language group setting above, we conducted NMT experiments on both the multilingual training from scratch (MLSC) (Johnson et al., 2017; Aharoni et al., 2019) and the mBART multilingual fine-tuning (mBART FT) (Tang et al., 2020) as baselines. We applied our proposed word-level contrastive learning in both MLSC and mBART FT, and compared with another strong baseline, word alignment based joint NMT training (+align) (Garg et al., 2019). For applying our method, we investigated the performance of joint training with word pairs extracted by both word2word (+w2w) and FastAlign (+FA).

Implementation We used mBART-large for mBART FT and transformer-base (Vaswani et al., 2017) for MLSC. See Appendix C for details.

4 Results and Analyses

BLEU Results We report case-sensitive tokenized BLEU (Papineni et al., 2002) results in Table 2 and 3. In Table 2, we observe that with our proposed training objectives, BLEU scores are comparable in 222_en-ja and 626_en-it-ja-nl-tr-vi while they are slightly improved in 626_en-tr-ro-et-my-kk. However, “+align” performs comparable or even worse compared with the baseline. Referring

Methods	en-tr		en-ro		en-et	
	→	←	→	←	→	←
MLSC	9.3	12.6	25.0	26.2	10.8	15.1
+align	9.0	12.4	24.6	26.5	10.7	14.6
+w2w (ours)	9.4	12.6	24.8	26.8	10.8	15.1
+FA (ours)	9.1	12.2	24.8	26.7	10.7	14.8
mBART FT	17.7	22.2	33.8	37.1	14.5	24.3
+align	17.5	21.9	33.8	36.7	15.2	24.3
+w2w (ours)	17.6	22.2	34.2	37.5	15.0	25.0
+FA (ours)	17.5	22.2	34.3	37.5	14.9	25.1

Methods	en-kk		en-my	
	→	←	→	←
MLSC	0.5	5.3	15.1	15.6
+align	0.4	5.4	15.0	15.3
+w2w (ours)	0.5	5.8	15.2	15.9
+FA (ours)	0.3	5.6	15.0	15.6
mBART FT	1.8	14.1	17.8	23.1
+align	1.8	14.0	16.9	22.1
+w2w (ours)	1.2	14.1	18.3	23.8
+FA (ours)	1.3	14.4	17.9	23.6

Table 3: **BLEU scores of 626_en-tr-ro-et-my-kk system.** Significantly better scores (Koehn, 2004) are in cyan and marginal improvements are in lightcyan.

to Table 3 for specific BLEUs on each language pair, we find that with our methods, translation performances are significantly improved for mBART FT while nontrivial improvements can merely be observed on en-ro and en-kk direction for MLSC.

Latent Encoder Alignment Property We now inspect which aspect of alignment-based methods impacts the translation performance. Previous work (Artetxe and Schwenk, 2019) show that the encoder of a strong multilingual NMT system is an ideal model for the bilingual sentence retrieval task. Inspired by this, we speculate that alignment-based objectives affect sentence retrieval performance, which further impacts the translation quality. We train MLSC and mBART FT and report the sentence retrieval precision and NMT loss during the training. Results are reported in Figure 2. We observe that the validation retrieval precision show the similar trend as the NMT loss. This indicates that during the normal many-to-many NMT training, encoder-side sentence-level retrieval precision is optimized along with the NMT loss.

Sentence Retrieval P@1 Correlates with BLEU According to the investigation of the encoder alignment property above, we verify the relationship between BLEU score and sentence retrieval precision on the validation set for each language pair. Results are shown in Figure 1. Cross-referencing

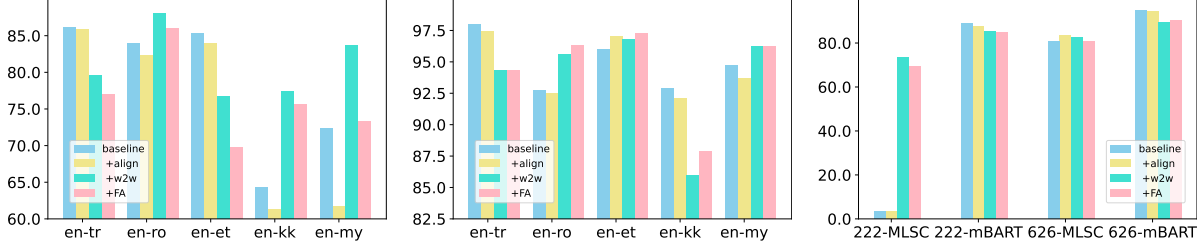


Figure 1: **Sentence retrieval P@1 on the validation set for each language pair.** *Left and middle* are the results on 626_en-tr-ro-et-my-kk. “626” in *right* subfigure denote 626_en-it-ja-nl-tr-vi. Refer to Appendix E for details.

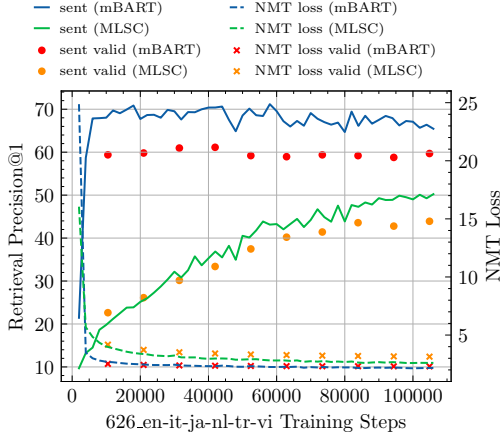


Figure 2: **NMT loss, sentence retrieval P@1 of the encoder in MLSC and mBART FT.** Average of the contextual embeddings on top of the encoder is used as the sentence embedding. We report the average score of forward and backward in-batch retrieval precisions.

the BLEU score in Table 3, we found that BLEU scores are improved when the encoder achieves gains on the sentence retrieval precision.¹ For example, we see increases of the retrieval P@1 on en-ro, en-et, and en-my on mBART FT (the middle of Figure 1) while BLEU scores are significantly improved on these three language pairs (Table 3). We further calculate the Pearson correlation coefficient between the BLEU changes and sentence retrieval P@1 changes for mBART+align, mBART+w2w, and mBART+FA in the 626_en-tr-ro-et-my-kk setting. Results are 0.79, 0.93, 0.90, respectively, demonstrating a strong correlation between translation quality and sentence retrieval precision.

Word-level Contrastive Objective and Sentence Retrieval P@1 With the word-level contrastive objective, we observed significant BLEU score improvements on language pairs such as en-ro, en-et and en-my as presented in Table 3. However, due

to the noises of extracted word pairs (Pan et al., 2021a) from word alignment toolkits that leads to insufficient supervision for improving sentence retrieval P@1, some language pairs such as en-kk do not show BLEU improvements. We found that for en-kk, numbers of extracted word pairs per sentence by word2word and FastAlign are 1.0 and 2.2, respectively. In contrast, the numbers are 4.2 and 20.7 for improved language pairs, calculated from Table 1. We expect this finding to provide new perspectives for improving many-to-many NMT.

Sentence-level Contrastive Objective We conducted the experiments for sentence-level contrastive objective (mRASP2) (Pan et al., 2021b) on 626_en-tr-ro-et-my-kk mBART FT. The average BLEU score of our +w2w is 20.89, which significantly outperforms mRASP’s 20.47 (last line in Table 8). Our word-level method outperforms the sentence-level method, indicating the sentence-level objective’s limitation. Moreover, we checked the sentence retrieval P@1 for mRASP2 (last line in Table 10) and found that it correlates with BLEU changes, indicating that sentence-level contrastive objective is suboptimal for language pairs with decreased retrieval precision.²

5 Conclusion

We proposed a word-level contrastive learning objective for many-to-many NMT. Experimental results showed that our proposed method leads to significantly better translation for several language pairs, which is then explained by analyses showing the relationship between BLEU scores and sentence retrieval performance of the NMT encoder. Future work can focus on: (1) further improving the encoder’s retrieval ability in many-to-many NMT; (2) contrastive objective’s feasibility in a massively multilingual scenario.

¹222_en-ja MLSC setting can hardly learn a well-aligned encoder while our methods improve the encoder sentence-level alignment quality without sacrificing BLEU scores.

²Note that the sentence-level contrastive objective incorporates sentences in multiple languages for contrastive loss. It does not necessarily improve the pair-wise retrieval precision.

References

- Roei Aharoni, Melvin Johnson, and Orhan Firat. 2019. [Massively multilingual neural machine translation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3874–3884, Minneapolis, Minnesota. Association for Computational Linguistics.
- Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George F. Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. 2019. [Massively multilingual neural machine translation in the wild: Findings and challenges](#). *CoRR*, abs/1907.05019.
- Mikel Artetxe and Holger Schwenk. 2019. [Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond](#). *Transactions of the Association for Computational Linguistics*, 7:597–610.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. [A simple framework for contrastive learning of visual representations](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research*, pages 1597–1607. PMLR.
- Yo Joong Choe, Kyubyong Park, and Dongwoo Kim. 2020. [word2word: A collection of bilingual lexicons for 3,564 language pairs](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3036–3045, Marseille, France. European Language Resources Association.
- Christos Christodoulopoulos and Mark Steedman. 2015. [A massively parallel corpus: the bible in 100 languages](#). *Lang. Resour. Evaluation*, 49(2):375–395.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: pre-training text encoders as discriminators rather than generators](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Chenchen Ding, Hnin Thu Zar Aye, Win Pa Pa, Khin Thandar Nwet, Khin Mar Soe, Masao Utiyama, and Eiichiro Sumita. 2020. [Towards burmese \(myanmar\) morphological analysis: Syllable-based tokenization and part-of-speech tagging](#). *ACM Trans. Asian Low Resour. Lang. Inf. Process.*, 19(1):5:1–5:34.
- 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, Atlanta, Georgia. Association for Computational Linguistics.
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016. [Multi-way, multilingual neural machine translation with a shared attention mechanism](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 866–875, San Diego, California. Association for Computational Linguistics.
- Sarthak Garg, Stephan Peitz, Udhyakumar Nallasamy, and Matthias Paulik. 2019. [Jointly learning to align and translate with transformer models](#). 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 4453–4462, Hong Kong, China. Association for Computational Linguistics.
- John M. Giorgi, Osvald Nitski, Bo Wang, and Gary D. Bader. 2021. [Declutr: Deep contrastive learning for unsupervised textual representations](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 879–895. Association for Computational Linguistics.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. 2021. [Supervised contrastive learning for pre-trained language model fine-tuning](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- 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](#). *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Philipp Koehn. 2004. [Statistical significance tests for machine translation evaluation](#). In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.
- Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu, Jiangtao Feng, Hao Zhou, and Lei Li. 2020. [Pre-training multilingual neural machine translation by leveraging alignment information](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2649–2663, Online. Association for Computational Linguistics.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. [Multilingual denoising pre-training for neural machine translation](#). *Transactions of the Association for Computational Linguistics*, 8:726–742.

- Hajime Morita, Daisuke Kawahara, and Sadao Kurohashi. 2015. [Morphological analysis for unsegmented languages using recurrent neural network language model](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2292–2297, Lisbon, Portugal. Association for Computational Linguistics.
- Toshiaki Nakazawa, Manabu Yaguchi, Kiyotaka Uchimoto, Masao Utiyama, Eiichiro Sumita, Sadao Kurohashi, and Hitoshi Isahara. 2016. [ASPEC: Asian scientific paper excerpt corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 2204–2208, Portorož, Slovenia. European Language Resources Association (ELRA).
- Lin Pan, Chung-Wei Hang, Haode Qi, Abhishek Shah, Saloni Potdar, and Mo Yu. 2021a. [Multilingual BERT post-pretraining alignment](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 210–219, Online. Association for Computational Linguistics.
- Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021b. [Contrastive learning for many-to-many multilingual neural machine translation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 244–258, Online. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Sukanta Sen, Kamal Kumar Gupta, Asif Ekbal, and Pushpak Bhattacharyya. 2019. [Multilingual unsupervised NMT using shared encoder and language-specific decoders](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3083–3089, Florence, Italy. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Edinburgh neural machine translation systems for WMT 16](#). In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 371–376, Berlin, Germany. Association for Computational Linguistics.
- Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. [Multilingual translation with extensible multilingual pretraining and finetuning](#). *CoRR*, abs/2008.00401.
- Arseny Tolmachev, Daisuke Kawahara, and Sadao Kurohashi. 2018. [Juman++: A morphological analysis toolkit for scriptio continua](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 54–59, Brussels, Belgium. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Xiangpeng Wei, Rongxiang Weng, Yue Hu, Luxi Xing, Heng Yu, and Weihua Luo. 2021. [On learning universal representations across languages](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

La. pair	Train	Valid	Test	OD Size
en-et	WMT18	WMT18	WMT18	10.7M
en-it	IWSLT17	IWSLT15	IWSLT16	13.6M
en-ja	IWSLT17	IWSLT15	IWSLT16	10.7M
en-kk	WMT19	WMT19	WMT19	851k
en-my	ALT	ALT	ALT	446k
en-nl	IWSLT17	IWSLT15	IWSLT16	12.7M
en-ro	WMT16	WMT16	WMT16	11.0M
en-tr	WMT17	WMT16	WMT16	11.1M
en-vi	IWSLT15	IWSLT13	IWSLT14	11.9M

Table 4: **Dataset statistics for each language pair.** “La. pair” means language pair and “OD Size” denotes the number of the out-of-domain sentence pairs used for training FastAlign.

Methods	en-ja	ja-en
MLSC	15.9	11.9
+align	16.3	11.5
+w2w (ours)	16.0	11.7
+FA (ours)	15.6	11.0
mBART FT	19.8	18.0
+align	19.6	17.5
+w2w (ours)	19.4	18.2
+FA (ours)	19.5	17.8

Table 5: **BLEU scores of 222_en-ja system.** Significantly better scores are in cyan and marginal improvements are in lightcyan. The significance test is done with Koehn (2004).

A Tokenization Settings

For Japanese, we use Jumanpp (Morita et al., 2015; Tolmachev et al., 2018) for segmentation and we following the setting in mBART (Liu et al., 2020) for other languages: `myseg.py` (Ding et al., 2020) is used for Burmese, Moses tokenization and special normalization is used for Romanian following (Sennrich et al., 2016),³ and Moses tokenization for other languages.⁴

B Datasets and Alignment Extraction

The datasets used for NMT training, validation and test are shown in Table 4. For the word alignment extraction using FastAlign, we also use out-of-domain parallel corpora to train the FastAlign jointly, aiming to obtain word alignments with less noise. The out-of-domain corpora for all the

³<https://github.com/rsennrich/wmt16-scripts>

⁴<https://github.com/moses-smt/mosesdecoder/blob/master/scripts/tokenizer/tokenizer.perl>

Methods	en-ja
MLSC	3.3
+align	3.5
+w2w (ours)	73.5
+FA (ours)	69.6
mBART FT	88.9
+align	87.4
+w2w (ours)	85.2
+FA (ours)	84.8

Table 6: **Sentence retrieval P@1 on the validation set for 222_en-ja.** The average of the contextual embeddings on top of the encoder is used as the sentence embedding. We report the average score of forward and backward retrieval precisions.

language pairs contain Tatoeba, Europarl, GlobalVoices, NewsCommentary, OpenSubtitles, TED, WikiMatrix, QED, GNOME, bible-uedin, and ASPEC (Nakazawa et al., 2016). We collect them from the OPUS project (Christodoulopoulos and Steedman, 2015) and WAT.⁵ The number of the out-of-domain parallel sentences for each language pair is shown in Table 4.

C Implementation Details

Following Tang et al. (2020), we set the oversampling temperature of 1.5 for all the settings. For MLSC, we set the dropout of 0.3 to avoid overfitting on small-scale training data. We used the batch size of 1,024 tokens for all the settings. For our word-level contrastive learning, we set the weight of 0.1, the temperature of 0.2, d' of 128, and a smaller dropout of 0.2 because our proposed objective serves as a regularization part. We followed the hyperparameter setting of Garg et al. (2019) for word alignment-based joint NMT training. We used 8 NVIDIA A100 for mBART FT and 8 TITAN Xp for MLSC model training. The model is validated every 1000 steps for 222_en-ja and 2000 steps for both two 626 settings. We do the early stopping if no improvement of the validation loss is observed for 8 checkpoints. The model with the best validation loss was used for evaluation.

D BLEU Scores

We report all the BLEU results of 222_en-ja, 626_en-it-ja-nl-tr-vi, and 626_en-tr-ro-et-my-kk in

⁵<https://lotus.kuee.kyoto-u.ac.jp/WAT/WAT2021/index.html>

Methods	en-ja		en-vi		en-it		en-nl		en-tr		Avg.
	→	←	→	←	→	←	→	←	→	←	
MLSC	15.4	11.8	29.6	28.6	27.5	32.7	29.1	36.4	11.6	14.9	23.76
+align	15.1	11.4	29.4	28.3	27.7	33.0	28.9	36.0	11.8	15.1	23.67
+w2w (ours)	15.3	11.6	29.7	28.2	27.6	32.4	28.6	35.8	10.8	14.4	23.44
+FA (ours)	15.5	11.6	29.6	28.0	27.8	33.2	29.1	35.9	11.2	14.9	23.68
mBART FT	17.8	17.0	34.1	35.7	32.5	38.0	32.6	41.6	18.7	23.1	29.11
+align	17.6	16.7	33.7	35.6	32.0	37.7	32.5	41.3	18.7	22.9	28.87
+w2w (ours)	17.6	17.2	34.2	35.7	32.5	38.2	32.1	41.7	18.7	22.9	29.08
+FA (ours)	17.5	17.7	34.0	35.2	32.4	37.9	32.3	41.4	18.6	23.1	29.01

Table 7: **BLEU scores of 626_en-it-ja-nl-tr-vi system.** Significantly better scores are in cyan and marginal improvements are in lightcyan. The significance test is done with [Koehn \(2004\)](#).

Methods	en-tr		en-ro		en-et		en-kk		en-my		Avg.
	→	←	→	←	→	←	→	←	→	←	
MLSC	9.3	12.6	25.0	26.2	10.8	15.1	0.5	5.3	15.1	15.6	13.55
+align	9.0	12.4	24.6	26.5	10.7	14.6	0.4	5.4	15.0	15.3	13.39
+w2w (ours)	9.4	12.6	24.8	26.8	10.8	15.1	0.5	5.8	15.2	15.9	13.69
+FA (ours)	9.1	12.2	24.8	26.7	10.7	14.8	0.3	5.6	15.0	15.6	13.48
mBART FT	17.7	22.2	33.8	37.1	14.5	24.3	1.8	14.1	17.8	23.1	20.64
+align	17.5	21.9	33.8	36.7	15.2	24.3	1.8	14.0	16.9	22.1	20.42
+w2w (ours)	17.6	22.2	34.2	37.5	15.0	25.0	1.2	14.1	18.3	23.8	20.89
+FA (ours)	17.5	22.2	34.3	37.5	14.9	25.1	1.3	14.4	17.9	23.6	20.87
+Sent (mRASP2)	17.2	22.0	34.0	36.8	14.1	24.2	1.8	13.7	17.5	23.4	20.47

Table 8: **BLEU scores of 626_en-tr-ro-et-my-kk system.** Significantly better scores are in cyan and marginal improvements are in lightcyan. The significance test is done with [Koehn \(2004\)](#).

Methods	en-ja	en-vi	en-it	en-nl	en-tr	Avg.
MLSC	52.7	84.6	91.0	85.7	89.7	80.9
+align	53.5	82.8	91.2	86.4	88.9	80.6
+w2w (ours)	73.4	85.7	91.4	84.7	83.1	83.7
+FA (ours)	71.3	84.9	91.3	83.8	82.0	82.7
mBART FT	87.1	96.2	97.3	94.6	98.5	94.7
+align	85.1	95.8	97.3	94.2	98.5	94.2
+w2w (ours)	81.6	91.4	94.7	90.8	89.6	89.6
+FA (ours)	82.6	92.3	95.0	91.7	90.4	90.4

Table 9: **Sentence retrieval P@1 on the validation set for 626_en-it-ja-nl-tr-vi.** Average of the contextual embeddings on top of the encoder is used as the sentence embedding. We report the average score of forward and backward retrieval precisions.

Methods	en-tr	en-ro	en-et	en-kk	en-my	Avg.
MLSC	86.2	84.0	85.4	64.4	72.4	78.5
+align	85.9	82.4	84.0	61.3	61.8	75.1
+w2w (ours)	79.6	88.1	76.8	77.4	83.7	81.1
+FA (ours)	77.0	86.1	69.8	75.7	73.4	76.4
mBART FT	98.0	92.7	96.0	92.9	94.7	94.9
+align	97.4	92.5	97.0	92.1	93.7	94.5
+w2w (ours)	94.3	95.6	96.8	86.0	96.2	93.8
+FA (ours)	94.3	96.3	97.3	87.9	96.2	94.4
+Sent (mRASP2)	94.6	95.5	89.0	89.6	95.6	92.9

Table 10: **Sentence retrieval P@1 on the validation set for 626_en-tr-ro-et-my-kk.** Average of the contextual embeddings on top of the encoder is used as the sentence embedding. We report the average score of forward and backward retrieval precisions.

Table 5, 7 and 8, respectively.

E Sentence Retrieval Precision

We report the sentence retrieval precisions for all the systems in Table 6, 9 and 10.