Multilingual Denoising Pre-training for Neural Machine Translation

Yinhan Liu^{‡*}, Jiatao Gu^{†*}, Naman Goyal^{†*}, Xian Li[†], Sergey Edunov[†], Marjan Ghazvininejad[†], Mike Lewis[†], and Luke Zettlemoyer[‡]

[†]Facebook AI

[‡]Birch Technology
[†]{jgu,naman,xianl,edunov,ghazvini,mikelewis,lsz}@fb.com
[‡]yinhan@birch.ai

Abstract

This paper demonstrates that multilingual denoising pre-training produces significant performance gains across a wide variety of machine translation (MT) tasks. We present mBART-a sequence-to-sequence denoising auto-encoder pre-trained on large-scale monolingual corpora in many languages using the BART objective (Lewis et al., 2019). mBART is the first method for pre-training a complete sequence-to-sequence model by denoising full texts in multiple languages, whereas previous approaches have focused only on the encoder, decoder, or reconstructing parts of the text. Pre-training a complete model allows it to be directly fine-tuned for supervised (both sentence-level and document-level) and unsupervised machine translation, with no taskspecific modifications. We demonstrate that adding mBART initialization produces performance gains in all but the highest-resource settings, including up to 12 BLEU points for low resource MT and over 5 BLEU points for many document-level and unsupervised models. We also show that it enables transfer to language pairs with no bi-text or that were not in the pre-training corpus, and present extensive analysis of which factors contribute the most to effective pre-training.¹

1 Introduction

Despite its wide adoption for other NLP tasks (Devlin et al., 2019; Liu et al., 2019; Yang et al., 2019b; Lewis et al., 2019; Raffel et al., 2019),

self-supervised pre-training is not yet common practice in machine translation (MT). Existing approaches (Lample and Conneau, 2019; Edunov et al., 2019; Lewis et al., 2019; Raffel et al., 2019) have been proposed either to partially pre-train the model or to only focus on English corpora. In this paper, we show that significant performance gains are possible by pre-training a complete autoregressive model with an objective that noises and reconstructs full texts across many languages.

In this work, we present *mBART*—a multilingual sequence-to-sequence (Seq2Seq) denoising auto-encoder. mBART is trained by applying the BART (Lewis et al., 2019) to large-scale monolingual corpora across many languages. The input texts are noised by masking phrases and permuting sentences, and a single Transformer (Vaswani et al., 2017) model is learned to recover the texts. Different from other pre-training approaches for MT (Lample and Conneau, 2019; Song et al., 2019), mBART pre-trains a complete autoregressive Seq2Seq model. mBART is trained once for all languages, providing a set of parameters that can be fine-tuned for any of the language pairs in both supervised and unsupervised settings, without any task-specific or language-specific modifications or initialization schemes.

Extensive experiments demonstrate that this simple approach works remarkably well. We first focus on existing MT benchmarks. For supervised sentence-level MT, mBART initialization leads to significant gains (up to 12 BLEU points) across low/medium-resource pairs (<10M bi-text pairs), without sacrificing performance in high-resource settings. These results further improve with back-translation (BT), setting a new state-of-the-art on WMT16 English-Romanian and the FloRes test sets. For document-level MT, our document-level pre-training improves results by up to 5.5

© 2020 Association for Computational Linguistics. Distributed under a CC-BY 4.0 license.

^{*} Equal contribution. Most of the work was done when the first author worked at Facebook.

¹Code and pre-trained models are available at https:// github.com/pytorch/fairseq/tree/master /examples/mbart.

Transactions of the Association for Computational Linguistics, vol. 8, pp. 726–742, 2020. https://doi.org/10.1162/tacl_a_00343 Action Editor: Stefan Riezler. Submission batch: 3/2020; Revision batch: 7/2020; Published 11/2020.



Figure 1: Sizes of the CC25 Corpus. A list of 25 languages ranked with monolingual corpus size.

BLEU points. For the unsupervised case, we see consistent gains and produce the first nondegenerate results for less related language pairs (e.g., 9.5 BLEU gain on Nepali-English). Previous pre-training schemes have only considered subsets of these applications, but we compare performance where possible and demonstrate that mBART consistently performs the best.

We also show that mBART enables new types of transfer across language pairs. For example, fine-tuning on bi-text in one language pair (e.g., Korean-English) creates a model that can translate from all other languages in the monolingual pretraining set (e.g., Italian-English), with no further training. We also show that languages not in the pre-training corpora can benefit from mBART, strongly suggesting that the initialization is at least partially language universal. Finally, we present a detailed analysis of which factors contribute the most to effective pre-training, including the number of languages and their overall similarity.

2 Multilingual Denoising Pre-training

We use the Common Crawl (CC) corpus ($\S2.1$) to pre-train BART models ($\S2.2$). Our experiments in the later sections involve fine-tuning a range of models pre-trained on different subsets ($\S2.3$).

2.1 Data: CC25 Corpus

Datasets We pre-train on 25 languages (CC25) extracted from the CC corpora (Wenzek et al., 2019; Conneau et al., 2019).² CC25 includes languages from different families and with varied amounts of text (Figure 1). Following Lample and Conneau (2019), we re-balanced the corpus by

up/down-sampling text from each language *i* with a ratio λ_i :

$$\lambda_i = \frac{1}{p_i} \cdot \frac{p_i^{\alpha}}{\sum_i p_i^{\alpha}},\tag{1}$$

where p_i is the percentage of each language in CC-25. We use the smoothing parameter $\alpha = 0.7$.

Pre-processing We tokenize with a sentencepiece model (SPM; Kudo and Richardson, 2018) learned on the full CC data that includes 250,000 subword tokens. Although not all of these languages are used for pre-training, this tokenization supports fine-tuning on additional languages. We do not apply additional preprocessing, such as true-casing or normalizing punctuation/ characters.

2.2 Model: mBART

Our models follow the BART (Lewis et al., 2019) Seq2Seq pre-training scheme, as reviewed in this section. Whereas BART was only pretrained for English, we systematically study the effects of pre-training on different sets of languages.

Architecture We use a standard Seq2Seq Transformer architecture (Vaswani et al., 2017), with 12 layers of encoder and 12 layers of decoder with model dimension of 1024 on 16 heads ($\sim 680M$ parameters). We include an additional layer-normalization layer on top of both the encoder and decoder, which we found stabilized training at FP16 precision.

Learning Our training data covers K languages: $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ where each \mathcal{D}_i is a collection of monolingual documents in language i. We (1) assume access to a noising function g, defined below, that corrupts text, and (2) train the model to predict the original text X given g(X). More formally, we aim to maximize \mathcal{L}_{θ} :

$$\mathcal{L}_{\theta} = \sum_{\mathcal{D}_i \in \mathcal{D}} \sum_{X \in \mathcal{D}_i} \log P(X|g(X);\theta) , \quad (2)$$

where X is an instance in language i and the distribution P is defined by the Seq2Seq model.

Noise Function Following Lewis et al. (2019), we use two types of noise in g. We first remove spans of text and replace them with a mask token. We mask 35% of the words in each instance by randomly sampling a span length according to a Poisson distribution ($\lambda = 3.5$). We also permute

²https://github.com/facebookresearch/cc _net.

the order of sentences within each instance. The decoder input is the original text with one position offset. A language id symbol <LID> is used as the initial token to predict the sentence. It is also possible to use other noise types, such as those in Lample et al. (2018c), but we leave the exploration of the optimal noising strategy to future work.

Instance Format For each instance of a batch, we sample a language id symbol <LID>, and we pack as many consecutive sentences as possible sampled from the corresponding corpus of <LID>, until either it hits the document boundary or reaches the 512 max token length. Sentences in the instance are separated by the end of sentence () token. Then, we append the selected <LID> token to represent the end of this instance. Pre-training at "multi sentence" level enables us to work on both sentence and document translation.

Optimization Our full model (including 25 languages) is trained on 256 Nvidia V100 GPUs (32GB) for 500K steps. The total batch size is around 128K tokens per GPU, matching BART (Lewis et al., 2019) configuration. We use the Adam optimizer ($\epsilon = 1e-6$, $\beta_2 = 0.98$) and linear learning rate decay scheduling. The total training time was approximately 2.5 weeks. We started the training with dropout 0.1 and reduced it to 0.05 at 250K steps and 0 at 400K steps. All experiments are done with Fairseq (Ott et al., 2019).

Reproducibility One potential issue of the proposed approach is the replicability problem due to the requirement of massive monolingual corpora and computational resources, with fine-grained selection on hyper-parameters during pre-training. It is likely to get slightly different fine-tuning performance if we re-train the system again. Tackling on this, we will release the pre-trained checkpoints as well as the code with full instructions for pre-training a new model.

Related Work: XLM(-R) and MASS There are several closely related approaches of multilingual pre-training for machine translation. XLM (Lample and Conneau, 2019) and XLM-R (Conneau et al., 2019) pretrain BERT (Devlin et al., 2019; Liu et al., 2019) in a multilingual fashion, and the resulted parameters can be used to initialize the translation model encoder. Different from XLM(-R), mBART simultaneously pretrains the encoder and the decoder due to the Seq2Seq setup, which is more natural to adapt to machine translation applications.

Similar to mBART, MASS (Song et al., 2019) is also a Seq2Seq-based pre-training technique with "word-masking". However, the decoder of MASS only predicted tokens that was masked in the encoder, whereas mBART reconstructs the full target sequence which allows to apply not only "masking" but any possible noise functions.

Furthermore, both XLM and MASS did not show evidence of the pre-trained models improving translation performance over two languages.

2.3 Pre-trained Models

To better measure the effects of different levels of multilinguality during pre-training, we built a range of models as follows:

- **mBART25** We pre-train a model on all 25 languages, using the setting described in §2.2.
- mBART06 To explore the effect of pretraining on related languages, we pretrain a model on a subset of six European languages: Ro, It, Cs, Fr, Es, and En. For a fair comparison, we use ~ 1/4 of the mBART25 batch size, which allows our model to have the same number of updates per language during pre-training.
- mBART02 We pre-train bilingual models, using English and one other language for four language pairs: En-De, En-Ro, En-It. We use a batch size of ~ 1/12 of that in the mBART25.
- **BART-En/Ro** To help establish a better understanding towards multilingual pretraining, we also train monolingual BART models on the En and Ro corpus only, respectively.
- **Random** As additional baselines, we will also include a comparison with a model randomly initialized without pre-training for each translation task. Because the sizes of different downstream datasets vary, we always grid-search the hyper-parameters (architecture, dropout, etc.) to find the best non-pretrained configuration.



Figure 2: Framework for our multilingual denoising pre-training (left) and fine-tuning on downstream MT tasks (right), where we use (1) sentence permutation and (2) word-span masking as the injected noise. A special language id token is added at both the encoder and decoder. One multilingual pre-trained model is used for all tasks.

All models use the same vocabulary (§2.1). Not all tokens will frequently occur in all pre-training corpora, but later experiments show that this large vocabulary can improve generalization in multilingual settings even for unseen languages.

2.4 Scaling-up Matters

Scaling-up the training data and model parameters has been a key factor in pre-training (Devlin et al., 2019; Conneau et al., 2019; Raffel et al., 2019). Compared to conventional semi-supervised methods (e.g., back-translation) and other pretraining for MT (Lample and Conneau, 2019; Song et al., 2019), we pre-train mBART on much more monolingual data with relatively deeper architecture. This scale, in combination with the new multi-lingual training, is central to our results (sections 3 to 5), although future work could more carefully study the relative contributions of each.

3 Sentence-level Machine Translation

This section shows that mBART pre-training provides consistent performance gains in low to medium resource sentence-level MT settings, including bi-text only and with back translation, and outperforms other existing pre-training schemes (\S 3.2). We also present a detailed analysis to understand better which factors contribute the most to these gains (\S 3.3), and show that pre-training can even improve performance for languages not present in the pre-training data (\S 3.4).

3.1 Experimental Settings

Datasets We gather 24 pairs of publicly available parallel corpora that cover all the languages in CC25 (Figure 1). Most pairs are from previous WMT (Gu, Kk, Tr, Ro, Et, Lt, Fi, Lv, Cs, Es, Zh, De, Ru, Fr \leftrightarrow En) and IWSLT (Vi, Ja, Ko, Nl, Ar, It \leftrightarrow En) competitions. We also use FLoRes pairs (Guzmán et al., 2019, En-Ne and En-Si), En-Hi from IITB (Kunchukuttan et al., 2017), and En-My from WAT19 (Ding et al., 2018, 2019). We divide the datasets into three categories—low resource (<1M sentence pairs), medium resource (>1M and <10M), and high resource (>10M).

Fine-tuning & Decoding We fine-tune mBART on a single pair of bi-text data, feeding the source language into the encoder and decoding the target language. As shown in Figure 2, we load the pre-trained weights and train the MT model on bi-texts with teacher forcing. For all directions, we train with 0.3 dropout, 0.2 label smoothing, 2500 warm-up steps, 3e-5 maximum learning rate. We use a maximum of 40K training updates for all low and medium resource pairs and 100K for high resource pairs. The final models are selected based on validation likelihood. We use beam-search with beam size 5 for decoding. Our initial experiments indicate that the fine-tuning process is generally stable with different seeds. Therefore, to reduce the total computation, all our results are reported with single execution. We validate the statistical significance with scripts from the *mosesdecoder*.³

³https://github.com/moses-smt/mosesdecoder /blob/master/scripts/analysis/bootstrap -hypothesis-difference-significance.pl.

| Languages Data Source | En- WM | ·Gu IT19 | En- WM | -Kk IT19 | En IWS | -Vi LT15 | En WM | -Tr [T17 | En IWS | -Ja LT17 | En- IWS | -Ko LT17 |
|--------------------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|
| Size | 10 |)K | 91 | K | 13 | 3K | 20 | 7K | 22 | 3K | 23 | 0K |
| Direction | \leftarrow | \rightarrow |
| Random | 0.0 | 0.0 | 0.8 | 0.2 | 23.6 | 24.8 | 12.2 | 9.5 | 10.4 | 12.3 | 15.3 | 16.3 |
| mBART25 | 0.3 | 0.1 | 7.4 | 2.5 | 36.1 | 35.4 | 22.5 | 17.8 | 19.1 | 19.4 | 24.6 | 22.6 |
| Languages | En | -NI | En | -Ar | En | -It | En- | My | En | -Ne | En | -Ro |
| Data Source | IWS | LT17 | IWS | LT17 | IWS | LT17 | WA | T19 | FLa | Res | WM | IT16 |
| Size | 23 | 7K | 25 | 0K | 25 | 0K | 25 | 9K | 56 | 4K | 60 | 8K |
| Direction | \leftarrow | \rightarrow |
| Random | 34.6 | 29.3 | 27.5 | 16.9 | 31.7 | 28.0 | 23.3 | 34.9 | 7.6 | 4.3 | 34.0 | 34.3 |
| mBART25 | 43.3 | 34.8 | 37.6 | 21.6 | 39.8 | 34.0 | 28.3 | 36.9 | 14.5 | 7.4 | 37.8 | 37.7 |
| Languages | En | -Si | En | -Hi | En | -Et | En | -Lt | En | -Fi | En | -Lv |
| Data Source | FLa | Res | IT | ТВ | WM | T18 | WM | T19 | WM | T17 | WM | IT17 |
| Size | 64 | 7K | 1.5 | 6M | 1.9 | 4M | 2.1 | 1 M | 2.6 | 6M | 4.5 | 0M |
| Direction | \leftarrow | \rightarrow |
| Random | 7.2 | 1.2 | 10.9 | 14.2 | 22.6 | 17.9 | 18.1 | 12.1 | 21.8 | 20.2 | 15.6 | 12.9 |
| mBART25 | 13.7 | 3.3 | 23.5 | 20.8 | 27.8 | 21.4 | 22.4 | 15.3 | 28.5 | 22.4 | 19.3 | 15.9 |

Table 1: Low/medium resource machine translation Pre-training consistently improves over a randomly initialized baseline, with particularly large gains on low resource language pairs (e.g., Vi-En).

3.2 Main Results

As shown in Table 1, initializing with the pretrained mBART25 weights shows gains on all the low and medium resource pairs when compared with randomly initialized baselines. We observe gains of 12 or more BLEU points on low resource pairs such as En-Vi, En-Tr, and noisily aligned pairs like En-Hi. Fine-tuning still fails in extremely low-resource cases such as En-Gu, which have $\sim 10k$ examples. In these settings, unsupervised translation is more appropriate, see $\S5.2$. For high resource cases (Table 2), we do not observe consistent gains, and pretraining slightly hurts performance when more than 25M parallel sentences are available. When a significant amount of bi-text data is given, we suspect that supervised training washes out the pre-trained weights.

Note that some reported runs of our baseline systems using the vanilla Transformers with randomly initialized weights have considerably noticeable gaps between the SoTA systems reported in the original competitions.⁴ The difference is mainly because we train and search

| Languages Size | Cs 11M | Es 15M | Zh 25M | De 28M | Ru 29M | Fr 41M |
|-------------------|------------------|------------------|------------------|------------------|------------------|---------------|
| RANDOM | 16.5 | 33.2 | 35.0 | 30.9 | 31.5 | 41.4 |
| мBART25 | 18.0 | 34.0 | 33.3 | 30.5 | 31.3 | 41.0 |

Table 2: High resource machine translation where all the datasets are from their latest WMT competitions. We only evaluate our models on En-X translation.

the hyper-parameters for baselines on officially provided bitext only without using any monolingual corpus or multilingual adaptation. For instance, the SoTA score for $En \rightarrow Gu$ is 28.2 in WMT19, compared with 0 in Table 1. It is basically because the quality of the original bitext data is low, and the SoTA systems commonly used additional languages such as Hi to boost the performance. Similar gaps can also be observed in pairs such as Kk-En and Lt-En, where Ru as the additional language is also crucial. The main purpose of this part is to discuss the effects of multilingual pre-training in a constrained bitext setting for a better comparison. We will include more discussions of combining

⁴http://matrix.statmt.org/.



Figure 3: Pre-training + back translation on FLoRes with two iterations of BT.

multilingual translation with pretraining in future work.

Plus Back-Translation Back-translation (BT; Sennrich et al., 2016) is a standard approach to augment bi-text with target-side monolingual data. We combine our pre-training with BT and test it on low resource language pairs—En-Si and En-Ne—using the FLoRes dataset (Guzmán et al., 2019). We use the same monolingual data as Guzmán et al. (2019) to generate BT data. Figure 3 shows that initializing the model with our mBART25 pre-trained parameters improves BLEU scores at each iteration of back translation, resulting in new state-of-the-art results in all four translation directions. It indicates that the pretrained mBART weights can be directly plugged into existing pipeline using BT.

Compared with Other **Pre-training** Approaches We also compare our pre-trained models with recent self-supervised pre-training methods, as shown in Table 3. We consider En-Ro translation, the only pair with established results. Our mBART model outperforms all the other pre-trained models, both with and without BT augmentation. We also show comparisons with the conventional BART model trained on the same En and Ro data only. Both have improvements over baselines, although worse than mBART results, indicating that pre-training in a multilingual setting is essential. Moreover, combining BT leads to additional gains, resulting in a new state-of-the-art for Ro-En translation.

3.3 Analysis

We also present additional analyses, to better quantify when our pre-training helps.

How many languages should you pre-train on? We investigate when it is helpful for pre-training to include languages other than the targeted language pair that will be used during fine tuning. Table 4

| Pre-traini | ng | Fine-tuning | | | | | |
|---------------------|-------|-----------------------|-------------------------------------|------|--|--|--|
| Model | Data | $En {\rightarrow} Ro$ | $\textbf{Ro}{\rightarrow}\text{En}$ | +BT | | | |
| RANDOM | None | 34.3 | 34.0 | 36.8 | | | |
| XLM (2019) | En Ro | - | 35.6 | 38.5 | | | |
| MASS (2019) | En Ro | - | - | 39.1 | | | |
| BART (2019) | En | - | - | 38.0 | | | |
| XLM-R (2019) | CC100 | 35.6 | 35.8 | - | | | |
| BART-En | En | 36.0 | 35.8 | 37.4 | | | |
| BART-Ro | Ro | 37.6 | 36.8 | 38.1 | | | |
| MBART02 | En Ro | 38.5 | 38.5 | 39.9 | | | |
| MBART25 | CC25 | 37.7 | 37.8 | 38.8 | | | |

Table 3: Comparison with other pre-training approaches on WMT16 Ro-En.

| Languages | De | Ro | It | My | En |
|-----------|------|------|------|------|-------|
| Size/GB | 66.6 | 61.4 | 30.2 | 1.6 | 300.8 |
| mBART02 | 31.3 | 38.5 | 39.7 | 36.5 | |
| mBART06 | _ | 38.5 | 39.3 | _ | |
| mBART25 | 30.5 | 37.7 | 39.8 | 36.9 | |

Table 4: Pretraining languages on En-X translation. The size refers to the size of monolingual data for X. The size of En is shown as reference. All the pretrained models were controlled to see the same number of English instances during training.

shows performance on four X-En pairs. Pretraining on more languages helps most when the target language monolingual data is limited (e.g., En-My, where the size of My is around 0.5%of En).

In contrast, when monolingual data is plentiful (De, Ro), pre-training on multiple languages slightly hurts the final results (<1 BLEU). In these cases, additional languages may reduce the capacity available for each test language. Additionally, the fact that mBART06 performs similar to mBART02 on Ro-En suggests that pre-training with similar languages is particularly helpful.



Figure 4: Fine-tuning curves for Ro-En along with Pre-training steps. Both mBART25 and mBART02 outperform the best baseline system after 25K steps.



Figure 5: Fine-tuning curves for En-De along with size of bitext. The x-axis is on a log scale.

How many pre-training steps are needed? We plot Ro-En BLEU score vs. Pre-training steps in Figure 4, where we take the saved checkpoints (every 25K steps) and apply the same fine-tuning process described in §3.1. Without any pre-training, our model overfits and performs much worse than the baseline. However, after just 25K steps (5% of training), both models outperform the best baseline. The models keep improving by over 3 BLEU for the rest of pre-training and have not fully converged after 500K steps. In addition, mBART25 is consistently slightly worse than mBART02, which confirms the observation in Table 4.

How much bi-text is needed? Tables 1 and 2 show that pre-training consistently improves for low and medium resource language pairs. To verify this trend, we plot performance for differing sized subsets of the En-De dataset. More precisely, we take the full En-De corpus (28M pairs) and randomly sample 10K, 50K, 100K, 500K, 1M, 5M, 10M datasets. We compare performance

without pre-training to the mBART02 results, as shown in Figure 5. The pre-trained model is able to achieve over 20 BLEU with only 10K training examples, whereas the baseline system scores 0. Unsurprisingly, increasing the size of bitext corpus improves both models. Our pre-trained model consistently outperforms the baseline models, but the gap reduces with increasing amounts of bi-text, especially after 10M sentence pairs. This result confirms our observation in §3.2 that our pre-training does not help translation in high-resource pairs.

3.4 Generalization to Languages NOT in Pre-training

In this section, we show that mBART can improve performance even for languages that did not appear in the pre-training corpora, suggesting that the pre-training has language universal aspects. Similar phenomena have also been reported in other multilingual pre-training approaches in other NLP applications (Pires et al., 2019; Wang et al., 2019; Artetxe et al., 2019).

Experimental Settings We report results finetuning for three pairs, NI-En, Ar-En, and De-NI, using the pre-trained mBART25, mBART06, and mBART02 (EnRo) models. The mBART06 and mBART02 models are not pre-trained on Ar, De or NI text, but all languages are in mBART25. Both De and NI are European languages and are related to En, Ro, and the other languages in the mBART06 pre-training data.

Results As shown in Table 5, we find large gains from pre-training on English-Romanian, even when translating a distantly related unseen language (Arabic) and two unseen languages (German and Dutch). The best results are achieved when pre-training includes both test languages, although pre-training on other languages is surprisingly competitive.

Unseen Vocabularies Arabic is distantly related to the languages in mBART02 and mBART06, and has a disjoint character set. This means that its word embeddings are largely not estimated during pre-training. However, we obtain similar improvements on Ar-En pairs to those on Nl-En. This result suggests that the pre-trained Transformer layers learn universal properties of language that generalize well even with minimal lexical overlap.

| | Monolingual | NI-En | En-Nl | Ar-En | En-Ar | NI-De | De-Nl |
|-----------------|-------------------|-------------|-------------|--------------|-------------|-------------|-------------|
| RANDOM | None | 34.6 (-8.7) | 29.3 (-5.5) | 27.5 (-10.1) | 16.9 (-4.7) | 21.3 (-6.4) | 20.9 (-5.2) |
| мBART02 | En Ro | 41.4 (-2.9) | 34.5 (-0.3) | 34.9 (-2.7) | 21.2 (-0.4) | 26.1 (-1.6) | 25.4 (-0.7) |
| м BART06 | En Ro Cs It Fr Es | 43.1 (-0.2) | 34.6 (-0.2) | 37.3 (-0.3) | 21.1 (-0.5) | 26.4 (-1.3) | 25.3 (-0.8) |
| мBART25 | All | 43.3 | 34.8 | 37.6 | 21.6 | 27.7 | 26.1 |

Table 5: Generalization to unseen languages Language transfer results, fine-tuning on language-pairs without pre-training on them. mBART25 uses all languages during pre-training, while other settings contain at least one unseen language pair. For each model, we also show the gap to mBART25 results.

Unseen Source or Target Languages Table 5 shows different performance when the unseen languages are on the source side, target side, or both sides. If both sides are unseen, the performance (in terms of difference from mBART25) is worse than where at least one language is seen during pre-training. Furthermore, although the En-X pairs perform similarly, mBART06 outperforms mBART02 on X-En pairs. Fine-tuning unseen languages on the source side is more difficult, and is worthy of extensive future study.

4 Document-level Machine Translation

We evaluate mBART on document-level machine translation tasks, where the goal is to translate segments of text that contain more than one sentence (up to an entire document). During pre-training, we use document fragments of up to 512 tokens, allowing the models to learn dependencies between sentences. We show that this pre-training significantly improves document-level translation.

4.1 Experimental Settings

Datasets We evaluate performance on two common document-level MT datasets: WMT19 En-De and TED15 Zh-En. For En-De, we use the document data from WMT19 to train our model, without any additional sentence-level data. The Zh-En dataset is from IWSLT 2014 and 2015 (Cettolo et al., 2012, 2015). Following Miculicich et al. (2018), we use 2010-2013 TED as the test set.

Pre-processing We pre-process with the approach used in pre-training. For each block, sentences are separated by end of sentence symbols () and the entire instance is ended with the specific language id (<LID>). On average, documents are split into 2–4 instances.

Fine-tuning & Decoding We use the same finetuning scheme as for sentence-level translation (§3.1), without using any task-specific techniques developed by previous work (Miculicich et al., 2018; Li et al., 2019), such as constrained contexts or restricted attention. For decoding, we simply pack the source sentences into blocks, and translate each instance block autoregressively. The model does not know how many sentences to generate in advance and decoding stops when <LID> is predicted. We use beam size 5 by default.

Baselines & Evaluation We train 4 models: a document-level (Doc-) MT model (§4.1) and a corresponded sentence-level (Sent-) MT model $(\S3.1)$ as the baseline, both with and without pre-training. We use mBART25 as the common pre-trained model for En-De and Zh-En. For En-De, even though our mBART25 Doc-MT model decodes multiple sentences together, the translated sentences can be aligned to the source sentences, which allows us to evaluate BLEU scores both on sentence-level (s-BLEU) and document-level (d-BLEU).⁵ For Zh-En, however, we cannot produce the same number of translated sentences as the reference due to alignment errors in the test data. We only provide the d-BLEU scores on this direction.

We also compare our models with Hierarchical Attention Networks (HAN, Miculicich et al., 2018) on Zh-En, which is the state-of-theart non-pretraining approach for document-level translation for this pair. They combine two layers of attention—first within and then across sentences.

⁵Standard BLEU scores match n-grams at sentence-level. We also consider document-level where we match n-grams over the whole document resulting in a slightly higher score.

⁽a) Sentence- and Document-level BLEU scores on En-De

| Model | Random | | mBART25 | | Madal | Random | mBART25 | HAN (2018) |
|---------|--------|--------|---------|--------|---------|--------|---------|------------|
| | s-BLEU | d-BLEU | s-BLEU | d-BLEU | Widdei | d-BLEU | d-BLEU | d-BLEU |
| Sent-MT | 34.5 | 35.9 | 36.4 | 38.0 | Sent-MT | 22.0 | 28.4 | _ |
| Doc-MT | × | 7.7 | 37.1 | 38.5 | Doc-MT | 3.2 | 29.6 | 24.0 |

Table 6: Document-level machine translation on En-De and Zh-En. (\times) The randomly initialized Doc-MT model cannot produce translations aligned to the original sentences, so only document evaluation is possible.

4.2 Main Results

Table 6 shows the main results for both En-De and Zh-En at both sentence-level and document-level.

Random vs. Pre-trained The MT models initialized with pre-trained weights outperform randomly initialized models by large margins, for both sentence-level and document-level training. Our mBART25 models (both Sent-MT and Doc-MT) also outperform HAN (Miculicich et al., 2018),⁶ despite the fact that they are not customized for document-level MT.

Sent-MT vs. Doc-MT For En-De and En-Zh, the mBART25 Doc-MT models outperform mBART25 fine-tuned at sentence-level by large margins, reversing the trend seen for models without pre-training. For both datasets, randomly initialized Doc-MT fails to work, resulting in much worse results than the sentence-level models. Such large performance gaps indicate that pre-training is *critical* for document level performance. It is in general difficult to collect high-quality document-level data in large quantities, suggesting that pre-training may be a strong strategy for future work. We also include a sampled example in Figure 6.

5 Unsupervised Machine Translation

In addition to supervised machine translation, we also evaluate our model on tasks where no bi-text is available for the target language pair. We define three types of *unsupervised* translation:

1. No bi-text of any kind. A common solution is to learn from back-translation (Artetxe et al., 2017; Lample et al., 2018c). We show that mBART provides a simple and effective initialization scheme for these methods (§5.1). 2. No bi-text for the target pair, but both languages appear in bi-text corpora with other pairs. This setup is common for multilingual MT systems (Johnson et al., 2017; Gu et al., 2019). In this paper, we limit our focus to building models for single language pairs, and leave discussions for multilingual MT to future work.

(b) Document-level BLEU scores on Zh-En

3. No bi-text for the target pair is available, but there is bi-text for translating from some other language into the target language. mBART supports effective transfer, even if the source language has no bi-text of any form (§5.2).

5.1 Unsupervised Machine Translation via Back-Translation

Datasets We evaluate our pre-trained models on En-De, En-Ne, and En-Si. En and De are both European languages sharing many sub-words, whereas Ne and Si are quite distinct from En. We use the same test sets as supervised benchmarks §3.1, and use the same pre-training data (CC25) for back-translation to avoid introducing new information.

Learning Following Lample and Conneau (XLM, 2019), we initialize the translation model with the mBART weights, and then learn to predict the monolingual sentences conditioned on source sentences generated by on-the-fly BT. Furthermore, we constrain mBART to only generating tokens in target language⁷ for the first 1000 steps of on-the-fly BT, to avoid it copying the source text.

Results Table 7 shows the unsupervised translation results compared with non-pretrained models, as well as models with existing pre-training methods. Our models achieve large gains over non-pretrained models for all directions, and

⁶d-BLEU is recomputed from the provided system output.

⁷We mask out the output probability of predicted tokens which appear less than 1% in the target monolingual corpus.

| SOURCE | 作为一名艺术家、联系对我来说是非常重要的。通过我的艺术作品我试着阐明人类不是与自然分隔开 而是每一件事都是互相联系的。大约10年前我第一次去了南极洲,我 也第一次看到了冰山,我感到敬畏,我的心快速地研动,头晕目眩,试着理解在我面前的这到底是什么。在我身边的冰山,浮出水面几乎200英尺,我只能感到很奇怪 这就是 一片雪花,覆盖在另一片雪花,年复一年形成的,冰山的形成是当它们从冰川能裂开或者从冰架上断裂开,每个冰山都有它们自己的想特个性。它们与其周边的环境和其 情况的互动具有一个鲜明的方式。有些冰山拒绝妥协 坚持到底,而另一些冰山就不能忍受 在一时剧烈激情喷涌下就水崩冰裂。当你看到冰山,很容易就想到 它们都是孤立 的,它们是独立的,单独一体的,更像是我们人类有时候对自身的看法。但现实远不止这个。随着冰山融化,我呼吸到 它古老的气味,随着冰山融化,它释放了富有矿 物质的鲜水 它们滋养了万物。我着手拍损这些冰山 好似我在拍摄我相先的骨像,了解到在这些个别的时刻 冰山原是以那样方式存在 但再也不会像那样存在了。当它们酸化 时,这绝不是死亡;也绝不是一个终结,而是通往生生不息之能的一个延续,我拍损过的冰山,有些冰是非常有轻一几千年年龄。有些冰超过十万年。我想给大家展示的 最后图片 是我在格陵兰岛的 Kekertsuatsiak上拍摄的一个冰山。这是一个非常难得的机会 大家实际上得以见证 一个冰山翻滚。所以这就如图所示。在左边你能看到一个小 船,这是一个约15英尺的船。我想让你注意 冰山的形状 它在水面上的变形。在这儿你看到它开始翻滚,小船移动到另一边,一个男人站在那里。这是一个平均尺寸的格陵 兰冰山,它浮出水面大约有120英尺高或者40米高。这视频是实时拍摄的。就像这冰山,它们最示给大家的是其个性的不同方面。谢谢。 |
|--------------------|--|
| Random SENT-MT | As an artist, connection is very important for me. Through my artwork, I try to illustrate that human beings are not separated from nature, but each one of them is interconnected. About 10 years ago, I first went to Antarctica, and I saw mountains for the first time. I felt awe. My heart was rapidly wiped out, trying to figure out what was going on in front of me. The iceberg that floats around me almost 200 feet of water, and I can only feel strange that this snow covers another piece of snow for a year and years. The form of the iceberg is that when they break apart from the glacier or breaking off the shelves of ice. Every iceberg has their own personality. They're interacting with their surrounding environment in a very different way. Some of the ice mountains refused to compromise, and some other mountains of ice can't endure, and the water collapses during a violent ice. And when you look at the iceberg, it's easy to think that they're all isolated, and they're independent, they're independent in these individual moments of ice is there in a way that's there, but it's never going to exist again. When they they not diseas, it's not the act, it's not deact, it's |
| Random DOC-MT | As an artist, as an artist. I'm going to focus on the glacier and the glacier and the glacier. There's a lot of ice in the ice, and there's a lot of ice in the ice, and there's a lot of ice in the ice, and there's a lot of ice in the ice, and there's a lot of ice in the ice, and there's a lot of ice in the ice and there's a lot of ice in the ice and there's a lot of ice in the ice and there's a lot of ice in the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ice that's melted from the ice that's melted from the ice of the ic |
| mBART25 SENT-MT | As an artist, connection is very important to me. And through my artwork, I try to illustrate that humans are not separate from nature, but that everything is connected to each other. I went to Antarctica for the first time about 10 years ago, and I saw the iceberg for the first time. I felt awe. And my heart was pounding, my head was pounding, trying to figure out what was going on in front of me. And the iceberg that was right next to me was floating almost 200 feet above the surface of the water, and I just couldn't help but feel that this was a snowflake covering another snowflake, forming over and over and over again. The formation of glaciers is when they break off from glaciers, or break off from ice shelves. Every iceberg has its own unique personality. They interact in a very distinctive way with the environment around them and with the circumstances around them. Some glaciers refused to compromise and insisted, and some glaciers couldn't stand the heat of passion as it poured down the ice. When you look at icebergs, it's easy to think that they're isolated, that they're independent, that they're separate, that they're more like the way we think about ourselves sometimes. But the reality is much more than that. And as the glacier mets, I breathe in its ancient smell. And as the glacier mets, it releases fresh water of minerals that nourish everything. I started photographing these icebergs like I was photographing my ancestors, and I learned that in these individual moments, the iceberg parstel, I photographed glaciers, and some of the mere very young thousands of years old. Some of the ice has been there for more than 100,000 years. And the last picture I want to show you is an iceberg that I photographed in Kekersuatsiak on the island. It's not about detert, it's not about the shape of the cieberg changes as it moves over the surface. And here you see it rolling, and the boat moves to the other side, and a man is standing there. This is an avery difficult opportunity to actually witness the rolli |
| mBART25 DOC-MT | And as an artist, connection is very important to me. Through my artwork, I try to convey the idea that humans are not separated from nature, but that everything is connected to each other. When I first went to Antarctica about 10 years ago, I saw for the first time icebergs. <u>And</u> I felt awe. My heart was shaking, my head was shaking, trying to understand what was in front of me. The icebergs around me were floating almost 200 feet above the surface of the water, and I could only feel how strange it was that this was a snowflake covering another snowflake, forming over and over again over and over again. <u>And</u> icebergs form when they break off from ice shelves. <u>And</u> each iceberg has its own unique personality. They interact in a very distinctive way with the environment around them and with the circumstances in which they're located. Some icebergs fuse to settle down, and some icebergs can't stand the heat of passion that pours down and breaks ice. <u>And</u> when you look at icebergs, it's easy to think that they're isolated, that they're independent, that they're individual, that they're more like the way we think about ourselves sometimes. But the reality is much more than that. As the icebergs like I'm photographing my ancestors, and I'm learning that in these individual moments, icebergs used to exist in that way and will never be the same again. When they melt, it's not about the oth, but it's about a continuation of a lifetime. <u>And</u> the icebergs I've photographed, some of them are very young thousands of years old. <u>And</u> went give that so very difficult opportunity for you to actually witness the rolling of a iceberg. <u>And</u> it floats about 12 potographed on Kekertsuatsiak in Greenland. <u>And</u> it's a very difficult opportunity for you to actually witness the rolling of a iceberg. <u>And</u> it floats about 12 peteotgraphed on Kekertsuatsiak in Greenland. <u>And</u> it's a very difficult opportunity for you to actually witness the rolling of a iceberg. <u>Sonere</u> it is. On the left you can see a little boat. It's a li |
| TARGET | As an artist, connection is very important to me. Through my work I'm trying to articulate that humans are not separate from nature and that everything is interconnected. I first went to Antarctica almost 10 years ago, where I saw my first icebergs. I was in awe. My heart beat fast, my head was dizzy, trying to comprehend what it was that stood in front of me. The icebergs around me were almost 200 feet out of the water, and I could only help but wonder that this was one snowflake on top of another snowflake, year after year. Icebergs are born when they calve off of glaciers or break off of ice shekes. Each iceberg has its own individual personality. They have a distinct way of interacting with their environment and their experiences. Some refuse to give up and hold on to the bitter end, while others can't take it anymore and crumble in a fit of dramatic passion. It's easy to think, when you look at an iceberg mets, I am breathing in its ancient atmosphere. As the iceberg metts, it is releasing mineral-rich fresh water that nourishes many forms of life. I approach photographing these icebergs as if I'm making portraits of my ancestors, knowing that in these individual moments they exist in that way and will never exist that way gain. It is not a death when they melt; it is not an end, but a continuation of the ice parts of an iceberg rolling. So here it is. You can see on the left side a small boat. That's about 15's out of take you to pay attention to the shape of the iceberg and where it is at the waterline. You can see here, it begins to roll, and the boat has moved to the other side, and the main standing there. This is an average-size Greenlandic iceberg. It's about 120 feet above the water, or 40 meters. And this video is real time. And just like that, the iceberg shows you a different side of its personality. Thank you. |

Figure 6: An example of document-level translation from mBART25 Sent-MT and Doc-MT, held out from the test set of TED15 Zh-En. The Doc-MT system produces much fluent and coherent translation, which is closer to the reference translation. For instance, Doc-MT model produces several "<u>And</u>" to connect sentences to make it reads better, while the Sent-MT model does not contain global knowledge and produce sentences independently. Additionally, both systems produce much better translations than models without pre-training where the non-pretrained Doc-MT model completely fails to produce readable translation output.

| | En-De | | En- | Ne | En-Si | | |
|-------------|--------------|---------------|--------------|---------------|--------------|---------------|--|
| | \leftarrow | \rightarrow | \leftarrow | \rightarrow | \leftarrow | \rightarrow | |
| Random | 21.0 | 17.2 | 0.0 | 0.0 | 0.0 | 0.0 | |
| XLM (2019) | 34.3 | 26.4 | 0.5 | 0.1 | 0.1 | 0.1 | |
| MASS (2019) | 35.2 | 28.3 | - | _ | - | - | |
| mBART | 34.0 | 29.8 | 10.0 | 4.4 | 8.2 | 3.9 | |

Table 7: Unsupervised MT via BT between dis-similar languages.

outperform XLM significantly for dissimilar pairs (En-Ne, En-Si) where the existing approaches completely fail. For En-De, our model also performs comparably against XLM and MASS.

5.2 Unsupervised Machine Translation via Language Transfer

We also report results when the target language appears in a bi-text with some other source language.

Datasets We only consider $X \rightarrow En$ translation, and choose the bitexts of 12 language pairs from §3.1, covering Indic languages (Ne, Hi, Si, Gu), European languages (Ro, It, Cs, Nl), East Asian languages (Zh, Ja, Ko), and Arabic (Ar).

Results The pre-trained mBART25 model is fine-tuned on each language pair, and then evaluated on the rest of pairs, as seen in Table 8. We also present the direct fine-tuning performance (§3) on the diagonal, for reference. We see transfer for all pairs with all fine-tuned models except from Gu-En where the supervised model completely fails (0.3 BLEU). In some cases we can achieve similar (Cs-En) or even much better (Ne-En, Gu-En) results compared with the supervised results. We also show an example of language transfer in Figure 7.

As a comparison, we also apply the same procedure on randomly initialized models without pretraining, which always ends up with ≈ 0 BLEU. This indicates that multilingual pre-training is essential and produces universal representations across languages, so that once the model learns to translate one language to En, it learns to translate all languages with similar representations.

When is language transfer useful? Table 8 also shows that the size of transfer effects varies with the similarity of different languages. First, for

most pairs, language transfer works better when fine-tuning is also conducted in the same language family, especially between Indic languages (Hi, Ne, Gu). However, significant vocabulary sharing is not required for effective transfer. For instance, Zh-En and It-En achieve the best transfer learning results on Ko-En and Ar-En, respectively. This is despite the low vocabulary overlap (even character overlap) between (Zh, Ko) and (It, Ar).

With BT We present a comparison of unsupervised MT with BT vs. language transfer in Table 9 where language transfer works better when there exists a close language translation to transfer from.

Moreover, we show promising results for combining these two techniques. We start from the best transferred model and apply (iterative) BT on the same monolingual corpus used in pre-training. Table 9 presents the results with 1 iteration of BT. We see improvements for all pairs. The complete analysis of both methods is left as future work.

6 Related Work

Self-supervised Learning for Text Generation This work inherits from the recent success brought by pre-training for NLP applications (Peters et al., 2018; Radford et al., 2018; Devlin et al., 2019; Yang et al., 2019b; Liu et al., 2019), especially for text generation (Radford et al., 2019; Song et al., 2019; Dong et al., 2019; Raffel et al., 2019; Lewis et al., 2019). The pre-trained models are usually used as the initialization for finetuning downstream tasks such as controllable language modeling (Shirish Keskar et al., 2019), summarization (Song et al., 2019; Liu and Lapata, 2019) and dialogue generation (Zhang et al., 2019).

Specifically for machine translation, unsupervised pre-training methods were also explored to improve the performance. Qi et al. (2018) investigated the application of pre-trained word embeddings for MT; Ramachandran et al. (2017) proposed to pre-train the encoder-decoder modules as two separate language models. Yang et al. (2019a); Zhu et al. (2020) explored fusion approaches to incorporate the pre-trained BERT weights to improve NMT training. In contrast to most prior work, we focus on pre-training one denoising autoencoder, and adapt the weights of the entire model for various MT applications.

| | | | | | | Fine | -tuning | Langua | ges | | | | |
|-----|-------|------|------|------|------|------|---------|--------------|------|------|------|------|------|
| | | Zh | Ja | Ko | Cs | Ro | NI | It | Ar | Hi | Ne | Si | Gu |
| Do | omain | News | TED | TED | News | News | TED | TED | TED | News | Wiki | Wiki | Wiki |
| | Ζн | 23.7 | 8.8 | 9.2 | 2.8 | 7.8 | 7.0 | 6.8 | 6.2 | 7.2 | 4.2 | 5.9 | 0.0 |
| ŝ | JA | 9.9 | 19.1 | 12.2 | 0.9 | 4.8 | 6.4 | 5.1 | 5.6 | 4.7 | 4.2 | 6.5 | 0.0 |
| age | Ко | 5.8 | 16.9 | 24.6 | 5.7 | 8.5 | 9.5 | 9.1 | 8.7 | 9.6 | 8.8 | 11.1 | 0.0 |
| gu | Cs | 9.3 | 15.1 | 17.2 | 21.6 | 19.5 | 17.0 | 16.7 | 16.9 | 13.2 | 15.1 | 16.4 | 0.0 |
| 'an | Ro | 16.2 | 18.7 | 17.9 | 23.0 | 37.8 | 22.3 | 21.6 | 22.6 | 16.4 | 18.5 | 22.1 | 0.0 |
| 10 | Nl | 14.4 | 30.4 | 32.3 | 21.2 | 27.0 | 43.3 | 3 4.1 | 31.0 | 24.6 | 23.3 | 27.3 | 0.0 |
| tin | Iт | 16.9 | 25.8 | 27.8 | 17.1 | 23.4 | 30.2 | 39.8 | 30.6 | 20.1 | 18.5 | 23.2 | 0.0 |
| les | Ar | 5.8 | 15.5 | 12.8 | 12.7 | 12.0 | 14.7 | 14.7 | 37.6 | 11.6 | 13.0 | 16.7 | 0.0 |
| - | Ηı | 3.2 | 10.1 | 9.9 | 5.8 | 6.7 | 6.1 | 5.0 | 7.6 | 23.5 | 14.5 | 13.0 | 0.0 |
| | NE | 2.1 | 6.7 | 6.5 | 5.0 | 4.3 | 3.0 | 2.2 | 5.2 | 17.9 | 14.5 | 10.8 | 0.0 |
| | Sı | 5.0 | 5.7 | 3.8 | 3.8 | 1.3 | 0.9 | 0.5 | 3.5 | 8.1 | 8.9 | 13.7 | 0.0 |
| | Gu | 8.2 | 8.5 | 4.7 | 5.4 | 3.5 | 2.1 | 0.0 | 6.2 | 13.8 | 13.5 | 12.8 | 0.3 |

Table 8: Unsupervised MT via language transfer on X-En translations. The model fine-tuned on one language pair is directly tested on another. We use gray color to show the direct fine-tuning results, and lightgray color to show language transfer within similar language groups. We **bold** the highest transferring score for each pair.

| SOURCE Ja | カナダやアメリカ その他の多くの先進国では 当たり前のことかもしれませんが 貧しい国々や 家父長社会、部族社会では 就 学とは女の子にとって 一大事です |
|------------------|---|
| TARGET En | It may be taken for granted in Canada, in America, in many developed countries, but in poor countries, in patriarchal societies, in tribal societies, it's a big event for the life of girl. |
| mBART25 Ja-En | In Canada, in the United States, and many other developed countries, it's taken for granted that in poor countries, in patriarchal societies, in tribal societies, education is very important for girls. |
| mBART25 Ko-En | It's commonplace in countries like Canada and the United States and many other先進 countries, but it's not commonplace in poor countries, in patriarchal societies, in clan societies, where schooling is a big deal for girls. |
| mBART25 Zh-En | It's commonplace in Canada, in the U.S., and in many other countries in the world, in poor countries, in patriarchal societies, in ethnic societies, that education is a priority for girls. |

Figure 7: An example of unsupervised MT via language transfer. mBART models finetuned with **Ko** or **Zh** are able to translate **Ja** sentence to **En** almost as correctly as in the supervised case.

| Source | online BT | Transfer | | Combined |
|--------|-----------|----------------|-----|----------|
| Ro | 30.5 | 23.0 (| Cs) | 33.9 |
| Ne | 10.0 | 17.9 (1 | Hi) | 22.1 |
| Zh | 11.3 | 9.2 (1 | Ko) | 15.0 |
| NI | 28.5 | 34.1 (| It) | 35.4 |

Table 9: BT vs. language transfer for unsupervised MT for X-En translations. For language transfer, we present the best transferring scores together with the language transferred from.

Multilinguality in NLP tasks This work is also related to the continual trend of multilingual language learning, including aligning multilingual word embeddings (Mikolov et al., 2013; Chen and Cardie, 2018; Lample et al., 2018b) into

universal space, and learning crosslingual models (Wada and Iwata, 2018; Lample and Conneau, 2019; Conneau et al., 2019) to exploit shared representations across languages.

For MT, the most relevant field is *multilingual translation* (Firat et al., 2016; Johnson et al., 2017; Aharoni et al., 2019; Arivazhagan et al., 2019) where the ultimate goal is to jointly train one translation model that translates multiple language directions at the same time, and shares representations to improve the translation performance on low-resource languages (Gu et al., 2018). In this paper, we focus on multilingualism in the pre-training stage and fine-tune the learned model in the standard bilingual scenario.

Compared with multilingual translation, we do not require parallel data across multiple languages but the targeted direction, which improves the scalability to low-resource languages and specific domains.

Document Translation As one of the key applications, our work is also related to previous efforts for incorporating document-level context into neural machine translation (Wang et al., 2017; Jean et al., 2017; Tiedemann and Scherrer, 2017; Miculicich et al., 2018; Tu et al., 2018). Li et al. (2019) is the most relevant work that also utilized pre-trained encoder (BERT) for handling longer context. However, the focus has been on designing new task-specific techniques, and doing sentence-level translation with a wider input context. To the best of our knowledge, our multilingual pre-trained model is the first that shows improved results on document-level translation with standard Seq2Seq models.

Unsupervised Translation This work also summarizes the previous efforts of learning to translate between languages without a direct parallel corpus. When no parallel data of any kind is available, Artetxe et al. (2017) and Lample et al. (2018a) proposed to jointly learn denoising auto-encoder and back-translation from both directions, which, however, required good initialization and only worked well on similar language pairs. Wu et al. (2019) solve the problem by mining sentences from Wikipedia and using them as weakly supervised translation pairs. Similar to Lample and Conneau (2019) and Song et al. (2019), we follow the first approach and treat our pre-trained model as the initialization step. We also investigate unsupervised translation using language transfer, which is similar to Pourdamghani et al. (2019), where the authors generate translationese of the source language and train a system on high-resource languages to correct these intermediate utterances. It is also closely related to Conneau et al. (2018) and Artetxe et al. (2019) for cross-lingual representation learning where we also show representation learned by mBART can be easily transferred between language without supervised data.

7 Conclusion

We demonstrate that multilingual de-noising pretraining is able to significantly improve both supervised and unsupervised machine translation at both the sentence level and document level. We analyze when and how pre-training is most effective and can be combined with other approaches such as back-translation. Our results also show the transfer learning ability of the learned representations from multilingual pre-training.

In future work, we will scale-up the current pre-training to more languages, for example, an mBART100 model. The size of our model makes it expensive to deploy in production—future work will explore pre-training more efficient models.

Acknowledgments

We thank Marc'Aurelio Ranzato, Guillaume Lample, Alexis Conneau, and Michael Auli for sharing their expertise on low-resource and unsupervised machine translation and Peng-Jen Chen and Jiajun Shen for details about FloRes and WAT datasets. We also thank our colleagues at FAIR and FAIAR for valuable feedback.

References

- Roee 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 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, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2017. Unsupervised neural machine translation. *arXiv preprint arXiv: 1710.11041*. **DOI:** https://doi.org/18653 /v1/D18-1399

- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2019. On the cross-lingual transferability of monolingual representations. *arXiv preprint arXiv:1910.11856*. **DOI:** https:// doi.org/10.18653/v1/2020.acl -main.421
- Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. Wit3: Web inventory of transcribed and translated talks. In *Conference* of European Association for Machine Translation, pages 261–268.
- Mauro Cettolo, Niehues Jan, Stüker Sebastian, Luisa Bentivogli, Roldano Cattoni, and Marcello Federico. 2015. The IWSLT 2015 evaluation campaign. In *International Workshop on Spoken Language Translation*.
- Xilun Chen and Claire Cardie. 2018. Unsupervised multilingual word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 261–270. Brussels, Belgium. Association for Computational Linguistics. **DOI:** https:// doi.org/10.18653/v1/D18-1024
- 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. arXiv preprint arXiv: 1911.02116. DOI: https://doi.org/10.18653/v1/2020.acl-main.747
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. **DOI:** https://doi .org/10.18653/v1/D18-1269
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *North American Association for Computational Linguistics (NAACL)*.
- Chenchen Ding, Hnin Thu Zar Aye, Win Pa Pa, Khin Thandar Nwet, Khin Mar Soe, Masao

Utiyama, and Eiichiro Sumita. 2019. Towards Burmese (Myanmar) morphological analysis: Syllable-based tokenization and part-of-speech tagging. ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP), 19(1):5. DOI: https://doi .org/10.1145/3325885

- Chenchen Ding, Masao Utiyama, and Eiichiro Sumita. 2018. NOVA: A feasible and flexible annotation system for joint tokenization and part-of-speech tagging. ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP), 18(2):17. DOI: https://doi.org/10.1145/3276773
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. *arXiv* preprint arXiv:1905.03197.
- Sergey Edunov, Alexei Baevski, and Michael Auli. 2019. Pre-trained language model representations for language generation. *arXiv preprint arXiv:1903.09722*. **DOI:** https://doi.org /10.18653/v1/N19-1409
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016. Multi-way, multilingual neural machine translation with a shared attention mechanism. In *NAACL*. **DOI:** https://doi.org/10.18653/v1/N16-1101
- Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O. K. Li. 2018. Universal neural machine translation for extremely low resource languages. 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 344–354. New Orleans, Louisiana. Association for Computational Linguistics.
- Jiatao Gu, Yong Wang, Kyunghyun Cho, and Victor O. K. Li. 2019. Improved zero-shot neural machine translation via ignoring spurious correlations. *arXiv preprint arXiv:1906.01181*.
- Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio

Ranzato. 2019. The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English. 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 6097–6110. Hong Kong, China. Association for Computational Linguistics. **DOI:** https://doi.org/10.18653/v1/D19 -1632

- Sébastien Jean, Stanislas Lauly, Orhan Firat, and Kyunghyun Cho. 2017. Does neural machine translation benefit from larger context? *CoRR*, abs/1704.05135.
- 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. Googles multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351. **DOI:** https://doi.org/10 .1162/tacl_a_00065
- Taku Kudo and John Richardson. 2018. Sentence-Piece: 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. DOI: https://doi.org/10.18653/v1/D18 -2012, PMID: 29382465
- Anoop Kunchukuttan, Pratik Mehta, and Pushpak Bhattacharyya. 2017. The IIT Bombay English-Hindi parallel corpus. *CoRR*, abs/1710.02855.
- Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining. *arXiv preprint arXiv:1901.07291*.
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018a. Unsupervised machine translation using monolingual corpora only. In *International Conference on Learning Representations*.
- Guillaume Lample, Alexis Conneau, Marc'Aurelio Ranzato, Ludovic Denoyer,

and Hervé Jégou. 2018b. Word translation without parallel data. In *International Conference on Learning Representations*.

- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018c. Phrase-based & neural unsupervised machine translation. *arXiv preprint arXiv:* 1804.07755. DOI: https://doi.org/10 .18653/v1/D18-1549
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*. **DOI:** https://doi.org /10.18653/v1/2020.acl-main.703
- Liangyou Li, Xin Jiang, and Qun Liu. 2019. Pretrained language models for documentlevel neural machine translation. *arXiv preprint arXiv:1911.03110*.
- Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*. **DOI:** https:// doi.org/10.18653/v1/D19-1387
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. ROBERTA: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2947–2954. Brussels, Belgium. Association for Computational Linguistics. **DOI**: https://doi.org/10.18653/v1/D18–1325
- Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. 2013. Exploiting similarities among languages for machine translation. *CoRR*, abs/1309.4168.
- 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 North American Association for Computational Linguistics (NAACL): System Demonstrations. **DOI:** https://doi.org /10.18653/v1/N19-4009

- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In North American Association for Computational Linguistics (NAACL).
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? *arXiv preprint arXiv:1906.01502*. **DOI:** https://doi.org/10.18653/v1/P19 -1493
- Nima Pourdamghani, Nada Aldarrab, Marjan Ghazvininejad, Kevin Knight, and Jonathan May. 2019. Translating translationese: A twostep approach to unsupervised machine translation. In ACL. DOI: https://doi.org /10.18653/v1/P19-1293
- Ye Qi, Devendra Singh Sachan, Matthieu Felix, Sarguna Janani Padmanabhan, and Graham Neubig. 2018. When and why are pretrained word embeddings useful for neural machine translation? *arXiv preprint arXiv:* 1804.06323. DOI: https://doi.org/10 .18653/v1/N18-2084
- Alec Radford, Karthik Narasimhan, Time Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning, OpenAI.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners, OpenAI.
- 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. *arXiv preprint arXiv:1910.10683*.
- Prajit Ramachandran, Peter J Liu, and Quoc Le. 2017. Unsupervised pretraining for sequence to sequence learning. In *Proceedings of the 2017*

Conference on Empirical Methods in Natural Language Processing, pages 383-391. DOI: https://doi.org/10.18653/v1/D17 -1039

- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 86–96. Berlin, Germany. Association for Computational Linguistics. **DOI:** https://doi .org/10.18653/v1/P16-1009
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv* preprint arXiv:1909.05858.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. MASS: Masked sequence to sequence pre-training for language generation. In *International Conference on Machine Learning (ICML)*.
- Jörg Tiedemann and Yves Scherrer. 2017. Neural machine translation with extended context. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 82–92. Copenhagen, Denmark. Association for Computational Linguistics. **DOI:** https://doi.org/10.18653/v1/W17 -4811
- Zhaopeng Tu, Yang Liu, Shuming Shi, and Tong Zhang. 2018. Learning to remember translation history with a continuous cache. *Transactions of the Association for Computational Linguistics*, 6:407–420. **DOI:** https://doi.org/10.1162/tacl_a_00029
- 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*.
- Takashi Wada and Tomoharu Iwata. 2018. Unsupervised cross-lingual word embedding by multilingual neural language models. *CoRR*, abs/1809.02306. DOI: https://doi.org /10.18653/v1/P19-1300

- 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. DOI: https://doi.org/10 .18653/v1/D17-1301
- Zihan Wang, Stephen Mayhew, Dan Roth, and others. 2019. Cross-lingual ability of multilingual bert: An empirical study. *arXiv preprint arXiv:1912.07840*.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzman, Armand Joulin, and Edouard Grave. 2019. CCNET: Extracting high quality monolingual datasets from web crawl data. *arXiv preprint arXiv:1911.00359*.
- Lijun Wu, Jinhua Zhu, Di He, Fei Gao, Xu Tan, Tao Qin, and Tie-Yan Liu. 2019. Machine translation with weakly paired bilingual documents.

- Jiacheng Yang, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Yong Yu, Weinan Zhang, and Lei Li. 2019a. Towards making the most of bert in neural machine translation. *arXiv preprint arXiv:1908.05672*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019b. XLNet: Generalized autoregressive pretraining for language understanding. arXiv preprint arXiv:1906.08237.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2019. DialoGPT: Large-scale generative pre-training for conversational response generation.
- Jinhua Zhu, Yingce Xia, Lijun Wu, Di He, Tao Qin, Wengang Zhou, Houqiang Li, and Tie-Yan Liu. 2020. Incorporating BERT into neural machine translation. *arXiv preprint arXiv:2002.06823*. **DOI:** https://doi .org/10.18653/v1/2020.acl-demos.30