

# Understanding Back-Translation at Scale

Sergey Edunov<sup>△</sup> Myle Ott<sup>△</sup> Michael Auli<sup>△</sup> David Grangier<sup>▽\*</sup>

<sup>△</sup>Facebook AI Research, Menlo Park, CA & New York, NY.

<sup>▽</sup>Google Brain, Mountain View, CA.

## Abstract

An effective method to improve neural machine translation with monolingual data is to augment the parallel training corpus with back-translations of target language sentences. This work broadens the understanding of back-translation and investigates a number of methods to generate synthetic source sentences. We find that in all but resource poor settings back-translations obtained via sampling or noised beam outputs are most effective. Our analysis shows that sampling or noisy synthetic data gives a much stronger training signal than data generated by beam or greedy search. We also compare how synthetic data compares to genuine bitext and study various domain effects. Finally, we scale to hundreds of millions of monolingual sentences and achieve a new state of the art of 35 BLEU on the WMT'14 English-German test set.

## 1 Introduction

Machine translation relies on the statistics of large parallel corpora, i.e. datasets of paired sentences in both the source and target language. However, bitext is limited and there is a much larger amount of monolingual data available. Monolingual data has been traditionally used to train language models which improved the fluency of statistical machine translation (Koehn, 2010).

In the context of neural machine translation (NMT; Bahdanau et al. 2015; Gehring et al. 2017; Vaswani et al. 2017), there has been extensive work to improve models with monolingual data, including language model fusion (Gulcehre et al., 2015, 2017), back-translation (Sennrich et al., 2016a) and dual learning (Cheng et al., 2016; He et al., 2016a). These methods have different advantages and can be combined to reach high accuracy (Hassan et al., 2018).

We focus on back-translation (BT) which operates in a semi-supervised setup where both bilingual and monolingual data in the target language are available. Back-translation first trains an intermediate system on the parallel data which is used to translate the target monolingual data into the source language. The result is a parallel corpus where the source side is *synthetic* machine translation output while the target is genuine text written by humans. The synthetic parallel corpus is then simply added to the real bitext in order to train a final system that will translate from the source to the target language. Although simple, this method has been shown to be helpful for phrase-based translation (Bojar and Tamchyna, 2011), NMT (Sennrich et al., 2016a; Poncelas et al., 2018) as well as unsupervised MT (Lample et al., 2018a).

In this paper, we investigate back-translation for neural machine translation at a large scale by adding hundreds of millions of back-translated sentences to the bitext. Our experiments are based on strong baseline models trained on the public bitext of the WMT competition. We extend previous analysis (Sennrich et al., 2016a; Poncelas et al., 2018) of back-translation in several ways. We provide a comprehensive analysis of different methods to generate synthetic source sentences and we show that this choice matters: sampling from the model distribution or noising beam outputs outperforms pure beam search, which is typically used, by 1.7 BLEU on average across several test sets. Our analysis shows that synthetic data based on sampling and noised beam search provides a stronger training signal than synthetic data based on argmax inference. We also study how adding synthetic data compares to adding real bitext in a controlled setup with the surprising finding that synthetic data can sometimes match the accuracy of real bitext. Our best setup achieves 35 BLEU on the WMT'14 English-German test set by rely-

\*Work done while at Facebook AI Research.

ing only on public WMT bitext as well as 226M monolingual sentences. This outperforms the system of DeepL by 1.7 BLEU who train on large amounts of high quality non-benchmark data. On WMT’14 English-French we achieve 45.6 BLEU.

## 2 Related work

This section describes prior work in machine translation with neural networks as well as semi-supervised machine translation.

### 2.1 Neural machine translation

We build upon recent work on neural machine translation which is typically a neural network with an encoder/decoder architecture. The encoder infers a continuous space representation of the source sentence, while the decoder is a neural language model conditioned on the encoder output. The parameters of both models are learned jointly to maximize the likelihood of the target sentences given the corresponding source sentences from a parallel corpus (Sutskever et al., 2014; Cho et al., 2014). At inference, a target sentence is generated by left-to-right decoding.

Different neural architectures have been proposed with the goal of improving efficiency and/or effectiveness. This includes recurrent networks (Sutskever et al., 2014; Bahdanau et al., 2015; Luong et al., 2015), convolutional networks (Kalchbrenner et al., 2016; Gehring et al., 2017; Kaiser et al., 2017) and transformer networks (Vaswani et al., 2017). Recent work relies on attention mechanisms where the encoder produces a sequence of vectors and, for each target token, the decoder attends to the most relevant part of the source through a context-dependent weighted-sum of the encoder vectors (Bahdanau et al., 2015; Luong et al., 2015). Attention has been refined with multi-hop attention (Gehring et al., 2017), self-attention (Vaswani et al., 2017; Paulus et al., 2018) and multi-head attention (Vaswani et al., 2017). We use a transformer architecture (Vaswani et al., 2017).

### 2.2 Semi-supervised NMT

Monolingual target data has been used to improve the fluency of machine translations since the early IBM models (Brown et al., 1990). In phrase-based systems, language models (LM) in the target language increase the score of fluent outputs during decoding (Koehn et al., 2003; Brants et al., 2007).

A similar strategy can be applied to NMT (He et al., 2016b). Besides improving accuracy during decoding, neural LM and NMT can benefit from deeper integration, e.g. by combining the hidden states of both models (Gulcehre et al., 2017). Neural architecture also allows multi-task learning and parameter sharing between MT and target-side LM (Domhan and Hieber, 2017).

Back-translation (BT) is an alternative to leverage monolingual data. BT is simple and easy to apply as it does not require modification to the MT training algorithms. It requires training a target-to-source system in order to generate additional synthetic parallel data from the monolingual target data. This data complements human bitext to train the desired source-to-target system. BT has been applied earlier to phrase-based systems (Bogjar and Tamchyna, 2011). For these systems, BT has also been successful in leveraging monolingual data for domain adaptation (Bertoldi and Federico, 2009; Lambert et al., 2011). Recently, BT has been shown beneficial for NMT (Sennrich et al., 2016a; Poncelas et al., 2018). It has been found to be particularly useful when parallel data is scarce (Karakanta et al., 2017).

Currey et al. (2017) show that low resource language pairs can also be improved with synthetic data where the source is simply a copy of the monolingual target data. Concurrently to our work, Imamura et al. (2018) show that sampling synthetic sources is more effective than beam search. Specifically, they sample multiple sources for each target whereas we draw only a single sample, opting to train on a larger number of target sentences instead. Hoang et al. (2018) and Cotterell and Kreutzer (2018) suggest an iterative procedure which continuously improves the quality of the back-translation and final systems. Niu et al. (2018) experiment with a multilingual model that does both the forward and backward translation which is continuously trained with new synthetic data.

There has also been work using source-side monolingual data (Zhang and Zong, 2016). Furthermore, Cheng et al. (2016); He et al. (2016a); Xia et al. (2017) show how monolingual text from both languages can be leveraged by extending back-translation to dual learning: when training both source-to-target and target-to-source models jointly, one can use back-translation in both directions and perform multiple rounds of BT. A simi-

lar idea is applied in unsupervised NMT (Lample et al., 2018a,b). Besides monolingual data, various approaches have been introduced to benefit from parallel data in other language pairs (Johnson et al., 2017; Firat et al., 2016a,b; Ha et al., 2016; Gu et al., 2018).

Data augmentation is an established technique in computer vision where a labeled dataset is supplemented with cropped or rotated input images. Recently, generative adversarial networks (GANs) have been successfully used to the same end (Antoniou et al., 2017; Perez and Wang, 2017) as well as models that learn distributions over image transformations (Haugberg et al., 2016).

### 3 Generating synthetic sources

Back-translation typically uses beam search (Sennrich et al., 2016a) or just greedy search (Lample et al., 2018a,b) to generate synthetic source sentences. Both are approximate algorithms to identify the maximum a-posteriori (MAP) output, i.e. the sentence with the largest estimated probability given an input. Beam is generally successful in finding high probability outputs (Ott et al., 2018a).

However, MAP prediction can lead to less rich translations (Ott et al., 2018a) since it always favors the most likely alternative in case of ambiguity. This is particularly problematic in tasks where there is a high level of uncertainty such as dialog (Serban et al., 2016) and story generation (Fan et al., 2018). We argue that this is also problematic for a data augmentation scheme such as back-translation. Beam and greedy focus on the head of the model distribution which results in very regular synthetic source sentences that do not properly cover the true data distribution.

As alternative, we consider sampling from the model distribution as well as adding noise to beam search outputs. First, we explore unrestricted sampling which generates outputs that are very diverse but sometimes highly unlikely. Second, we investigate sampling restricted to the most likely words (Graves, 2013; Ott et al., 2018a; Fan et al., 2018). At each time step, we select the  $k$  most likely tokens from the output distribution, re-normalize and then sample from this restricted set. This is a middle ground between MAP and unrestricted sampling.

As a third alternative, we apply noising Lample et al. (2018a) to beam search outputs. Adding noise to input sentences has been very benefi-

cial for the autoencoder setups of (Lample et al., 2018a; Hill et al., 2016) which is inspired by denoising autoencoders (Vincent et al., 2008). In particular, we transform source sentences with three types of noise: deleting words with probability 0.1, replacing words by a filler token with probability 0.1, and swapping words which is implemented as a random permutation over the tokens, drawn from the uniform distribution but restricted to swapping words no further than three positions apart.

## 4 Experimental setup

### 4.1 Datasets

The majority of our experiments are based on data from the WMT’18 English-German news translation task. We train on all available bitext excluding the ParaCrawl corpus and remove sentences longer than 250 words as well as sentence-pairs with a source/target length ratio exceeding 1.5. This results in 5.18M sentence pairs. For the back-translation experiments we use the German monolingual newscrawl data distributed with WMT’18 comprising 226M sentences after removing duplicates. We tokenize all data with the Moses tokenizer (Koehn et al., 2007) and learn a joint source and target Byte-Pair-Encoding (BPE; Sennrich et al., 2016) with 35K types. We develop on newstest2012 and report final results on newstest2013-2017; additionally we consider a held-out set from the training data of 52K sentence-pairs.

We also experiment on the larger WMT’14 English-French task which we filter in the same way as WMT’18 English-German. This results in 35.7M sentence-pairs for training and we learn a joint BPE vocabulary of 44K types. As monolingual data we use newscrawl2010-2014, comprising 31M sentences after language identification (Lui and Baldwin, 2012). We use newstest2012 as development set and report final results on newstest2013-2015.

The majority of results in this paper are in terms of case-sensitive tokenized BLEU (Papineni et al., 2002) but we also report test accuracy with de-tokenized BLEU using sacreBLEU (Post, 2018).

### 4.2 Model and hyperparameters

We re-implemented the Transformer model in pytorch using the fairseq toolkit.<sup>1</sup> All experiments

<sup>1</sup>Code available at <https://github.com/pytorch/fairseq>

are based on the Big Transformer architecture with 6 blocks in the encoder and decoder. We use the same hyper-parameters for all experiments, i.e., word representations of size 1024, feed-forward layers with inner dimension 4096. Dropout is set to 0.3 for En-De and 0.1 for En-Fr, we use 16 attention heads, and we average the checkpoints of the last ten epochs. Models are optimized with Adam (Kingma and Ba, 2015) using  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$ , and  $\epsilon = 1e - 8$  and we use the same learning rate schedule as Vaswani et al. (2017). All models use label smoothing with a uniform prior distribution over the vocabulary  $\epsilon = 0.1$  (Szegedy et al., 2015; Pereyra et al., 2017). We run experiments on DGX-1 machines with 8 Nvidia V100 GPUs and machines are interconnected by Infini-band. Experiments are run on 16 machines and we perform 30K synchronous updates. We also use the NCCL2 library and the torch distributed package for inter-GPU communication. We train models with 16-bit floating point operations, following Ott et al. (2018b). For final evaluation, we generate translations with a beam of size 5 and with no length penalty.

## 5 Results

Our evaluation first compares the accuracy of back-translation generation methods (§5.1) and analyzes the results (§5.2). Next, we simulate a low-resource setup to experiment further with different generation methods (§5.3). We also compare synthetic bitext to genuine parallel data and examine domain effects arising in back-translation (§5.4). We also measure the effect of upsampling bitext during training (§5.5). Finally, we scale to a very large setup of up to 226M monolingual sentences and compare to previous research (§5.6).

### 5.1 Synthetic data generation methods

We first investigate different methods to generate synthetic source translations given a back-translation model, i.e., a model trained in the reverse language direction (Section 5.1). We consider two types of MAP prediction: greedy search (greedy) and beam search with beam size 5 (beam). Non-MAP methods include unrestricted sampling from the model distribution (sampling), restricting sampling to the  $k$  highest scoring outputs at every time step with  $k = 10$  (top10) as well as adding noise to the beam outputs (beam+noise). Restricted sampling is a middle-ground between

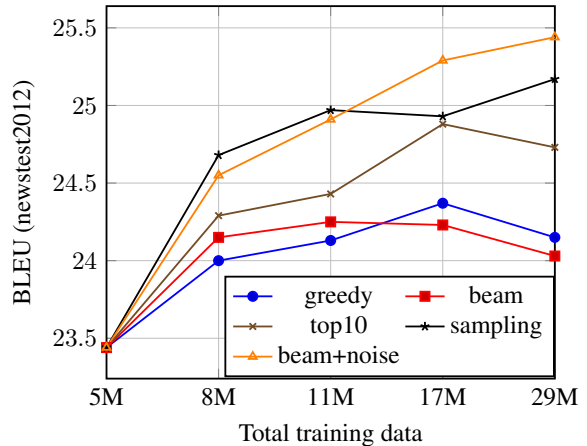


Figure 1: Accuracy of models trained on different amounts of back-translated data obtained with greedy search, beam search ( $k = 5$ ), randomly sampling from the model distribution, restricting sampling over the ten most likely words (top10), and by adding noise to the beam outputs (beam+noise). Results based on newstest2012 of WMT English-German translation.

beam search and unrestricted sampling, it is less likely to pick very low scoring outputs but still preserves some randomness. Preliminary experiments with top5, top20, top50 gave similar results to top10.

We also vary the amount of synthetic data and perform 30K updates during training for the bitext only, 50K updates when adding 3M synthetic sentences, 75K updates for 6M and 12M sentences and 100K updates for 24M sentences. For each setting, this corresponds to enough updates to reach convergence in terms of held-out loss. In our 128 GPU setup, training of the final models takes 3h 20min for the bitext only model, 7h 30min for 6M and 12M synthetic sentences, and 10h 15min for 24M sentences. During training we also sample the bitext more frequently than the synthetic data and we analyze the effect of this in more detail in §5.5.

Figure 1 shows that sampling and beam+noise outperform the MAP methods (pure beam search and greedy) by 0.8-1.1 BLEU. Sampling and beam+noise improve over bitext-only (5M) by between 1.7-2 BLEU in the largest data setting. Restricted sampling (top10) performs better than beam and greedy but is not as effective as unrestricted sampling (sampling) or beam+noise.

Table 1 shows results on a wider range of



	news2013	news2014	news2015	news2016	news2017	Average
bitext	27.84	30.88	31.82	34.98	29.46	31.00
+ beam	27.82	32.33	32.20	35.43	31.11	31.78
+ greedy	27.67	32.55	32.57	35.74	31.25	31.96
+ top10	28.25	33.94	34.00	36.45	32.08	32.94
+ sampling	28.81	34.46	34.87	37.08	32.35	33.51
+ beam+noise	29.28	33.53	33.79	37.89	32.66	33.43

Table 1: Tokenized BLEU on various test sets of WMT English-German when adding 24M synthetic sentence pairs obtained by various generation methods to a 5.2M sentence-pair bitext (cf. Figure 1).

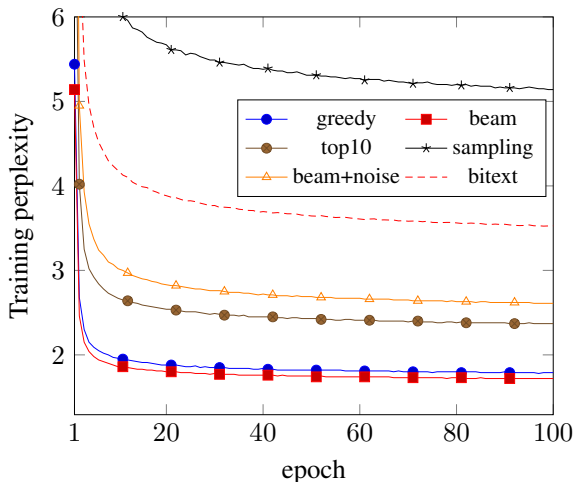


Figure 2: Training perplexity (PPL) per epoch for different synthetic data. We separately report PPL on the synthetic data and the bitext. Bitext PPL is averaged over all generation methods.

test sets (newstest2013-2017). Sampling and beam+noise perform roughly equal and we adopt sampling for the remaining experiments.

## 5.2 Analysis of generation methods

The previous experiment showed that synthetic source sentences generated via sampling and beam with noise perform significantly better than those obtained by pure MAP methods. Why is this?

Beam search focuses on very likely outputs which reduces the diversity and richness of the generated source translations. Adding noise to beam outputs and sampling do not have this problem: Noisy source sentences make it harder to predict the target translations which may help learning, similar to denoising autoencoders (Vincent et al., 2008). Sampling is known to better approximate the data distribution which is richer than the argmax model outputs (Ott et al., 2018a). There-

Perplexity	
human data	75.34
beam	72.42
sampling	500.17
top10	87.15
beam+noise	2823.73

Table 2: Perplexity of source data as assigned by a language model (5-gram Kneser-Ney). Data generated by beam search is most predictable.

fore, sampling is also more likely to provide a richer training signal than argmax sequences.

To get a better sense of the training signal provided by each method, we compare the loss on the training data for each method. We report the cross entropy loss averaged over all tokens and separate the loss over the synthetic data and the real bitext data. Specifically, we choose the setup with 24M synthetic sentences. At the end of each epoch we measure the loss over 500K sentence pairs sub-sampled from the synthetic data as well as an equally sized subset of the bitext. For each generation method we choose the same sentences except for the bitext which is disjoint from the synthetic data. This means that losses over the synthetic data are measured over the *same* target tokens because the generation methods only differ in the source sentences. We found it helpful to up-sample the frequency with which we observe the bitext compared to the synthetic data (§5.5) but we do not upsample for this experiment to keep conditions as similar as possible. We assume that when the training loss is low, then the model can easily fit the training data without extracting much learning signal compared to data which is harder to fit.

Figure 2 shows that synthetic data based on

source	Diese gegenstzlichen Auffassungen von Fairness liegen nicht nur der politischen Debatte zugrunde.
reference	These competing principles of fairness underlie not only the political debate.
beam	These conflicting interpretations of fairness are not solely based on the political debate.
sample	<i>Mr President</i> , these contradictory interpretations of fairness are not based solely on the political debate.
top10	Those conflicting interpretations of fairness are not solely at the heart of the political debate.
beam+noise	conflicting BLANK interpretations BLANK are of not BLANK based on the political debate.

Table 3: Example where sampling produces inadequate outputs. "Mr President," is not in the source. BLANK means that a word has been replaced by a filler token.

greedy or beam is much easier to fit compared to data from sampling, top10, beam+noise and the bitext. In fact, the perplexity on beam data falls below 2 after only 5 epochs. Except for sampling, we find that the perplexity on the training data is somewhat correlated to the end-model accuracy (cf. Figure 1) and that all methods except sampling have a lower loss than real bitext.

These results suggest that synthetic data obtained with argmax inference does not provide as rich a training signal as sampling or adding noise. We conjecture that the regularity of synthetic data obtained with argmax inference is not optimal. Sampling and noised argmax both expose the model to a wider range of source sentences which makes the model more robust to reordering and substitutions that happen naturally, even if the model of reordering and substitution through noising is not very realistic.

Next we analyze the richness of synthetic outputs and train a language model on real human text and score synthetic source sentences generated by beam search, sampling, top10 and beam+noise. We hypothesize that data that is very regular should be more predictable by the language model and therefore receive low perplexity. We eliminate a possible domain mismatch effect between the language model training data and the synthetic data by splitting the parallel corpus into three non-overlapping parts:

1. On 640K sentences pairs, we train a back-translation model,
2. On 4.1M sentence pairs, we take the source side and train a 5-gram Kneser-Ney language model (Heafield et al., 2013),

3. On the remaining 450K sentences, we apply the back-translation system using beam, sampling and top10 generation.

For the last set, we have genuine source sentences as well as synthetic sources from different generation techniques. We report the perplexity of our language model on all versions of the source data in Table 2. The results show that beam outputs receive higher probability by the language model compared to sampling, beam+noise and real source sentences. This indicates that beam search outputs are not as rich as sampling outputs or beam+noise. This lack of variability probably explains in part why back-translations from pure beam search provide a weaker training signal than alternatives.

Closer inspection of the synthetic sources (Table 3) reveals that sampled and noised beam outputs are sometimes not very adequate, much more so than MAP outputs, e.g., sampling often introduces target words which have no counterpart in the source. This happens because sampling sometimes picks highly unlikely outputs which are harder to fit (cf. Figure 2).

### 5.3 Low resource vs. high resource setup

The experiments so far are based on a setup with a large bilingual corpus. However, in resource poor settings the back-translation model is of much lower quality. Are non-MAP methods still more effective in such a setup? To answer this question, we simulate such setups by sub-sampling the training data to either 80K sentence-pairs or 640K sentence-pairs and then add synthetic data from sampling and beam search. We compare these smaller setups to our original 5.2M sentence bitext configuration. The accuracy of the

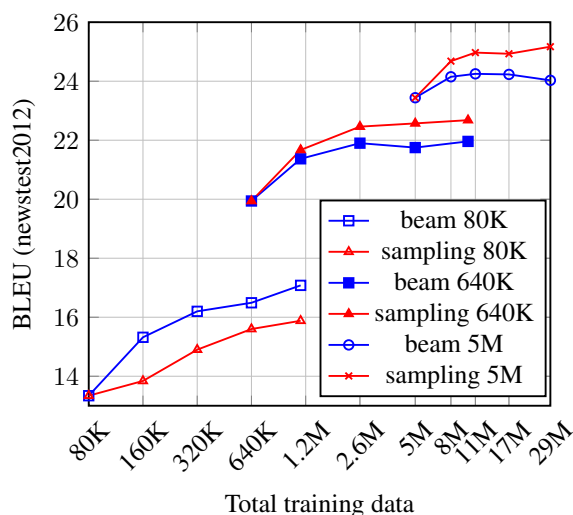


Figure 3: BLEU when adding synthetic data from beam and sampling to bitext systems with 80K, 640K and 5M sentence pairs.

German-English back-translation systems steadily increases with more training data: On newstest2012 we measure 13.5 BLEU for 80K bitext, 24.3 BLEU for 640K and 28.3 BLEU for 5M.

Figure 3 shows that sampling is more effective than beam for larger setups (640K and 5.2M bitexts) while the opposite is true for resource poor settings (80K bitext). This is likely because the back-translations in the 80K setup are of very poor quality and the noise of sampling and beam+noise is too detrimental for this brittle low-resource setting. When the setup is very small the very regular MAP outputs still provide useful training signal while the noise from sampling becomes harmful.

#### 5.4 Domain of synthetic data

Next, we turn to two different questions: How does real human bitext compare to synthetic data in terms of final model accuracy? And how does the domain of the monolingual data affect results?

To answer these questions, we subsample 640K sentence-pairs of the bitext and train a back-translation system on this set. To train a forward model, we consider three alternative types of data to add to this 640K training set. We either add:

- the remaining parallel data (bitext),
- the back-translated target side of the remaining parallel data (BT-bitext),
- back-translated newscrawl data (BT-news).

The back-translated data is generated via sampling. This setup allows us to compare synthetic data to genuine data since BT-bitext and bitext share the same target side. It also allows us to estimate the value of BT data for domain adaptation since the newscrawl corpus (BT-news) is pure news whereas the bitext is a mixture of europarl and commoncrawl with only a small news-commentary portion. To assess domain adaptation effects, we measure accuracy on two held-out sets:

- newstest2012, i.e. pure newswire data.
- a held-out set of the WMT training data (valid-mixed), which is a mixture of europarl, commoncrawl and the small news-commentary portion.

Figure 4 shows the results on both validation sets. Most strikingly, BT-news performs almost as well as bitext on newstest2012 (Figure 4a) and improves the baseline (640K) by 2.6 BLEU. BT-bitext improves by 2.2 BLEU, achieving 83% of the improvement with real bitext. This shows that synthetic data can be nearly as effective as real human translated data when the domains match.

Figure 4b shows the accuracy on valid-mixed, the mixed domain valid set. The accuracy of BT-news is not as good as before since the domain of the BT data and the test set do not match. However, BT-news still improves the baseline by up to 1.2 BLEU. On the other hand, BT-bitext matches the domain of valid-mixed and improves by 2.7 BLEU. This trails the real bitext by only 1.3 BLEU and corresponds to 67% of the gain achieved with real human bitext.

In summary, synthetic data performs remarkably well, coming close to the improvements achieved with real bitext for newswire test data, or trailing real bitext by only 1.3 BLEU for valid-mixed. In absence of a large parallel corpus for news, back-translation therefore offers a simple, yet very effective domain adaptation technique.

#### 5.5 Upsampling the bitext

We found it beneficial to adjust the ratio of bitext to synthetic data observed during training. In particular, we tuned the rate at which we sample data from the bitext compared to synthetic data. For example, in a setup of 5M bitext sentences and 10M synthetic sentences, an upsampling rate of 2 means that we double the frequency at which we

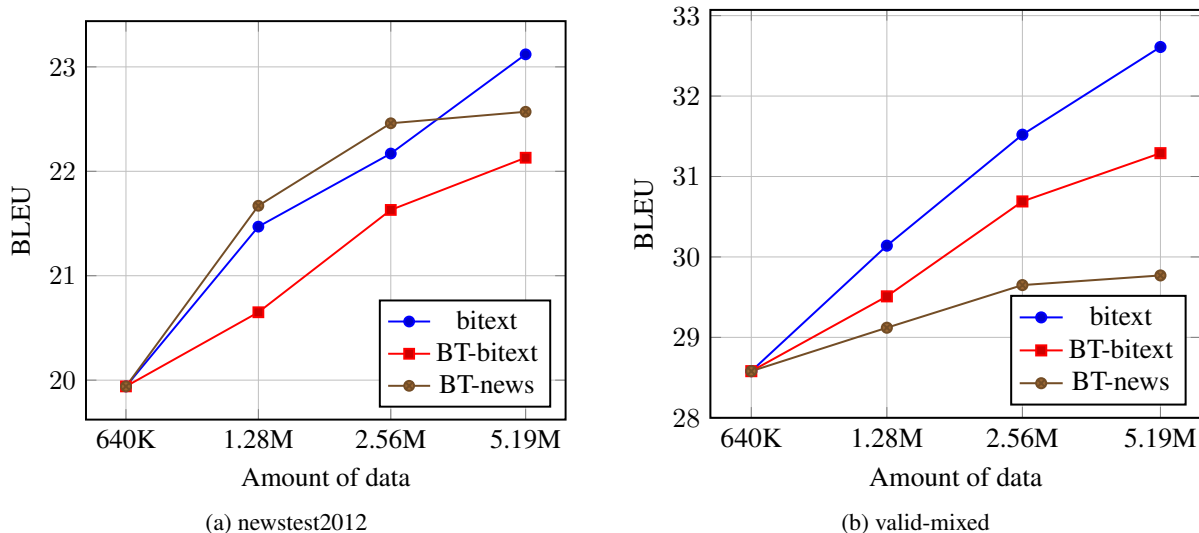


Figure 4: Accuracy on (a) newstest2012 and (b) a mixed domain valid set when growing a 640K bitext corpus with (i) real parallel data (bitext), (ii) a back-translated version of the target side of the bitext (BT-bitext), (iii) or back-translated newscrawl data (BT-news).

visit bitext, i.e. training batches contain on average an equal amount of bitext and synthetic data as opposed to 1/3 bitext and 2/3 synthetic data.

Figure 5 shows the accuracy of various upsampling rates for different generation methods in a setup with 5M bitext sentences and 24M synthetic sentences. Beam and greedy benefit a lot from higher rates which results in training more on the bitext data. This is likely because synthetic beam and greedy data does not provide as much training signal as the bitext which has more variation and is harder to fit. On the other hand, sampling and beam+noise require no upsampling of the bitext, which is likely because the synthetic data is already hard enough to fit and thus provides a strong training signal (§5.2).

## 5.6 Large scale results

To confirm our findings we experiment on WMT’14 English-French translation where we show results on newstest2013-2015. We augment the large bitext of 35.7M sentence pairs by 31M newscrawl sentences generated by sampling. To train this system we perform 300K training updates in 27h 40min on 128 GPUs; we do not up-sample the bitext for this experiment. Table 4 shows tokenized BLEU and Table 5 shows detokenized BLEU.<sup>2</sup> To our knowledge, our baseline

<sup>2</sup>sacreBLEU signatures: BLEU+case.mixed+lang.en-fr+numrefs.1+smooth.exp+test.SET+tok.13a+version.1.2.7 with SET ∈ {wmt13, wmt14/full, wmt15}

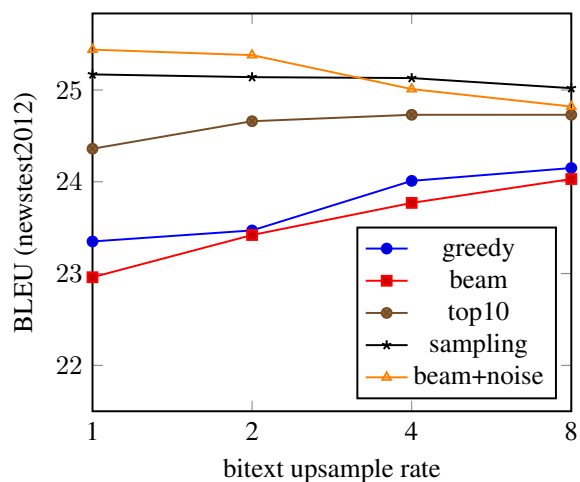


Figure 5: Accuracy when changing the rate at which the bitext is upsampled during training. Rates larger than one mean that the bitext is observed more often than actually present in the combined bitext and synthetic training corpus.

is the best reported result in the literature for newstest2014, and back-translation further improves upon this by 2.6 BLEU (tokenized).

Finally, for WMT English-German we train on all 226M available monolingual training sentences and perform 250K updates in 22.5 hours on 128 GPUs. We upsample the bitext with a rate of 16 so that we observe every bitext sentence

<sup>3</sup>sacreBLEU signatures: BLEU+case.mixed+lang.en-LANG+numrefs.1+smooth.exp+test.wmt14/full+tok.13a+version.1.2.7 with LANG ∈ {de,fr}



	news13	news14	news15
bitext	36.97	42.90	39.92
+sampling	<b>37.85</b>	<b>45.60</b>	<b>43.95</b>

Table 4: Tokenized BLEU on various test sets for WMT English-French translation.

	news13	news14	news15
bitext	35.30	41.03	38.31
+sampling	<b>36.13</b>	<b>43.84</b>	<b>40.91</b>

Table 5: De-tokenized BLEU (sacreBLEU) on various test sets for WMT English-French.

16 times more often than each monolingual sentence. This results in a new state of the art of 35 BLEU on newstest2014 by using only WMT benchmark data. For comparison, DeepL, a commercial translation engine relying on high quality bilingual training data, achieves 33.3 tokenized BLEU.<sup>4</sup> Table 6 summarizes our results and compares to other work in the literature. This shows that back-translation with sampling can result in high-quality translation models based on benchmark data only.

## 6 Conclusions and future work

Back-translation is a very effective data augmentation technique for neural machine translation. Generating synthetic sources by sampling or by adding noise to beam outputs leads to higher accuracy than argmax inference which is typically used. In particular, sampling and noised beam outperforms pure beam by 1.7 BLEU on average on newstest2013-2017 for WMT English-German translation. Both methods provide a richer training signal for all but resource poor setups. We also find that synthetic data can achieve up to 83% of the performance attainable with real bitext. Finally, we achieve a new state of the art result of 35 BLEU on the WMT’14 English-German test set by using publicly available benchmark data only.

In future work, we would like to investigate an end-to-end approach where the back-translation model is optimized to output synthetic sources that are most helpful to the final forward model.

<sup>4</sup><https://www.deepl.com/press.html>

	En-De	En-Fr
a. Gehring et al. (2017)	25.2	40.5
b. Vaswani et al. (2017)	28.4	41.0
c. Ahmed et al. (2017)	28.9	41.4
d. Shaw et al. (2018)	29.2	41.5
DeepL	33.3	<b>45.9</b>
Our result	<b>35.0</b>	45.6
<i>detok. sacreBLEU</i> <sup>3</sup>	33.8	43.8

Table 6: BLEU on newstest2014 for WMT English-German (En-De) and English-French (En-Fr). The first four results use only WMT bitext (WMT’14, except for b, c, d in En-De which train on WMT’16). DeepL uses proprietary high-quality bitext and our result relies on back-translation with 226M newscrawl sentences for En-De and 31M for En-Fr. We also show de-tokenized BLEU (SacreBLEU).

## References

- Karim Ahmed, Nitish Shirish Keskar, and Richard Socher. 2017. Weighted transformer network for machine translation. *arxiv*, 1711.02132.
- Antreas Antoniou, Amos J. Storkey, and Harrison Edwards. 2017. Data augmentation generative adversarial networks. *arXiv*, abs/1711.04340.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations (ICLR)*.
- Nicola Bertoldi and Marcello Federico. 2009. Domain adaptation for statistical machine translation with monolingual resources. In *Workshop on Statistical Machine Translation (WMT)*.
- Ondrej Bojar and Ales Tamchyna. 2011. Improving translation model by monolingual data. In *Workshop on Statistical Machine Translation (WMT)*.
- Thorsten Brants, Ashok C. Popat, Peng Xu, Franz Josef Och, and Jeffrey Dean. 2007. Large language models in machine translation. In *Conference on Natural Language Learning (CoNLL)*.
- Peter F. Brown, John Cocke, Stephen Della Pietra, Vincent J. Della Pietra, Frederick Jelinek, John D. Lafferty, Robert L. Mercer, and Paul S. Roossin. 1990. A statistical approach to machine translation. *Computational Linguistics*, 16:79–85.
- Yong Cheng, Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Semi-supervised learning for neural machine translation. In *Conference of the Association for Computational Linguistics (ACL)*.

- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Ryan Cotterell and Julia Kreutzer. 2018. Explaining and generalizing back-translation through wake-sleep. *arXiv preprint arXiv:1806.04402*.
- Anna Currey, Antonio Valerio Miceli Barone, and Kenneth Heafield. 2017. Copied Monolingual Data Improves Low-Resource Neural Machine Translation. In *Proc. of WMT*.
- Tobias Domhan and Felix Hieber. 2017. Using target-side monolingual data for neural machine translation through multi-task learning. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Angela Fan, Yann Dauphin, and Mike Lewis. 2018. Hierarchical neural story generation. In *Conference of the Association for Computational Linguistics (ACL)*.
- Orhan Firat, Kyunghyun Cho, and Yoshua Bengio. 2016a. Multi-way, multilingual neural machine translation with a shared attention mechanism. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Orhan Firat, Baskaran Sankaran, Yaser Al-Onaizan, Fatos T. Yarman-Vural, and Kyunghyun Cho. 2016b. Zero-resource translation with multi-lingual neural machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *International Conference of Machine Learning (ICML)*.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv*, 1308.0850.
- Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O. K. Li. 2018. Universal neural machine translation for extremely low resource languages. *arXiv*, 1802.05368.
- Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loic Barrault, Hui-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. *arXiv*, 1503.03535.
- Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, and Yoshua Bengio. 2017. On integrating a language model into neural machine translation. *Computer Speech & Language*, 45:137–148.
- Thanh-Le Ha, Jan Niehues, and Alexander H. Waibel. 2016. Toward multilingual neural machine translation with universal encoder and decoder. *arXiv*, 1611.04798.
- Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, et al. 2018. Achieving human parity on automatic chinese to english news translation. *arXiv*, 1803.05567.
- Soren Hauberg, Oren Freifeld, Anders Boesen Lindbo Larsen, John W. Fisher, and Lars Kai Hansen. 2016. Dreaming more data: Class-dependent distributions over diffeomorphisms for learned data augmentation. In *AISTATS*.
- Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tieyan Liu, and Wei-Ying Ma. 2016a. Dual learning for machine translation. In *Conference on Advances in Neural Information Processing Systems (NIPS)*.
- Wei He, Zhongjun He, Hua Wu, and Haifeng Wang. 2016b. Improved neural machine translation with smt features. In *Conference of the Association for the Advancement of Artificial Intelligence (AAAI)*, pages 151–157.
- Kenneth Heafield, Ivan Pouzyrevsky, Jonathan H. Clark, and Philipp Koehn. 2013. Scalable Modified Kneser-Ney Language Model Estimation. In *Conference of the Association for Computational Linguistics (ACL)*.
- Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. Learning distributed representations of sentences from unlabelled data. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. Iterative back-translation for neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24.
- Kenji Imamura, Atsushi Fujita, and Eiichiro Sumita. 2018. Enhancement of encoder and attention using target monolingual corpora in neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 55–63.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda B. Viégas, Martin Wattenberg, Gregory S. 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 (TACL)*, 5:339–351.
- Lukasz Kaiser, Aidan N. Gomez, and François Chollet. 2017. Depthwise separable convolutions for neural machine translation. *CoRR*, abs/1706.03059.

- Nal Kalchbrenner, Lasse Espeholt, Karen Simonyan, Aäron van den Oord, Alex Graves, and Koray Kavukcuoglu. 2016. Neural machine translation in linear time. *CoRR*, abs/1610.10099.
- Alina Karakanta, Jon Dehdari, and Josef van Genabith. 2017. Neural machine translation for low-resource languages without parallel corpora. *Machine Translation*, pages 1–23.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*.
- Philipp Koehn. 2010. *Statistical machine translation*. Cambridge University Press.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *ACL Demo Session*.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Patrik Lambert, Holger Schwenk, Christophe Servan, and Sadaf Abdul-Rauf. 2011. Investigations on translation model adaptation using monolingual data. In *Workshop on Statistical Machine Translation (WMT)*.
- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018a. Unsupervised machine translation using monolingual corpora only. In *International Conference on Learning Representations (ICLR)*.
- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018b. Phrase-based & neural unsupervised machine translation. *arXiv*, 1803.05567.
- Marco Lui and Timothy Baldwin. 2012. langid.py: An off-the-shelf language identification tool. In *Proceedings of the ACL 2012 system demonstrations*, pages 25–30. Association for Computational Linguistics.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Xing Niu, Michael Denkowski, and Marine Carpuat. 2018. Bi-directional neural machine translation with synthetic parallel data. *arXiv preprint arXiv:1805.11213*.
- Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018a. Analyzing uncertainty in neural machine translation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 3956–3965.
- Myle Ott, Sergey Edunov, David Grangier, and Michael Auli. 2018b. Scaling neural machine translation.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Conference of the Association for Computational Linguistics (ACL)*.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations (ICLR)*.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton. 2017. Regularizing neural networks by penalizing confident output distributions. In *International Conference on Learning Representations (ICLR) Workshop*.
- Luis Perez and Jason Wang. 2017. The effectiveness of data augmentation in image classification using deep learning. *arXiv*, 1712.04621.
- Alberto Poncelas, Dimitar Sht. Shterionov, Andy Way, Gideon Maillette de Buy Wenniger, and Peyman Passban. 2018. Investigating backtranslation in neural machine translation. *arXiv*, 1804.06189.
- Matt Post. 2018. A call for clarity in reporting bleu scores. *arXiv*, 1804.08771.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Improving neural machine translation models with monolingual data. *Conference of the Association for Computational Linguistics (ACL)*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In *Conference of the Association for Computational Linguistics (ACL)*.
- Iulian Serban, Alessandro Sordani, Yoshua Bengio, Aaron C. Courville, and Joelle Pineau. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Conference of the Association for the Advancement of Artificial Intelligence (AAAI)*.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-attention with relative position representations. In *Proc. of NAACL*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Conference on Advances in Neural Information Processing Systems (NIPS)*.

- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2015. Rethinking the Inception Architecture for Computer Vision. *arXiv preprint arXiv:1512.00567*.
- 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 *Conference on Advances in Neural Information Processing Systems (NIPS)*.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, , and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *International Conference on Machine Learning (ICML)*.
- Yingce Xia, Tao Qin, Wei Chen, Jiang Bian, Nenghai Yu, and Tie-Yan Liu. 2017. Dual supervised learning. In *International Conference on Machine Learning (ICML)*.
- Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.