# How Far can 100 Samples Go? Unlocking Zero-Shot Translation with Tiny Multi-Parallel Data

**Di Wu    Shaomu Tan    Yan Meng    David Stap    Christof Monz**
Language Technology Lab
University of Amsterdam
{d.wu, s.tan, y.meng, d.stap, c.monz}@uva.nl

## Abstract

Zero-shot translation aims to translate between language pairs not seen during training in Multilingual Machine Translation (MMT) and is widely considered an open problem. A common, albeit resource-consuming, solution is to add as many related translation directions as possible to the training corpus. In this paper, we show that for an English-centric model, surprisingly large zero-shot improvements can be achieved by simply fine-tuning with a very small amount of multi-parallel data. For example, on the EC30 dataset, we obtain up to +21.7 ChrF++ non-English overall improvements (870 directions) by using only 100 multi-parallel samples while preserving English-centric translation quality. This performance exceeds M2M100 by an average of 5.9 ChrF++ in the involved non-English directions. When investigating the size effect of fine-tuning data on translation quality, we found that already a small, randomly sampled set of fine-tuning directions is sufficient to achieve comparable improvements. The resulting non-English performance is close to the complete translation upper bound. Even in a minimal setting—fine-tuning with only one single sample—the well-known off-target issue is almost completely resolved, explaining parts—but not all—of the observed improvements in translation quality.[1]

## 1 Introduction

The zero-shot capability shown by Multilingual Machine Translation (MMT) (Johnson et al., 2017) is of considerable significance, particularly in the context of translating between low-resource or distant language pairs. However, even for systems trained on large-scale data, the zero-shot performance is still far from sufficient (Tan and Monz, 2023), especially when scaling up the number of involved languages. Substantial efforts (Zhang et al., 2020; Pan et al., 2021; Gu and Feng, 2022; Mao

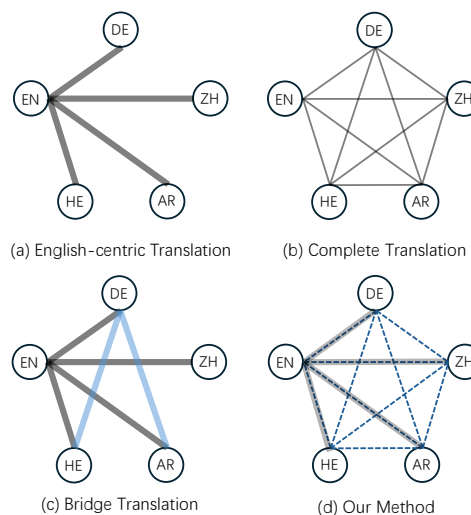[1]Models and codes are released: https://github.com/moore3930/MultiParallelFinetuning4MMT



Figure 1: (a) English-centric training data is normally readily available but can only cover a few real-world directions, while (b) complete translation (Freitag and Firat, 2020) aims to cover all but suffers from the small data scale. (c) Mining partial non-English data as the bridge languages shows promising zero-shot improvements but is also resource-consuming when scaling up. (d) We show that substantial overall improvements can be achieved by fine-tuning an English-centric model with tiny extra multi-parallel data, which is readily available, like NTREX (Federmann et al., 2022).

et al., 2023) have been dedicated to improving the zero-shot capabilities of models trained on readily available, predominantly English-centric corpora.

To fully cover translation directions, Freitag and Firat (2020) propose mining multi-parallel (multi-way aligned) samples to extend the training set from English-centric to a complete multilingual one as shown in Figure 1-(b). Non-English translation quality in this setting indeed increases substantially. However, such a setting is far away from real-world practice when scaling up. As shown in Freitag and Firat (2020), to solely extend the training set from bilingual aligned to all languages (6-way) involved in their case, the amount of available data drops from 123M to 10K, which is insufficient.

To reflect worldwide translation needs, Fan et al. (2021) build and open-source a training dataset covering 100 languages through industry-scale mining. In addition to English-centric data, supervised data for thousands of bridge language pairs is mined and included, organized based on language families. The MMT model trained on the resulting data, M2M100, exhibits clear improvements in many non-English directions. This work drives a simple but resource-consuming solution for real-world demand: mining as much training data as possible to bridge non-English language pairs, at least at the language family level.

In this paper, we take a step back and again look at the readily available English-centric model. We empirically show that the zero-shot ability of an English-centric model can be easily unlocked via fine-tuning with a tiny amount of multi-parallel data, which is much simpler and more efficient than the extensive training data mining and bridging done by earlier work. E.g., after fine-tuning the English-centric system trained on the EC30 dataset with 100 multi-parallel samples from NTREX (Federmann et al., 2022), we observed 21.7 ChrF++ zero-shot improvements across 870 language directions. This performance surpasses M2M100 (Fan et al., 2021) by 5.9 ChrF++ (See Table 16), despite M2M100 being trained on billions of sentence pairs spanning thousands of non-English directions.

Furthermore, we investigate the size effect of fine-tuning data: 1) Surprisingly, even when fine-tuning using a randomly sampled 10% of directions, the overall improvements are nearly the same as that of full-direction fine-tuning. 2) The improvements brought by very small fine-tuning datasets only slightly lag behind the upper bound (complete translation) while preserving English-centric capabilities, showing great practical potential. 3) Even with just one single multi-parallel sample for fine-tuning, the well-known off-target problem (Zhang et al., 2020; Yang et al., 2021; Sennrich et al., 2023), is easily addressed, reducing the off-target rate from 51.8% to 1.9%. However, not all improvements in translation quality can solely be attributed to lower off-target rates as we also see clear improvements in cases where translations are already in the correct target language.

Given the high efficiency and practicality, we encourage the community to consider fine-tuning with tiny readily available multi-parallel data, like NTREX (Federmann et al., 2022), as a strong baseline for zero-shot translation. Our findings challenge some prevailing assumptions about MMT systems based on the clear evidence proposed in this paper that: 1) the off-target issue is likely overestimated and can be easily handled, and 2) the potential of single-language-centric (like English) MMT models is substantial and often overlooked. We hope these insights will inspire new discussions and explorations within the community.

## 2 Related Work

The zero-shot translation capability of MMT is associated with multilingualism, following the hypothesis of universal representation or interlingua. Arivazhagan et al. (2019) view zero-shot translation as a domain adaptation problem (Ben-David et al., 2006) in MMT, and apply auxiliary losses to explicitly incentivize the model to learn and use domain- (language-) invariant representation. Liu et al. (2021) attribute the low quality of zero-shot MT to the positional correspondence to input tokens, which hinders modeling language-agnostic representation. Pan et al. (2021) and Stap et al. (2023) use a contrastive loss to close the representational gap between different languages. Some other approaches aim to harness the capabilities of pretrained multilingual models for zero-shot translation. Chen et al. (2022) employ multilingual pretrained encoders to extend bilingual translation to many-to-one translation. Recently, some work has focused on leveraging pre-trained large language models for multilingual translation (Zhang et al., 2023; Moslem et al., 2023). Despite the inclusion of the so-called "emergent abilities" (Wei et al., 2022) triggered by zero-shot prompting, we categorize these works as following a similar line.

This paper focuses on a data-centric approach for comprehensively improving zero-shot performance. We empirically show that a well-trained English-centric model can be easily boosted for overall zero-shot capability via fine-tuning with minimal data, even if only covering a small portion of translation directions (10%). This allows us to leverage multi-parallel data, which is hard to obtain in large quantities, resulting in a highly efficient and practical solution to overall zero-shot translation. We note that Maillard et al. (2023) also show that integrating small high-quality data (6K samples) into the training corpus can have a big impact on low-resource translation systems, especially when combined with back translation (Sen-

nrich et al., 2016). However, we contend that the reasons for the effectiveness differ: 1) as shown in Section 3.6, the substantial improvements persist when using data built from the training set, where the influence from domain or quality level is eliminated, and 2) our method can work with extremely minimal fine-tuning data (100 or even a single sample).

## 3 Experiments

In this paper, we propose a simple approach that largely improves overall zero-shot translation quality, i.e., fine-tuning an English-centric model in multiple directions with sentence pairs constructed from small, readily available multi-parallel datasets, like NTREX (Federmann et al., 2022). We refer to such a process as "fine-tuning with multi-parallel data" or "multi-parallel fine-tuning".

### 3.1 Fine-Tuning Data Construction

Given a multi-parallel dataset comprising $N$ distinct languages, each with $K$ samples, we can generate pairwise data in all $N \times (N-1)$ possible directions. Note that acquiring large quantities of multi-parallel data poses challenges due to many professional human translators being involved. However, horizontally expanding a readily available multi-parallel dataset to include one more language is straightforward. It simply requires annotating $K$ additional samples for the new language based on the current dataset, with $2 \times N$ new translation directions indirectly covered.

### 3.2 Datasets

**NTREX-128.** NTREX[2] (Federmann et al., 2022) is initially proposed as an evaluation dataset, expanding multilingual testing for translation from English into 128 target languages, which consists of 1997 samples per language and mainly focus on the News domain. Given the multi-parallel organization of NTREX data, we can easily build arbitrary pairwise data across 128 languages. In this paper, we leverage NTREX to create our fine-tuning datasets and conduct experiments to highlight the big impact of such a tiny amount of data.

**Europarl-8.** Europarl[3] (Koehn, 2005) consists of 20 English-centric language pairs from the proceedings of the European Parliament, with sizes ranging from 399K to 2M. A characteristic of Europarl is

that part of the samples are multi-way aligned. In this paper, we select the eight most resource-rich languages, i.e., EN, DA, DE, ES, FI, FR, IT, and NL, to mine a fully multi-parallel dataset named Europarl-8 via aligning multi-way sentences with the same English part. This results in about 1.2M fully multi-parallel data instances, where each sentence has 7 counterparts in other languages.

**EC30.** To ensure a more diverse and inclusive large-scale evaluation, we follow Tan and Monz (2023); Wu and Monz (2023) and use the EC30 dataset, which is built from WMT (Bojar et al., 2017) and OPUS (Tiedemann, 2012) corpora. EC30 comprises 61 million English-centric bilingual sentences for training, encompassing 30 non-English languages with diverse resource levels (High: 5M, Medium: 1M, Low: 100K). Each resource group includes languages from 5 families with multiple writing systems.

**Evaluation Benchmark.** For all of the experiments in this paper, we evaluate translations via the Flores-101 benchmark (Goyal et al., 2022). Flores comprises 3001 sentences sourced from English Wikipedia, which covers a variety of topics and domains and is translated into 101 languages by professional translators. We use *dev* and *devtest* as the validation and test dataset, consisting of 997 and 1012 samples, respectively. All results are evaluated on three widely used metrics, namely, ChrF++ (Popović, 2017), SacreBLEU (Post, 2018)[4], and COMET (Rei et al., 2020), to demonstrate the consistency of improvements across a broad spectrum of evaluation metrics. A more detailed description of the datasets is provided in Appendix A.1.

### 3.3 Experimental Setup

#### 3.3.1 Training Setting

For experiments on the EC30 dataset, we use Transformer-Big with 16 attention heads, 1,024 embedding dimensions, and 4,096 feedforward dimensions. For Europarl-8, we use a smaller backbone, as the training data is smaller, where a standard 6-layer encoder, 6-layer decoder transformer model is applied with 4 attention heads, 512 embedding dimensions, and 1,024 feedforward dimensions. In total, 447M and 64M training parameters are involved for the two models. More detailed training settings are provided in Appendix A.2.

| | Model | H-H | H-M | H-L | M-H | M-M | M-L | L-H | L-M | L-L | AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Two-Tag** | Baseline | 11.0 | 14.8 | 10.6 | 11.3 | 14.9 | 10.5 | 13.7 | 17.4 | 10.9 | 12.8 |
| | Boost-100 | 37.9 | 39.1 | 30.6 | 38.0 | 38.8 | 30.5 | 33.9 | 34.8 | 26.8 | 34.5 |
| | Boost-All | 38.6 | 39.9 | 32.0 | 38.7 | 39.5 | 31.8 | 34.9 | 35.9 | 28.4 | 35.5 |
| | Δ-100 | +26.9 | +24.3 | +20.0 | +26.7 | +23.9 | +20.0 | +20.2 | +17.4 | +15.9 | +21.7 |
| | Δ-All | +27.6 | +25.1 | +21.4 | +27.4 | +24.6 | +21.3 | +21.2 | +18.5 | +17.5 | +22.7 |
| **One-Tag** | Baseline | 28.0 | 30.4 | 20.0 | 27.9 | 29.5 | 19.7 | 24.0 | 26.0 | 16.4 | 24.7 |
| | Boost-100 | 36.6 | 37.4 | 29.0 | 36.5 | 36.9 | 28.8 | 32.1 | 32.7 | 24.8 | 32.8 |
| | Boost-All | 37.2 | 38.2 | 30.6 | 37.3 | 37.8 | 30.5 | 33.4 | 34.0 | 26.9 | 34.0 |
| | Δ-100 | +8.6 | +7.0 | +9.0 | +8.6 | +7.4 | +9.1 | +8.1 | +6.7 | +8.4 | +8.1 |
| | Δ-All | +9.2 | +7.8 | +10.6 | +9.4 | +8.3 | +10.8 | +9.4 | +8.0 | +10.5 | +9.3 |

Table 1: Zero-shot performance (ChrF++, 870 directions) trained on the EC30 dataset (61M English-centric sentence pairs), grouped by **H**igh-, **M**edium, and **L**ow-resource, respectively. Δ-100 and Δ-All mean the corresponding performance changes after fine-tuning compared to the baselines. Large improvements are observed in both One-Tag and Two-Tag settings after fine-tuning with 100 samples. Results in SacreBLEU and COMET are provided in Table 12 and Table 14, respectively. The comparison with M2M100 is provided in Table 16.

| | Model | High | | Medium | | Low | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | AVG |
| **Two-Tag** | Baseline | 52.5 | 57.5 | 53.9 | 56.5 | 42.5 | 49.8 | 49.6 | 54.6 | 52.1 |
| | Boost-100 | 51.9 | 56.6 | 53.6 | 57.0 | 42.6 | 50.2 | 49.4 | 54.6 | 52.0 |
| | Boost-All | 51.9 | 56.0 | 53.6 | 56.6 | 43.2 | 50.9 | 49.6 | 54.5 | 52.0 |
| | Δ-100 | -0.6 | -0.9 | -0.3 | +0.5 | +0.1 | +0.4 | -0.2 | 0.0 | -0.1 |
| | Δ-All | -0.6 | -1.5 | -0.3 | +0.1 | +0.7 | +1.1 | 0.0 | -0.1 | -0.1 |
| **One-Tag** | Baseline | 52.6 | 57.0 | 54.0 | 56.2 | 42.9 | 49.6 | 49.8 | 54.3 | 52.1 |
| | Boost-100 | 52.0 | 56.1 | 53.7 | 56.5 | 43.2 | 49.9 | 49.6 | 54.2 | 51.9 |
| | Boost-All | 51.7 | 55.7 | 53.5 | 56.2 | 43.5 | 50.3 | 49.6 | 54.1 | 51.8 |
| | Δ-100 | -0.6 | -0.9 | -0.3 | +0.3 | +0.3 | +0.3 | -0.2 | -0.1 | -0.2 |
| | Δ-All | -0.9 | -1.3 | -0.5 | 0.0 | +0.6 | +0.7 | -0.3 | -0.2 | -0.3 |

Table 2: English-centric performance (ChrF++, 60 directions) trained on the EC30 dataset (61M English-centric sentence pairs). EN-X and X-EN denote the average out-of- and into-English translation performance for each resource group, respectively. It's easy to see that the performance changes (Δ-100 and Δ-All) after fine-tuning are negligible. Results in SacreBLEU and COMET are provided in Table 13 and Table 15, respectively.

### 3.3.2 Fine-Tuning Setting

We use full-parameter fine-tuning and keep our setup as simple as possible to highlight generalizability. We reset all running statuses, including optimizer, lr scheduler, and data loaders. Also, the fine-tuning parameters are aligned with those in the training period, except for the experiments in Section 4, where we set batch accumulation to 1 as extremely small fine-tuning data is used.

Note that, because we only fine-tune with a small amount of data, the process is highly efficient: most fine-tuning experiments in this paper were completed within 1 GPU hour.

### 3.4 Large-Scale Experiments on EC30

In this section, we show how far a tiny amount of multi-parallel data can improve the zero-shot capability of an already well-trained large-scale English-centric MMT system. We conduct experiments on EC30, involving 30 English-centric and 870 zero-shot directions. We build fine-tuning data based on NTREX to cover all except Occitan-related directions—Occitan is not included in NTREX—as described in Section 3.2. It is noteworthy that MMT systems typically use two language tag strategies: 1) the one-tag strategy, i.e., adding the target language IDs to the encoder input, which is shown by Wu et al. (2021) to be more effective for zero-shot translation, or 2) the two-tag strategy, i.e., adding source and target language IDs to the encoder and decoder input, respectively, which is often applied to recent large-scale MMT systems (Fan et al., 2021; Pan et al., 2021; Costa-jussà et al., 2022). To show comprehensive results, we conducted experiments in both settings.

Table 1 shows the zero-shot performance across 9 resource groups. Boost-All refers to using all 1997 multi-parallel samples from NTREX to con-
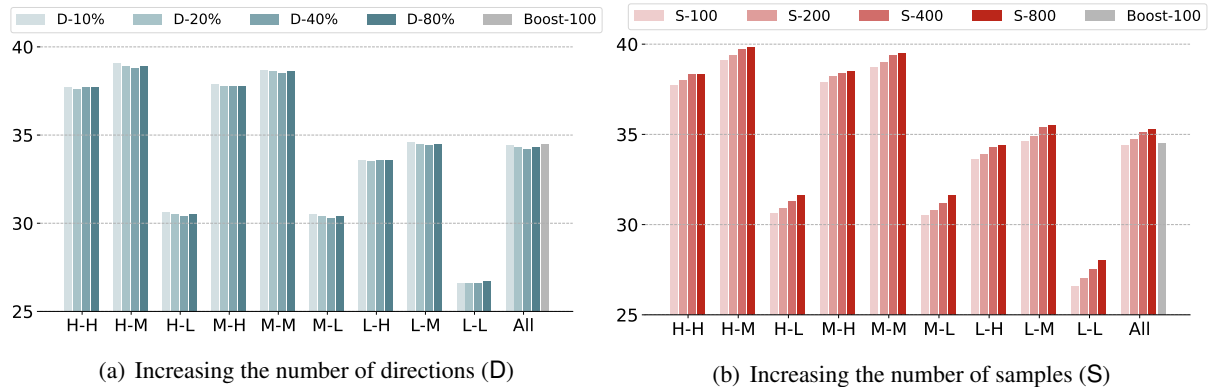
(a) Increasing the number of directions (D)    (b) Increasing the number of samples (S)

Figure 2: Zero-shot performance (ChrF++) on EC30 for each scaling step, grouped by **H**igh-, **M**edium, and **L**ow-resource, respectively. (a) When we randomly selected {10%, 20%, 40%, 80%} of directions and fine-tuned with 100 samples, the overall zero-shot performance stayed nearly unchanged. However, (b) when we fixed 10% of the directions and increased the fine-tuning samples from 100 to 800, we observed consistent improvements across all resource groups. Note that S-100 and D-10% are identical. Meanwhile, the directions or samples involved in the previous step are always subsets of those in the subsequent step.

struct pair-wise fine-tuning data for all directions (including English-centric ones), while Boost-100 means using only 100 randomly sampled samples, rather than 1997, to construct the fine-tuning data.

We find that 1) fine-tuning with tiny data leads to very strong overall improvements for both tagging strategies, with up to +9.3 and +22.7 average ChrF++ point gains, respectively. 2) The zero-shot capability of the two-tag baseline lags behind the one-tag baseline, in line with Wu et al. (2021). However, after fine-tuning with multi-parallel data, the overall performance in the two-tag setting consistently outperforms the one-tag setting for each group, yielding an average margin of +1.5 ChrF++ (35.5 v.s., 34.0). Consistent improvements also hold for other metrics, see Appendix A.3.

In Table 6, we show that 854 out of 870 zero-shot directions get strong boosts (more than 10.0 ChrF++). The resulting zero-shot performance even averagely surpasses the industrial system, M2M100, by 5.9 ChrF++ (See Table 16), despite M2M100 being trained on billions of sentence pairs spanning thousands of non-English directions.

In Table 2, we show the impact of fine-tuning on English-centric directions: The trade-off effect mainly occurs in the high-resource group. However, the influence on the medium- and low-resource groups is negligible or even positive, especially for the low-resource part, resulting in nearly unchanged overall English-centric performance. For instance, fine-tuning with 100 multi-parallel samples on the two-tag model yields +21.7 ChrF++ zero-shot gains, with negligible drops in averaged

English-centric performance (-0.1 ChrF++).

It is noteworthy that fine-tuning with just 100 samples achieves comparable improvements to using the entire NTREX dataset (+21.7 v.s. +22.7 in Table 1), even though the latter's size is 20 times larger. This diminishing effectiveness naturally leads us to ask (i) whether more fine-tuning data or more fine-tuning directions is important and (ii) how close can our method come to the upper-bound improvements.

We answer both questions in Section 3.5 and 3.6, respectively. If not specified, we employ the two-tag strategy in subsequent experiments because of its higher zero-shot and English-centric performance after fine-tuning.

## 3.5 More Data or More Directions?

In Section 3.4, we showed that fine-tuning an English-centric model with a small amount of bitext derived from NTREX (covering all directions) yields substantial zero-shot improvements. A natural assumption is that the improvement in each direction is triggered by the corresponding directional data. In this section, we investigate whether this is true, i.e., what will happen if we only cover a subset of translation directions during fine-tuning.

We conduct experiments on the same English-centric model trained on EC30, see Section 3.4, and control the scale of the fine-tuning data in the following two settings: (a) We conducted a random sampling of 100 multi-parallel NTREX sentences to construct pairwise data to cover all 870 direc-
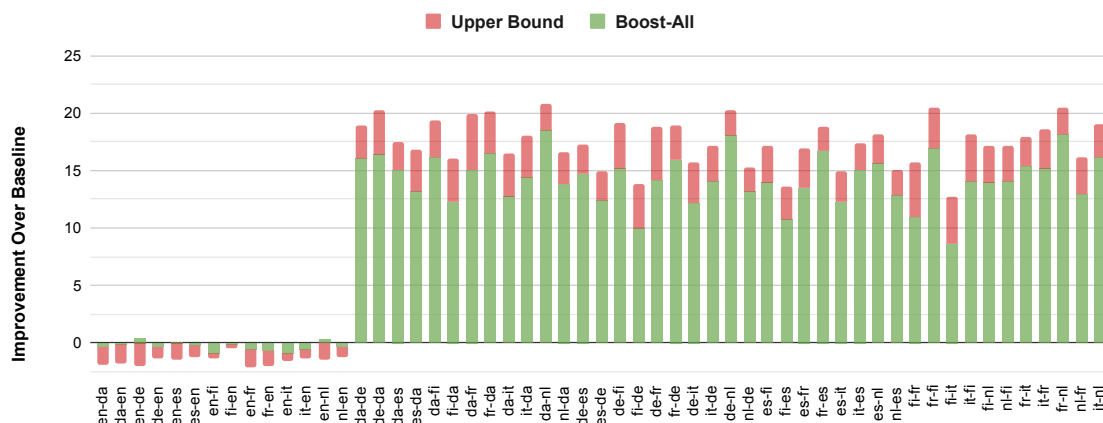
Figure 3: ChrF++ improvements of the upper bound and boosted models over the English-centric baseline on the Europarl-8 dataset. It is clear that the overall non-English capability of the boosted model is close to the upper bound (complete translation), meanwhile, it also holds the performance in English-centric directions.

tions[5]. Then, we randomly sampled {10%, 20%, 40%, 80%} directions for fine-tuning. (b) We fixed the 10% of directions as mentioned in (a) and conducted a random sampling of {100, 200, 400, 800} multi-parallel NTREX instances to construct the fine-tuning set for the corresponding directions.

Note that the bitext size in settings (a) and (b) for each scaling step is kept identical, e.g., to facilitate a fair comparison with the setting of fine-tuning with 100 multi-parallel samples in 80% directions, we also consider fine-tuning in 10% directions with 800 multi-parallel samples. Meanwhile, the directions or samples involved in the previous step are always subsets of those in the subsequent step.

In Figure 2, we show all of the corresponding fine-tuning results. Surprisingly, when fixing the size of the multi-parallel samples to 100 and then increasing the fine-tuning directions from 10% to 80%, no improvement is observed for any resource group (Figure 2-a). Fine-tuning in randomly sampled 10% directions using 100 samples achieves comparable overall results to fine-tuning in all directions (Boost-100). However, when we fix the directions to 10% and increase the multi-parallel sample size from 100 to 800, consistent improvements for all groups can be observed (Figure 2-b). This shows that the overall improvements are not sensitive to the number of directions, at least when the directions extend to a certain scale, like 10%.

In Appendix A.4.1, we further show that when we limited the fine-tuning direction set to fall in a specific family (Germanic), overall improvements

only slightly lag behind that of fully fine-tuning, showing surprising boosting effects.

### 3.6 How Close to the Upper Bound?

In this section, we show to what extent our fine-tuning method can approximate an upper bound. Here, we consider the performance of complete translation, i.e., training with fully multi-parallel data, as the "upper bound", since identical scales of non-English bitext cover all the directions that the English-centric counterparts cannot cover.

We conduct experiments on Europarl-8, see Section 3.2, where 8-way aligned data are available. Both English-centric and complete translation models are trained based on it. Note that we reuse the former's vocabulary for the latter to ensure a fair comparison. Also, we present the results after fine-tuning the English-centric model using full-direction pairwise data constructed from all of the 1,997 samples from NTREX.

In Figure 3, we show that 1) the upper-bound performance, i.e., that of the complete model, surpasses the baseline by a large margin, resulting in +17.4 average ChrF++ gains for non-English directions. 2) However, the boosted model's performance closely approaches the upper bound in all non-English directions. 3) For the 14 English-centric directions, both the upper bound and the boosted models exhibit degradation compared to the baseline, which reveals trade-off effects from English-centric to non-English directions. However, the boosted method only slightly degrades for a few English-centric directions (e.g., en-fi and en-it in Figure 3), whereas the upper bound model
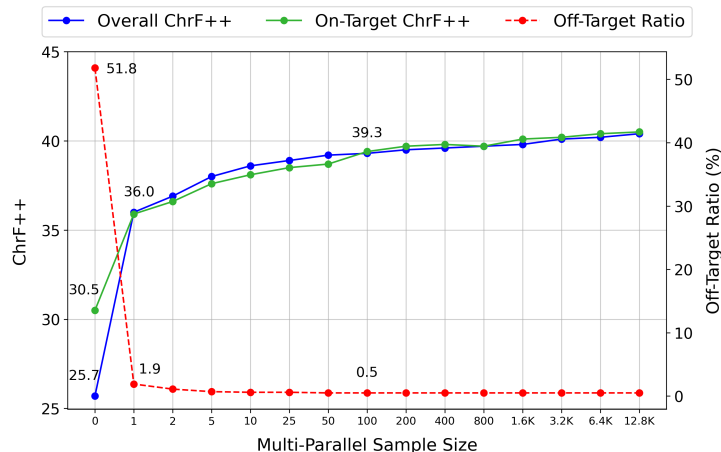
---

15097

Figure 4: Zero-shot performance and off-target ratio on Europarl-8 at each scaling step. Surprisingly, fine-tuning with one sample almost completely handles the off-target problem (from 51.8 % to 1.9%). The green solid line denotes the quality improvements of the translation samples that have no off-target issue, where the gains from one single sample are also large (+5.5 ChrF++). Mean and variance for each scaling step are provided in Table 17.

drops for most. Detailed scores, including those in other metrics, are provided in Appendix A.5. In short, the boosted model achieves strong non-English gains (+14.3 ChrF++, 42 directions) with a negligible cost in English-centric directions.

## 4 Analysis

### 4.1 Off-Target and Fine-Tuning Data Size

The off-target problem has been viewed as a primary cause that impairs the zero-shot capability (Zhang et al., 2020; Yang et al., 2021; Sennrich et al., 2023; Chen et al., 2023). In this section, we delve into the impact of fine-tuning data size on off-target ratios and final performance. Moreover, we disentangle the gains of the already on-target translations, showing the extent to which the enhancements are beyond alleviating the off-target issue. Note that since Europarl-8 is fully multi-parallel, we can readily build the corresponding full-direction fine-tuning data at different scales.

Here, we sample multiple sets of multi-parallel instances from the training set of Europarl-8, ranging from 1 to 12.8K with different seeds 3 times. The average results in ChrF++ after fine-tuning for each scaling step are provided in Figure 4. Also, we report the corresponding off-target ratio evaluated by fastText[6] (Joulin et al., 2017) following previous works (Yang et al., 2021; Costa-jussà et al., 2022).

In Figure 4, the blue solid line shows the overall zero-shot performance at each scaling step, where the starting point (fine-tuning with 0 samples) denotes the performance of the original English-

centric model. Notably, a high off-target ratio (51.8%) exists at this point. Surprisingly, even fine-tuning with just one multi-parallel sample, very strong overall zero-shot improvements can be obtained (from 25.7 to 36.0 ChrF++). Meanwhile, the off-target issue is almost completely resolved, dropping from 51.8% to 1.9%. Increasing from 1 to 100 samples, we can still observe clear zero-shot capability boosting (from 36.0 to 39.3 ChrF++), while the off-target change is marginal. Further scaling up fine-tuning data from the point of 100 samples shows nearly linear performance gains.

We further disentangle the evaluation set into on-target parts[7] for each language direction, where each source sample is already translated into the correct direction when evaluating the English-centric model. In Figure 4, the green solid line denotes the average performance on the on-target part. It is easy to see that the improvements brought from one single fine-tuning sample are still strong (+5.5 ChrF++), even after isolating the impact of off-target issues. As the number of fine-tuning samples increases, on-target improvements closely follow the trend of the overall improvements, further showing that the overall improvements are not only due to resolving the off-target issue.

Given the clear findings mentioned above, we question previous works (Zhang et al., 2020; Sennrich et al., 2023; Chen et al., 2023) that attribute off-target as a challenging problem and primary source for the low zero-shot performance.

---

[6]https://github.com/facebookresearch/fastText

[7]Note that the on-target sample size varies across directions, with an average of 488 samples per direction.

| Model | H-H | H-M | H-L | M-H | M-M | M-L | L-H | L-M | L-L | ZS-AVG | EN-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Baseline | 11.0 | 14.8 | 10.6 | 11.3 | 14.9 | 10.5 | 13.7 | 17.4 | 10.9 | 12.8 | 52.1 |
| (a) Multi-Parallel | 37.6 | 38.9 | 30.1 | 36.7 | 37.3 | 29.0 | 32.1 | 32.9 | 24.8 | 33.3 | 50.9 |
| (b) Multi-Directional | 37.6 | 39.1 | 30.3 | 37.1 | 38.0 | 29.7 | 33.0 | 34.1 | 26.0 | 33.9 | 51.5 |

Table 3: Decoupling multi-parallel and multi-directional fine-tuning on the EC30 dataset. ZS-AVG and EN-AVG denote the average results of the zero-shot and English-centric performance in ChrF++, respectively. In all groups, the performance in setting (a) closely trails but never surpasses that in setting (b), showing that the boosting effects do not depend on multi-way semantic equivalence.

## 4.2 Does Multi-Parallelism Matter?

We have shown remarkable boosting effects obtained from a tiny amount of multi-parallel data. But, does the data have to be *multi*-parallel? In this section, we explore whether utilizing multi-parallel data, instead of just pairwise data, for fine-tuning is vital for significant enhancements.

To this end, we design experiments to decouple the impacts of multi-parallel signals. Let's consider creating fine-tuning data in 5 languages from NTREX, involving 20 directions. We carefully control the resulting bitext distribution: Firstly, we randomly map the 1,997 multi-parallel samples into 10 buckets, ensuring an even allocation of approximately 100 samples per bucket. Then, we construct pairwise data in the following two ways:

(a) Multi-Parallel: We constructed pairwise samples in the 20 directions using the multi-parallel data in one randomly picked bucket.

(b) Multi-Directional: For each bucket, we construct fine-tuning samples for only one specific language pair (2 directions), also resulting in the 20 translation directions.

Note that the size of the bitext in settings (a) and (b) are identical. In (a), each sentence has semantically equivalent counterparts in all other 4 languages. However, in (b), each sentence has only one counterpart, resulting in simple pairwise data.

To cover different language families and resource levels, we choose DE, FR, RU, HE, and AR as these 5 languages. In Table 3, we show the results that fine-tune the EC30-based English-centric model with data in (a) multi-parallel and (b) multi-directional settings, respectively. Firstly, compared to the baseline model, clear improvements can be observed for zero-shot translation in both settings. Meanwhile, in all groups, the performance in setting (a) closely trails but never surpasses that in setting (b). It shows that the boosting effects do NOT depend on multi-way semantic equivalence,

| Setting | EN-X | X-EN | Zero-Shot | Off-Target (%) |
|---|---|---|---|---|
| Baseline | 49.8 | 51.3 | 25.7 | 51.8 |
| Numbers | 49.9 | 51.5 | 26.8 | 46.1 |
| Words | 48.3 | 50.7 | 35.8 | 3.6 |
| NTREX | 49.5 | 50.9 | 40.0 | 0.5 |

Table 4: The results (ChrF++) after fine-tuning. "Numbers", "Words", and "NTREX" denote different types of fine-tuning data (see Section 4.3).

implying simple multi-directional data is sufficient in case fully multi-parallel samples do not exist.

## 4.3 The Role of Semantic and Syntactic Information in Fine-Tuning Data

Considering that a small amount of fine-tuning data, e.g., 100 or even one single sample, can still substantially enhance overall zero-shot performance, a related question arises: To what extent do these improvements stem from the intrinsic information inherent in the data itself? In this section, we provide some insights into the role that the semantics and syntactic of fine-tuning data play in the surprising improvements for zero-shot translation.

We choose the English-centric model trained on Europarl-8 as our baseline (see Section 3.6), and fine-tune it on three datasets as follows:

**Number Pairs.** For each direction, we perform a uniform sampling of digits (ranging from 1 to 1000) multiple times, concatenating them to a certain length. Then, it is replicated for both the source and target sides, forming a number translation sample, as shown in Figure 6. Given no semantic other than numerical information is contained in this setting, we try to check whether the improvements stem from factors other than the data itself, e.g., the tags.

**Word Pairs.** We utilize bilingual dictionaries from MUSE [8] (Lample et al., 2018) to build word pairs for all of the directions. MUSE contains 110 English-centric bilingual dictionaries and all languages of Europarl-8 are included. We first select

the intersection of English words for the 7 involved English-centric dictionaries. Then, we extend them by mapping words paired with the same English words together[9]. E.g., given an EN-DE pair {*bike*, *Fahrrad*} and an EN-NL pair {*bike*, *fiets*}, we can build a new DE-NL word pair {*Fahrrad*, *fiets*}. Finally, we built 28 dictionaries (16,737 word pairs for each) covering all 56 directions.

**Sentence Pairs.** We use 100 randomly selected multi-parallel samples from NTREX to construct pairwise data covering all directions.

To ensure a fair comparison, we maintain similar surface information across the three datasets, such as aligning the number of tokens in the number-pair and word-pair datasets with the English portion in the sentence-pair dataset. Table 4 shows the corresponding fine-tuning results: 1) Fine-tuning with number pairs results in marginal improvements. Conversely, fine-tuning with word pairs leads to noticeable zero-shot improvements (+10.1 ChrF++). Simultaneously, the off-target ratio also decreases to an acceptable level. This means that semantic information, particularly at the lexical level, plays an important role here. 2) When using sentence-pair data (NTREX) to fine-tune, considerable further improvements compared to word-pair counterparts can be observed, showing that syntactic-level information also matters.

## 5 Conclusion

In this paper, we show that the zero-shot performance of an English-centric MMT model can be easily boosted by a tiny amount of multi-parallel data. On EC30, +21.7 ChrF++ average gains can be achieved by fine-tuning using 100 samples from NTREX, meanwhile preserving the English-centric performance, see Section 3.4. More surprisingly, we show that fine-tuning on a small portion (10%) of directions can achieve comparable improvements to full-direction fine-tuning, see Section 3.5, which are even close to the ideal but impractical upper-bound model, see Section 3.6.

In terms of using language tags, we show that fine-tuning can address the two-tag model's performance degradation in zero-shot directions (Wu et al., 2021). Moreover, the final performance substantially surpasses that of the one-tag model across multiple metrics, see Section 3.4.

We also question earlier findings (Zhang et al., 2020; Yang et al., 2021; Sennrich et al., 2023; Chen et al., 2023) that consider the off-target issue as a challenging problem for MMT. This paper shows that the off-target issue can be easily addressed by fine-tuning with even a single sample. Moreover, we also observe clear gains on the already on-target translations after fine-tuning with a few samples, implying that the off-target issue is not the primary source, but more like a symptom, of low zero-shot quality, see Section 4.1.

By decoupling the impacts of multi-parallel signals, we demonstrate that fine-tuning data does not have to be *multi*-parallel to achieve significant improvements, see Section 4.2. However, we still recommend using multi-parallel data, as it is already sufficient, albeit small-scale, and convenient for building fine-tuning data in each involved direction. Lastly, we shed some light on the impact of different types of signals from the fine-tuning data on the final performance, see Section 4.3.

Given the clear advantages of our proposed method, we encourage the community to consider fine-tuning as a strong baseline for zero-shot translation, especially in the two-tag setting. Furthermore, we suggest that the community reconsider the real bottlenecks of MMT systems in light of the evidence presented in this paper: 1) the off-target issue is overestimated and can be easily managed, and 2) the potential of single-language-centric (like English) MMT models, is substantial and often overlooked.

## Limitations

Multi-parallel data are normally built in a way that translates the same English data into multiple other languages by professional human translators. Hence, in the resulting non-English fine-tuning data, both the source and target side are translated instead of using the original text. This may exacerbate potential drawbacks in certain directions, such as translationese.

## Broader Impact

MMT systems have significant progress recently. However, potential challenges such as mistranslation or fairness problems still exist, e.g., the generation ability is not guaranteed to be fair across languages or demographic features, which may run the risk of reinforcing societal biases.

---

[9]Specifically, for each one-to-many mapping that exists, we randomly select a one-to-one mapping.

## Acknowledgement

## References

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Roee Aharoni, Melvin Johnson, and Wolfgang Macherey. 2019. The missing ingredient in zero-shot neural machine translation. *arXiv preprint arXiv:1903.07091*.

Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. 2006. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Proceedings of the Second Conference on Machine Translation*, pages 169–214, Copenhagen, Denmark. Association for Computational Linguistics.

Guanhua Chen, Shuming Ma, Yun Chen, Dongdong Zhang, Jia Pan, Wenping Wang, and Furu Wei. 2022. Towards making the most of cross-lingual transfer for zero-shot neural machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 142–157, Dublin, Ireland. Association for Computational Linguistics.

Liang Chen, Shuming Ma, Dongdong Zhang, Furu Wei, and Baobao Chang. 2023. On the off-target problem of zero-shot multilingual neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9542–9558, Toronto, Canada. Association for Computational Linguistics.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. *The Journal of Machine Learning Research*, 22(1):4839–4886.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Markus Freitag and Orhan Firat. 2020. Complete multilingual neural machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 550–560, Online. Association for Computational Linguistics.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Shuhao Gu and Yang Feng. 2022. Improving zero-shot multilingual translation with universal representations and cross-mapping. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6492–6504, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

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

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand.

Guillaume Lample, Alexis Conneau, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. Word translation without parallel data. In *International Conference on Learning Representations*.

Danni Liu, Jan Niehues, James Cross, Francisco Guzmán, and Xian Li. 2021. Improving zero-shot translation by disentangling positional information. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages

1259–1273, Online. Association for Computational Linguistics.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Zhuoyuan Mao, Raj Dabre, Qianying Liu, Haiyue Song, Chenhui Chu, and Sadao Kurohashi. 2023. Exploring the impact of layer normalization for zero-shot neural machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1300–1316, Toronto, Canada. Association for Computational Linguistics.

Yasmin Moslem, Rejwanul Haque, and Andy Way. 2023. Adaptive machine translation with large language models. *arXiv preprint arXiv:2301.13294*.

Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-many multilingual neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 244–258, Online. Association for Computational Linguistics.

Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.

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

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

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.

Rico Sennrich, Jannis Vamvas, and Alireza Mohammadshahi. 2023. Mitigating hallucinations and off-target machine translation with source-contrastive and language-contrastive decoding. *arXiv preprint arXiv:2309.07098*.

David Stap, Vlad Niculae, and Christof Monz. 2023. Viewing knowledge transfer in multilingual machine translation through a representational lens. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14973–14987, Singapore. Association for Computational Linguistics.

Shaomu Tan and Christof Monz. 2023. Towards a better understanding of variations in zero-shot neural machine translation performance. *arXiv preprint arXiv:2310.10385*.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.

Di Wu and Christof Monz. 2023. Beyond shared vocabulary: Increasing representational word similarities across languages for multilingual machine translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9749–9764, Singapore. Association for Computational Linguistics.

Liwei Wu, Shanbo Cheng, Mingxuan Wang, and Lei Li. 2021. Language tags matter for zero-shot neural machine translation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3001–3007, Online. Association for Computational Linguistics.

Yilin Yang, Akiko Eriguchi, Alexandre Muzio, Prasad Tadepalli, Stefan Lee, and Hany Hassan. 2021. Improving multilingual translation by representation and gradient regularization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7266–7279, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Biao Zhang, Philip Williams, Ivan Titov, and Rico Sennrich. 2020. Improving massively multilingual neural machine translation and zero-shot translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1628–1639, Online. Association for Computational Linguistics.

Xuan Zhang, Navid Rajabi, Kevin Duh, and Philipp Koehn. 2023. Machine translation with large language models: Prompting, few-shot learning, and fine-tuning with QLoRA. In *Proceedings of the Eighth Conference on Machine Translation*, pages 468–481, Singapore. Association for Computational Linguistics.

# A Appendix

## A.1 Detailed Dataset Description

**EC30.** In Table 5, we list the details of the EC40 dataset. We conducted experiments on EC30, a subset of EC40, where we excluded the data of 10 super low-resource languages, resulting in 30 English-centric language pairs with a total of 61M pairwise data. Each resource group consists of languages from 5 families with multiple writing systems.

## A.2 Training Setting

For all of the English-centric training, the learning rate is 5e-4 with 4,000 warmup steps and a *inverse sqrt* decay schedule. All dropout rates and label smoothing are set to 0.1. In the case of EC30 and Europarl-8, the batch size is set as 8,196 tokens, accumulating gradients 20 and 8 times, respectively. Also, data from different language pairs are sampled with a temperature of 5.0 and 2.0, respectively. The same temperature is applied to both BPE building and MMT training periods. We train all models with an early-stopping strategy[10] and evaluate by using the best checkpoint as selected based on the loss on the development set.

For fine-tuning, all parameters are kept the same as those in training, except for 1) we set batch accumulation as 1 in Section 4 as extremely small fine-tuning data is used, and 2) we set patience as 3 for quick experiments.

Note that we use 4 A6000 GPU cards for English-centric training with FP16 optimization, which means the actual batch size is also 4 times bigger. For fine-tuning, we use a single A6000 GPU card.

## A.3 Detailed Results on EC30

We report our detailed results in 970 directions (including English-centric and zero-shot directions) on EC30 datasets for both one-tag and two-tag models. The results are measured by 3 widely used metrics, i.e., ChrF++, SacreBLEU, and COMET.

Table 6, Table 7 and Table 8 show the specific performance of the two-tag model in each direction measured by ChrF++, SacreBLEU, and COMET, respectively. In each table, we report the corresponding performance of the baseline, boost-100,

and boost-all models. We also report the corresponding results in one-tag setting in Table 9, Table 10, and Table 11, respectively. The results grouped by resource level can be found in Table 1, 12, and 14 for ChrF++, SacreBLEU, and COMET, respectively.

We also report the influence fine-tuning brings on English-centric directions, which can be found in Table 2, 13, and 15 for ChrF++, SacreBLEU, and COMET, respectively.

## A.4 More Data or More Directions?

### A.4.1 Limited Fine-tuning Direction Set

To further investigate the surprising boosting effects that partial directional data brings, we limit the fine-tuning direction set to fall in a specific family and check the corresponding influence across language families. Here, we limit the fine-tuning set within *Germanic* (including English) and also use NTREX to build pairwise samples to cover all of the possible 42 translation directions.

Figure 5 summarizes the zero-shot performance across language groups. It is easy to see that even when limiting fine-tuning to a specific language family (Boost-Germanic), the overall performance remains comparable to full fine-tuning (Boost-All). More specifically, Boost-Germanic achieves a slight improvement over Boost-All in Germanic directions, meanwhile slightly lagging behind in all other groups, which is also intuitive. However, the gap between the two settings is still small. This finding further demonstrates the insensitivity of the directional data during fine-tuning. Detailed results, including those in other metrics, are provided in Table 18.

### A.4.2 Detailed Results: More Data or More Directions?

Table 18 shows the detailed results when fine-tuning with Germanic data and all of the NTREX data.

## A.5 Detailed Results: How Close to the Upper Bound?

In Table 19, 20 and 21, we show the detailed results for the baseline, boosted, and upper bound models on the Europarl-8 dataset in the metric of ChrF++, COMET, and SacreBLEU, respectively.

## A.6 Number Pairs

The synthetic number pairs and word pairs are illustrated in Figure 6.

---

[10]Patience is set to {10, 20}, i.e., training stops if performance on the validation set does not improve for the last {10, 20} checkpoints, with 1,000 steps between checkpoints.

| Resource | Germanic | | | Romance | | | Slavic | | | Indo-Aryan | | | Afo-Asiatic | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ISO | Language | Script | ISO | Language | Script | ISO | Language | Script | ISO | Language | Script | ISO | Language | Script |
| High (5M) | DE | German | Latin | FR | French | Latin | RU | Russian | Cyrillic | HI | Hindi | Devanagari | AR | Arabic | Arabic |
| | NL | Dutch | Latin | ES | Spain | Latin | CS | Czech | Latin | BN | Bengali | Bengali | HE | Hebrew | Hebrew |
| Med (1M) | SV | Swedish | Latin | IT | Italian | Latin | PL | Polish | Latin | KN | Kannada | Devanagari | MT | Maltese | Latin |
| | DA | Danish | Latin | PT | Portuguese | Latin | BG | Bulgarian | Cyrillic | MR | Marathi | Devanagari | HA | Hausa* | Latin |
| Low (100K) | AF | Afrikaans | Latin | RO | Romanian | Latin | UK | Ukrainian | Cyrillic | SD | Sindhi | Arabic | TI | Tigrinya | Ethiopic |
| | LB | Luxembourgish | Latin | OC | Occitan | Latin | SR | Serbian | Latin | GU | Gujarati | Devanagari | AM | Amharic | Ethiopic |
| eLow (50K) | NO | Norwegian | Latin | AST | Asturian | Latin | BE | Belarusian | Cyrillic | NE | Nepali | Devanagari | KAB | Kabyle* | Latin |
| | IC | Icelandic | Latin | CA | Catalan | Latin | BS | Bosnian | Latin | UR | Urdu | Arabic | So | Somali | Latin |

Table 5: Details of the EC40 dataset. Numbers in the table represent the number of sentences, e.g., 5M denotes exactly 5,000,000 sentences. Two exceptions are Hausa and Kabyle, where the size is 334K and 18K, respectively.



(a) Within Each Family  (b) Out of Germanic  (c) Into Germanic

Figure 5: Zero-shot performance (ChrF++) on EC30. Boost-All means fully fine-tuning, while Boost-Germanic means partially fine-tuning using Germanic languages. (a) shows the average performance evaluated within a specific language group, where both the source and target languages belong. (b) and (c) show the average performance in out-of-Germanic and into-Germanic directions, respectively. Detailed results are provided in Table 18.



Table 6: The ChrF++ performance on EC30 for the baseline, boost-100, and boost-all models in the two-tag fashion, respectively. 970 directions of results are shown, including both English-centric and zero-shot ones. Overall, in all 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 845 out of 870 zero-shot directions, the gain exceeds 10.0.



Table 7: The SacreBLEU performance on EC30 for the baseline, boost-100, and boost-all models in the two-tag fashion, respectively. 970 directions of results are shown, including English-centric and zero-shot ones. Overall, in all 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 702 out of 870 zero-shot directions, the gain exceeds 5.0.

Table 8: The COMET performance on EC30 for the baseline, boost-100, and boost-all models in the two-tag fashion, respectively. 970 directions of results are shown, including English-centric and zero-shot ones. Overall, in 869 out of 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 636 out of 870 zero-shot directions, the gain exceeds 10.0.

Table 9: The ChrF++ performance on EC30 for the baseline, boost-100, and boost-all models in the one-tag fashion, respectively. 970 directions of results are shown, including English-centric and zero-shot ones. Overall, in all 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 564 out of 870 zero-shot directions, the gain exceeds 5.0.

Table 10: The SacreBLEU performance on EC30 for the baseline, boost-100, and boost-all models in the one-tag fashion, respectively. 970 directions of results are shown, including English-centric and zero-shot ones. Overall, in 869 out of 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 423 out of 870 zero-shot directions, the gain exceeds 3.0.

Table 11: The COMET performance on EC30 for the baseline, boost-100, and boost-all models in the one-tag fashion, respectively. 970 directions of results are shown, including English-centric and zero-shot ones. Overall, in 866 out of 870 zero-shot directions, the boost-100 model achieves better performance compared to the baseline. Moreover, in 501 out of 870 zero-shot directions, the gain exceeds 5.0.

|  | Model | H-H | H-M | H-L | M-H | M-M | M-L | L-H | L-M | L-L | AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Two-Tag | Baseline | 1.5 | 2.2 | 1.3 | 1.8 | 2.4 | 1.3 | 2.6 | 3.1 | 1.3 | 2.0 |
|  | Boost-100 | 13.4 | 14.2 | 9.0 | 13.8 | 14.2 | 9.2 | 10.8 | 11.2 | 6.8 | 11.4 |
|  | Boost-All | 14.0 | 14.9 | 10.0 | 14.4 | 14.9 | 10.1 | 11.6 | 12.1 | 7.8 | 12.2 |
|  | Δ-100 | +11.9 | +12.0 | +7.7 | +12.0 | +11.8 | +7.9 | +8.2 | +8.1 | +5.5 | +9.4 |
|  | Δ-All | +12.5 | +12.7 | +8.7 | +12.6 | +12.5 | +8.8 | +9.0 | +9.0 | +6.5 | +10.2 |
| One-Tag | Baseline | 8.4 | 9.0 | 4.6 | 8.8 | 8.9 | 4.7 | 6.4 | 6.7 | 3.1 | 6.7 |
|  | Boost-100 | 12.4 | 12.8 | 8.0 | 12.7 | 12.6 | 8.1 | 9.5 | 9.7 | 5.7 | 10.2 |
|  | Boost-All | 12.9 | 13.4 | 9.0 | 13.2 | 13.3 | 9.1 | 10.4 | 10.6 | 6.8 | 10.9 |
|  | Δ-100 | +4.0 | +3.8 | +3.4 | +3.9 | +3.7 | +3.4 | +3.1 | +3.0 | +2.6 | +3.5 |
|  | Δ-All | +4.5 | +4.4 | +4.4 | +4.4 | +4.4 | +4.4 | +4.0 | +3.9 | +3.7 | +4.2 |

Table 12: Zero-shot performance (SacreBLEU) on the EC30 dataset (61M sentence pairs) with two language tag strategies, grouped by **H**igh-, **M**edium, and **L**ow-resource, respectively. Δ-100 and Δ-All mean the corresponding performance changes compared with the baselines.

|  | Model | High | | Medium | | Low | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  |  | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | AVG |
| Two-Tag | Baseline | 28.0 | 31.4 | 29.5 | 31.4 | 18.7 | 25.8 | 25.4 | 29.5 | 27.5 |
|  | Boost-100 | 27.3 | 29.9 | 29.0 | 31.1 | 18.8 | 25.0 | 25.0 | 28.7 | 26.9 |
|  | Boost-All | 27.1 | 29.7 | 28.7 | 31.1 | 19.3 | 25.8 | 25.0 | 28.9 | 27.0 |
|  | Δ-100 | -0.7 | -1.5 | -0.5 | -0.3 | +0.1 | -0.8 | -0.4 | -0.8 | -0.6 |
|  | Δ-All | -0.9 | -1.7 | -0.8 | -0.3 | +0.6 | 0.0 | -0.4 | -0.6 | -0.5 |
| One-Tag | Baseline | 28.4 | 30.9 | 29.8 | 30.9 | 19.5 | 25.3 | 25.9 | 29.0 | 27.5 |
|  | Boost-100 | 27.6 | 30.5 | 29.6 | 31.4 | 19.4 | 25.5 | 25.5 | 29.1 | 27.3 |
|  | Boost-All | 27.0 | 30.0 | 28.7 | 31.2 | 19.6 | 25.5 | 25.1 | 28.9 | 27.0 |
|  | Δ-100 | -0.8 | -0.4 | -0.2 | +0.5 | -0.1 | +0.2 | -0.4 | +0.1 | -0.2 |
|  | Δ-All | -1.4 | -0.9 | -1.1 | +0.3 | +0.1 | +0.2 | -0.8 | -0.1 | -0.5 |

Table 13: English-centric performance (SacreBLEU) on the EC30 dataset (61M sentence pairs). EN-X and X-EN denote the average out-of- and into-English translation performance of each resource group, respectively.

|  | Model | H-H | H-M | H-L | M-H | M-M | M-L | L-H | L-M | L-L | AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Two-Tag | Baseline | 51.0 | 52.4 | 47.3 | 48.6 | 49.9 | 45.0 | 46.4 | 49.3 | 43.6 | 48.2 |
|  | Boost-100 | 70.6 | 71.2 | 64.0 | 69.1 | 69.7 | 63.1 | 63.1 | 64.8 | 58.7 | 66.0 |
|  | Boost-All | 71.3 | 71.8 | 65.7 | 69.9 | 70.4 | 64.8 | 64.4 | 66.0 | 60.9 | 67.3 |
|  | Δ-100 | +19.6 | +18.8 | +16.7 | +20.5 | +19.8 | +18.1 | +16.7 | +15.5 | +15.1 | +17.8 |
|  | Δ-All | +20.3 | +19.4 | +18.4 | +21.3 | +20.5 | +19.8 | +18.0 | +16.7 | +17.3 | +19.1 |
| One-Tag | Baseline | 61.2 | 62.6 | 54.4 | 58.8 | 60.2 | 52.2 | 51.9 | 54.9 | 47.8 | 56.0 |
|  | Boost-100 | 68.1 | 68.8 | 61.7 | 66.5 | 67.1 | 60.5 | 60.2 | 62.1 | 55.8 | 63.4 |
|  | Boost-All | 68.9 | 69.5 | 63.6 | 67.6 | 68.0 | 62.7 | 62.1 | 63.7 | 58.7 | 65.0 |
|  | Δ-100 | +6.9 | +6.2 | +7.3 | +7.7 | +6.9 | +8.3 | +8.3 | +7.2 | +8.0 | +7.4 |
|  | Δ-All | +7.7 | +6.9 | +9.2 | +8.8 | +7.8 | +10.5 | +10.2 | +8.8 | +10.9 | +9.0 |

Table 14: Zero-shot performance (COMET) on the EC30 dataset (61M sentence pairs) with two language tag strategies, grouped by **H**igh-, **M**edium, and **L**ow-resource, respectively. Δ-100 and Δ-All mean the corresponding performance changes compared with the baselines.

| | Model | High | | Medium | | Low | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | EN-X | X-EN | AVG |
| Two-Tag | Baseline | 82.3 | 83.8 | 81.0 | 80.0 | 73.7 | 73.3 | 79.0 | 79.0 | 79.0 |
| | Boost-100 | 82.0 | 83.4 | 81.1 | 81.3 | 73.7 | 74.4 | 78.9 | 79.7 | 79.3 |
| | Boost-All | 82.1 | 83.0 | 81.0 | 81.1 | 74.3 | 75.2 | 79.1 | 79.8 | 79.5 |
| | $\Delta$-100 | -0.3 | -0.4 | +0.1 | +1.3 | 0.0 | +1.1 | -0.1 | +0.7 | +0.3 |
| | $\Delta$-All | -0.2 | -0.8 | 0.0 | +1.1 | +0.6 | +1.9 | +0.1 | +0.8 | +0.5 |
| One-Tag | Baseline | 82.0 | 83.3 | 80.6 | 79.6 | 73.2 | 72.6 | 78.6 | 78.5 | 78.5 |
| | Boost-100 | 81.8 | 83.0 | 80.8 | 80.9 | 73.7 | 74.2 | 78.8 | 79.4 | 79.1 |
| | Boost-All | 81.6 | 82.7 | 80.6 | 80.7 | 73.9 | 74.7 | 78.7 | 79.4 | 79.0 |
| | $\Delta$-100 | -0.2 | -0.3 | +0.2 | +1.3 | +0.5 | +1.6 | +0.2 | +0.9 | +0.6 |
| | $\Delta$-All | -0.4 | -0.6 | 0.0 | +1.1 | +0.7 | +2.1 | +0.1 | +0.9 | +0.5 |

Table 15: English-centric performance (COMET) on the EC30 dataset (61M sentence pairs). EN-X and X-EN denote the average out-of- and into-English translation performance of each resource group, respectively.

| | Model | H-H | H-M | H-L | M-H | M-M | M-L | L-H | L-M | L-L | AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Two-Tag | M2M100 | 40.0 | 35.6 | 24.8 | 35.0 | 30.3 | 21.8 | 33.1 | 29.9 | 20.3 | 30.1 |
| | Baseline | 11.0 | 14.8 | 10.6 | 11.3 | 14.9 | 10.5 | 13.7 | 17.4 | 10.9 | 12.8 |
| | Boost-100 | 37.9 | 38.9 | 32.2 | 37.6 | 38.0 | 31.7 | 34.5 | 35.2 | 28.6 | 35.0 |
| | Boost-All | 38.6 | 39.6 | 33.5 | 38.3 | 38.6 | 33.0 | 35.6 | 36.3 | 30.3 | 36.0 |
| | $\Delta_{M2M}$-100 | -2.1 | +3.3 | +7.4 | +2.6 | +7.7 | +9.9 | +1.4 | +5.3 | +8.3 | +4.9 |
| | $\Delta_{M2M}$-All | -1.4 | +4.0 | +8.7 | +3.3 | +8.3 | +11.2 | +2.5 | +6.4 | +10.0 | +5.9 |

Table 16: The performance (ChrF++) of our methods compared to M2M100 (418M parameters) in the involved directions of EC30. Note that some scores slightly differ from those presented in Table 1. This disparity arises because M2M100 lacks support for the medium-resource language *mt* and the low-resource language *ti*. To ensure a fair comparison, we evaluate in 756 zero-shot directions, where these two languages are excluded. $\Delta_{M2M}$-100 and $\Delta_{M2M}$-All denote the performance gap between M2M100 and our two boosted models (Boost-100 and Boost-All).

| Data Size | 1 | 2 | 5 | 10 | 25 | 50 | 100 |
|---|---|---|---|---|---|---|---|
| AVG (ChrF++) | 36.0 | 36.9 | 38.0 | 38.6 | 38.9 | 39.2 | 39.3 |
| STDEV | 0.5 | 0.5 | 0.2 | 0.2 | 0.2 | 0.1 | 0.1 |

Table 17: We report the mean and variance of the overall zero-shot performance for each step in Figure 4 after running with three random seeds. Across sample sizes ranging from 1 to 100, we note a consistently small standard deviation. Furthermore, with increasing sample sizes, the standard deviation diminishes further, indicating a stable performance across observations.

| Metric | Model | Within Family | | | | | Out of Germanic | | | | Into Germanic | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Germanic | Romance | Slavic | Aryan | Asiatic | Romance | Slavic | Aryan | Asiatic | Romance | Slavic | Aryan | Asiatic |
| ChrF++ | Baseline | 20.8 | 23.5 | 12.1 | 8.0 | 7.2 | 21.0 | 10.5 | 5.8 | 7.3 | 20.9 | 20.0 | 16.8 | 17.8 |
| | Boost-All | 44.6 | **45.6** | **39.2** | **34.2** | **27.4** | **42.2** | **38.0** | **34.1** | **30.9** | **42.6** | **40.6** | **34.8** | **36.4** |
| | Boost-Germanic | **45.0** | 44.5 | 38.5 | 30.3 | 25.0 | 41.3 | 37.4 | 32.4 | 29.4 | 42.4 | 40.5 | 33.2 | 35.9 |
| Comet | Baseline | 54.3 | 58.5 | 50.4 | 45.9 | 41.6 | 56.7 | 47.6 | 39.5 | 41.3 | 54.9 | 55.4 | 52.8 | 50.0 |
| | Boost-All | 68.2 | **71.6** | **74.0** | **71.2** | **66.2** | **66.3** | **69.2** | **69.5** | **68.2** | **67.3** | **67.3** | **63.4** | **61.1** |
| | Boost-Germanic | **68.5** | 70.8 | 72.8 | 68.0 | 63.3 | 65.8 | 68.3 | 68.7 | 66.8 | 66.9 | 66.9 | 61.8 | 60.1 |
| SacreBLEU | Baseline | 2.8 | 3.7 | 2.0 | 3.7 | 0.8 | 2.7 | 1.6 | 2.5 | 1.1 | 2.6 | 2.2 | 1.6 | 1.7 |
| | Boost-All | 17.5 | **18.7** | **13.5** | **15.9** | **6.5** | **16.4** | **12.8** | **15.3** | **8.6** | **15.1** | **13.8** | **9.7** | **11.3** |
| | Boost-Germanic | **17.8** | 17.4 | 12.7 | 14.3 | 5.3 | 15.3 | 12.2 | 14.1 | 7.6 | 14.8 | 13.6 | 9.1 | 10.9 |

Table 18: Zero-shot performance on the EC30 dataset. Boost-All means fine-tuning using all of the NTREX data, while Boost-Germanic means partially fine-tuning using Germanic languages. The results are in three groups: 1) "Within Family" shows the performance within a specific language group, where both the source and target languages belong. 2) "Out of Germanic" shows the average performance that is translated out of Germanic languages, e.g., from Germanic to Romance. and 3) "Into Germanic" shows the average performance that is translated into Germanic languages, e.g., from Romance to Germanic.

| ChrF++ | en | da | de | es | fi | fr | it | nl |
|---|---|---|---|---|---|---|---|---|
| en | - | 55.5/55.1/53.6 | 50.6/51.0/49.0 | 46.2/46.3/44.8 | 44.2/43.3/42.8 | 57.4/56.8/55.3 | 47.9/47.0/46.3 | 46.5/46.8/45.3 |
| da | 57.2/57.1/55.4 | - | 28.3/44.3/47.2 | 24.8/39.9/42.3 | 22.4/38.5/41.8 | 30.5/45.6/50.4 | 26.8/39.5/43.3 | 23.4/41.9/44.2 |
| de | 53.4/53.1/52.0 | 28.5/44.9/48.7 | - | 24.6/39.3/41.9 | 21.4/36.6/40.5 | 30.2/44.4/49.0 | 26.9/39.1/42.6 | 23.4/41.4/43.7 |
| es | 48.0/47.8/46.8 | 26.5/39.7/43.3 | 25.9/38.3/40.8 | - | 20.9/34.8/38.1 | 30.5/44.0/47.4 | 28.3/40.6/43.2 | 22.6/38.2/40.7 |
| fi | 45.9/45.8/45.4 | 26.7/39.0/42.7 | 26.8/36.8/40.6 | 25.2/35.9/38.8 | - | 28.4/39.4/44.1 | 26.6/35.2/39.3 | 22.5/36.4/39.6 |
| fr | 56.0/55.3/54.0 | 27.3/43.8/47.4 | 26.8/42.7/45.7 | 25.7/42.4/44.5 | 20.7/37.6/41.2 | - | 28.4/43.8/46.3 | 23.0/41.2/43.5 |
| it | 50.2/49.6/48.9 | 26.4/40.8/44.4 | 25.8/39.8/43.0 | 26.5/41.6/43.9 | 21.1/35.1/39.3 | 31.2/46.4/49.8 | - | 22.9/39.0/41.9 |
| nl | 48.2/47.9/47.0 | 27.0/40.8/43.6 | 26.7/39.9/42.0 | 24.5/37.3/39.5 | 20.7/34.8/37.8 | 28.4/41.3/44.5 | 26.1/36.8/39.5 | - |

Table 19: The detailed ChrF++ results for the baseline, boosted, and upper bound models on the Europarl-8 dataset are presented, encompassing 14 English-centric directions and 42 zero-shot directions.

| COMET | en | da | de | es | fi | fr | it | nl |
|---|---|---|---|---|---|---|---|---|
| en | - | 79.5/78.8/76.9 | 72.8/72.5/70.4 | 75.1/74.9/71.6 | 77.5/76.6/75.7 | 75.3/73.8/72.0 | 76.1/74.9/72.8 | 75.4/75.0/73.0 |
| da | 78.6/78.7/77.6 | - | 49.7/64.8/71.8 | 55.3/65.3/71.4 | 50.4/68.5/75.6 | 52.7/61.6/70.3 | 53.1/64.6/72.4 | 52.4/68.1/73.5 |
| de | 76.6/76.6/75.8 | 56.1/71.6/76.1 | - | 54.2/63.6/69.7 | 49.0/66.2/73.9 | 51.4/61.7/68.6 | 52.6/64.2/71.4 | 51.0/68.5/73.7 |
| es | 75.5/75.5/74.6 | 54.8/66.6/73.7 | 46.9/60.7/67.3 | - | 48.2/66.5/73.4 | 53.4/65.5/72.1 | 54.5/69.6/75.8 | 50.3/64.6/70.7 |
| fi | 75.2/74.3/75.0 | 55.6/66.5/73.6 | 47.9/59.3/67.4 | 54.5/62.5/69.2 | - | 51.7/59.5/68.6 | 52.7/63.0/69.9 | 51.7/63.1/69.9 |
| fr | 79.1/79.1/77.6 | 57.4/69.6/75.5 | 48.6/63.7/70.6 | 57.9/70.9/75.2 | 50.3/68.9/75.6 | - | 56.3/72.5/77.2 | 52.7/68.4/73.2 |
| it | 76.6/76.6/75.8 | 55.8/67.8/74.2 | 47.4/61.3/69.2 | 56.6/70.4/74.8 | 48.8/66.6/74.3 | 54.8/67.3/73.4 | - | 50.8/65.6/71.8 |
| nl | 75.7/75.6/75.0 | 55.6/69.2/75.0 | 48.6/64.1/70.2 | 54.0/64.2/69.5 | 48.3/66.0/72.9 | 51.4/61.0/68.4 | 51.7/64.1/70.6 | - |

Table 20: The detailed COMET results for the baseline, boosted, and upper bound models on the Europarl-8 dataset are presented, encompassing 14 English-centric directions and 42 zero-shot directions.

| SacreBLEU | en | da | de | es | fi | fr | it | nl |
|---|---|---|---|---|---|---|---|---|
| en | - | 30.4/30.4/28.1 | 23.0/23.3/21.6 | 18.9/19.2/17.9 | 13.7/12.8/12.9 | 32.6/30.8/30.6 | 19.9/18.9/18.7 | 17.6/18.0/16.9 |
| da | 31.1/31.0/28.3 | - | 6.4/16.4/19.5 | 4.2/13.7/15.9 | 3.0/9.7/12.0 | 7.9/18.4/24.3 | 4.9/12.5/15.8 | 3.3/13.4/15.1 |
| de | 26.5/26.5/24.3 | 7.0/18.9/22.8 | - | 4.3/13.0/15.7 | 2.8/8.6/11.1 | 7.6/17.6/22.9 | 5.4/12.5/15.1 | 3.5/13.3/15.3 |
| es | 18.4/18.5/17.4 | 5.0/13.0/15.9 | 4.2/10.2/12.1 | - | 2.0/6.5/8.2 | 6.9/15.1/20.3 | 5.5/11.8/15.1 | 2.7/9.5/11.4 |
| fi | 18.5/18.4/17.9 | 5.7/13.7/16.8 | 4.8/10.1/13.1 | 4.5/10.5/12.7 | - | 6.1/13.5/18.0 | 4.5/9.2/12.7 | 2.8/9.0/11.3 |
| fr | 28.8/28.4/26.2 | 5.8/17.3/20.6 | 5.4/14.2/17.5 | 4.6/15.8/17.8 | 2.4/9.3/11.5 | - | 6.1/16.0/18.6 | 3.2/12.5/14.5 |
| it | 20.7/20.6/19.4 | 5.1/13.8/16.7 | 4.5/11.4/14.3 | 5.0/14.2/17.0 | 2.3/7.1/9.5 | 7.9/18.4/22.8 | - | 2.9/10.0/12.1 |
| nl | 20.7/20.5/19.2 | 5.6/14.6/17.0 | 5.1/12.2/13.7 | 4.2/11.4/13.3 | 2.4/7.4/8.9 | 6.6/14.7/18.0 | 4.3/10.3/12.1 | - |

Table 21: The detailed SacreBLEU results for the baseline, boosted, and upper bound models on the Europarl-8 dataset are presented, encompassing 14 English-centric directions and 42 zero-shot directions.

Figure 6: Illustration of number pairs.

| Source | Target |
|---|---|
| 961 271 4 137 146 37 124 498 28 73 | 961 271 4 137 146 37 124 498 28 73 |
| 323 135 42 121 324 499 17 11 37 33 | 323 135 42 121 324 499 17 11 37 33 |
| 713 93 331 225 51 7 375 223 12 192 | 713 93 331 225 51 7 375 223 12 192 |
| ... | ... |