# UniLoRA: Unified LoRA Fine-tuning for Any-to-Any Machine Translation with Low-resource Corpus

## Anonymous ACL submission

## Abstract

It is a great challenge for any-to-any machine translation in achieving robust performance across diverse language pairs, primarily due to the scarcity of non-English parallel corpora. Existing approaches often rely on massive multilingual datasets or cascaded translation pipelines, which introduce inefficiencies, computational costs, and error propagation. Models trained on limited language pairs struggle to generalize to unseen language directions, hindering practical deployment in real-world multilingual scenarios. To address these limitations, we propose UniLoRA, a novel instruction fine-tuning framework for Large Language Models (LLMs) that enables efficient any-to-any translation with minimal reliance on limited multilingual parallel data. Our approach leverages English-centric parallel corpora alongside limited multilingual translation examples to align cross-lingual representations, effectively bridging language gaps without requiring exhaustive language-specific supervision. UniLoRA employs parameter-efficient Low-Rank Adaptation (LoRA) modules alongside Mixture-of-Experts (MoE) framework to enable dynamic adaptation to arbitrary translation directions. Experiments demonstrate that our approach achieves competitive performance on diverse translation directions. This work provides a resource-efficient paradigm for democratizing high-quality any-to-any translation capabilities across linguistically diverse environments. Our code is available at: `https://anonymous.4open.science/r/UniL-1BD1/`.

## 1 Introduction

Recent large language models have demonstrated exceptional performance across numerous NLP tasks while exhibiting robust multilingual capabilities (Grattafiori et al., 2024; Üstün et al., 2024). In particular, LLMs-based Multilingual Neural Machine Translation (MNMT) systems have attained broad language coverage while maintaining high translation quality (Zeng et al., 2024). These systems, however, exhibit performance disparity across translation directions, especially in non-English-involved directions (Zhu et al., 2024; Xu et al., 2025). This disparity largely arises from the scarcity of **non-English Parallel Corpora**, which constitute only a minor fraction of publicly available datasets, making it difficult to develop robust any-to-any MNMT models (Arivazhagan et al., 2019; Schwenk et al., 2021; Kreutzer et al., 2022).

To address this issue, existing approaches mainly adopt two divergent paradigms. The first involves training multilingual models on massive aggregated datasets with synthetic data (Chen et al., 2017; Fan et al., 2021; NLLB Team et al., 2022), but these approaches suffer from the inferior quality of the synthetic parallel corpus. The second strategy utilizes cascaded pivot-based pipelines. For example, the source language is translated into the pivot language (e.g., English), and then into the target language as seen in Figure 1(a). Though reducing the dependency of the non-English parallel corpus, this pivot-based paradigm still suffers from the uncertainty of the pivot language selection and the error propagation issue (Liu et al., 2018; Zhang et al., 2020).

The key issue of MNMT is to establish a unified any-to-any translation with as little multilingual parallel corpus as possible. With the general representation ability of LLMs, most researchers currently turn to study the any-representation-any MNMT as seen in Figure 1(b). Based on LLMs, researchers can fine-tune a Low-Rank Adaptor (LoRA) (Hu et al., 2022) for the specialized translation direction as seen in Figure 2(a) (Chen et al., 2024; Gao et al., 2024). The straightforward extension for any-to-any translation is to fine-tune one LoRA network for any translation direction. We argue that this extension will introduce **a non-uniform fine-tune representation** and suffer from
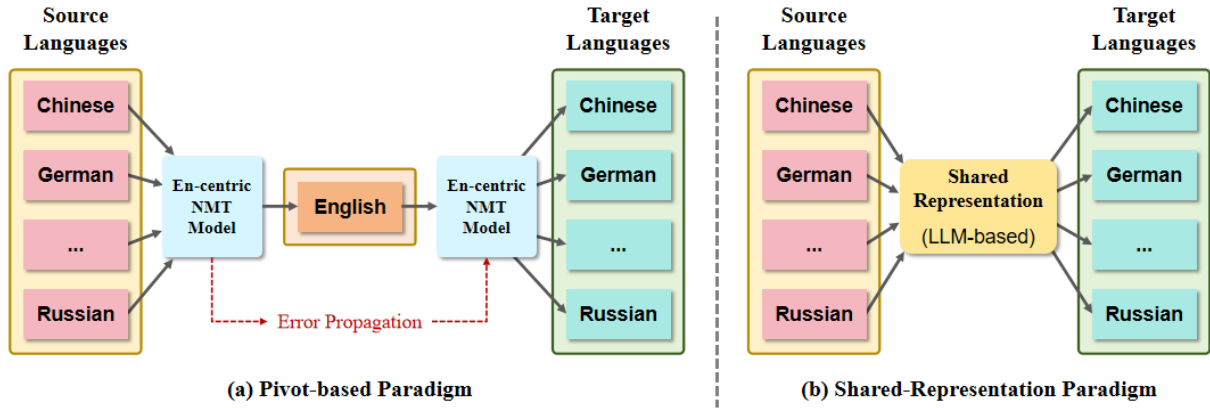
Figure 1: Paradigms for achieving any-to-any translation: (a) Pivot-based paradigm, where translations between language pairs are mediated through a pivot language (typically English). (b) Representation paradigm based on LLM, where languages share a unified representation space, facilitating direct translation between any language pair.

**unscaled LoRA network increase** with the language increase. Specifically, though the backbone LLMs generate a uniform representation for all languages, the multiple LoRAs (Zadouri et al., 2024) will generate a specific (non-uniform) representation for each translation direction. Meanwhile, the LoRA networks will scale up heavily with the increase of language pairs.

To solve these problems, we present **UniLoRA**, which fine-tunes LLMs with unified multiple LoRA networks as illustrated in Figure 2(c). Inspired by Mixture-of-Experts (MoE) (Shazeer et al., 2017), our UniLoRA approach assigns a language with only one language expert, and the matrices $A$ and $B$ in LoRA are represented as the input encoder and output decoder of the assigned language, respectively. This language expert's design greatly alleviates the LoRA scaling-up problem. To unify the fine-tuning representation, a UniCore Module is introduced to merge the multiple LoRA embeddings. Furthermore, a two-stage fine-tune strategy is employed with English-centric parallel corpora in the first stage and limited multilingual corpora in the second stage. We conduct extensive experiments on our two-stage UniLoRA approach, which demonstrates that UniLoRA with hundreds of multilingual parallel sentences outperforms the SoTA MNMT systems with billions of parallel sequences for training. Our main contributions are summarized as follows:

- We introduce the UniLoRA approach, which treats one LoRA network for a specific language. Through the UniCore Module, our UniLoRA approach generates uniform representations from both LLMs backbones and the fine-tuned LoRA network.

- Our UniLoRA approach eliminates dependency on massive non-English parallel corpora through English-involved fine-tuning and multilingual activation with hundreds of high-quality sentence pairs.

- Extensive experiments demonstrate that our UniLoRA approach outperforms existing LoRA-based approaches in quality metrics, enabling fine-tuned general-purpose LLMs to achieve competitive performance of sophisticated MNMT models.

## 2 Related Works

While contemporary multilingual translation resources demonstrate broad linguistic coverage with substantial parallel data availability (e.g., OPUS (Tiedemann, 2012), IWSLT datasets (Cettolo et al., 2017) and WMT datasets (Specia et al., 2021; Kocmi et al., 2022)), these resources predominantly feature English-centric alignment. Uncommon translation directions exhibit scarce or entirely unavailable coverage in publicly accessible datasets, with existing materials for such directions often exhibiting notable quality degradation (Fan et al., 2021). Therefore, synthetic data augmentation addresses this by enhancing NMT training for under-resourced directions (Zhang et al., 2018), with innovations including: pseudo-corpus refinement (Zhang and Matsumoto, 2019; Adjeisah et al., 2021), monolingual-to-parallel expansion (Marie
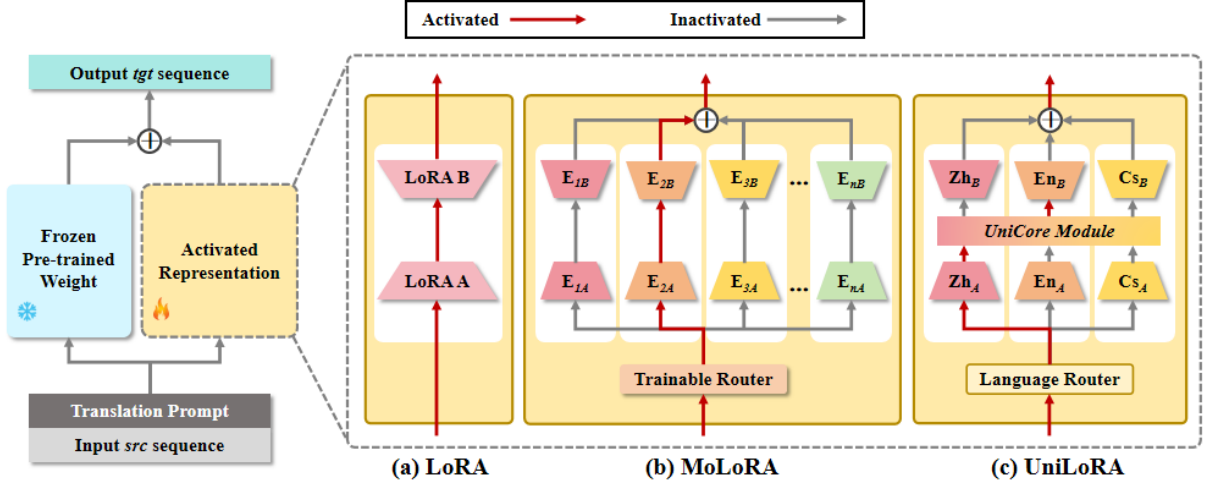
2

Figure 2: Comparison of LLM-based NMT fine-tuning approaches. (a) Standard LoRA for NMT, (b) MoLoRA framework, a straightforward extension for MNMT, and (c) our proposed UniLoRA framework. The UniLoRA architecture distinguishes itself through interconnected language-specific expert design and a shared UniCore module that facilitates unified cross-lingual transfer across all language pairs.

and Fujita, 2021), and noise-reduced generation via graph-prompting LLMs (Pan et al., 2024). Recent works further optimize augmentation strategies during fine-tuning LLMs on the MNMT task (Liu et al., 2023; Lu et al., 2024). In contrast to data-centric approaches, inspired by (Xu et al., 2024), our framework introduces a staged fine-tuning protocol that achieves any-to-any translation proficiency with minimal parallel data requirements, circumventing both synthetic data generation overhead and English-centric bias.

In addition, achieving balanced translation quality across language pairs remains a central challenge in multilingual translation (Tan et al., 2019). Recent innovations in MoE frameworks address this through specialized parameter allocation: typology-aware language group routing (Li et al., 2023), dynamic path optimization (Kudugunta et al., 2021), and task-specific expert decomposition (Tourni and Naskar, 2024). Token-level feature mixing via smoothed gating networks further enhances language-specific feature representation (Liu et al., 2022). Diverging from these structural adaptations, our work mainly focuses on parameter-efficient fine-tuning and eliminating reliance on scarce training data.

## 3 Methodology

### 3.1 Preliminaries

Our approach builds upon two core components: the low-rank adaptation paradigm shown in Figure 2(a), and its extension, Mixture-of-LoRAs

(MoLoRA) (Zhu et al., 2023; Zadouri et al., 2024), as illustrated in Figure 2(b).

When employing the LoRA adapter, the pre-trained LLMs' weight matrix $W_0$ remains frozen, while a trainable low-rank decomposition matrix $\Delta W = BA$ is superimposed onto the selected linear layers. This decomposition consists of two low-rank matrices: $A \in R^{r \times d_i}$ (*LoRA A*) and $B \in R^{d_o \times r}$ (*LoRA B*), where $r \ll min(d_i, d_o)$. The updated forward computation can be formulated as:

$$y = (\Delta W + W_0)x = (BA + W_0)x, \quad (1)$$

where $x \in R^{d_i}$ and $y \in R^{d_o}$ denotes the input and output sequence, respectively.

Building upon the LoRA method, the MoLoRA method integrates the MoE paradigm. A MoLoRA module consists of $N$ LoRA experts, denoted as $E_1, E_2, \ldots, E_n$, which are used to adapt the pre-trained layer during fine-tuning. Each expert $E_i$ is decomposed into two trainable low-rank matrices: $E_{iA}$ and $E_{iB}$, corresponding to the *LoRA A* and *LoRA B* components, respectively. The MoLoRA module further incorporates a trainable token-level expert router $\theta^{MoL}$, which computes routing weights $s_i^{MoL}$ for each expert $E_i$. The routing weight can be calculated as:

$$s_i^{MoL} = \theta^{MoL}(x)_i = softmax(W^{MoL}x)_i, \quad (2)$$

where $W^{MoL} \in R^{N \times d_i}$ is the router's weight matrix. The final output $y$ is computed by aggregating
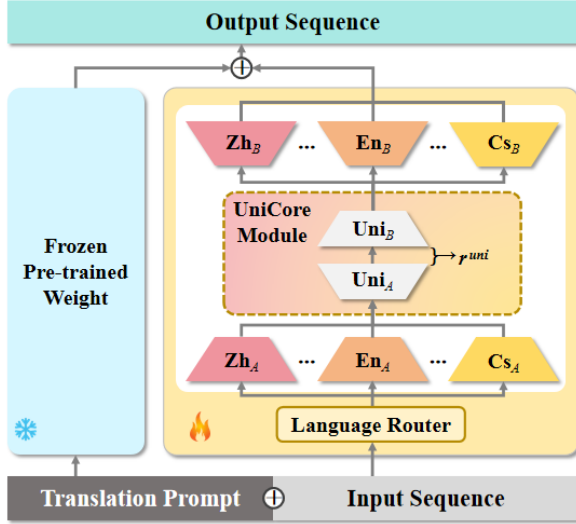
3

Figure 3: Overview of the proposed UniLoRA approach. The shared UniCore module facilitating cross-lingual transfer comprises LoRA $A$ and LoRA $B$ components, which are structurally interconnected with their respective counterparts in each language-specific expert module.

the contributions of all LoRA experts:

$$y = W_0 x + \sum_{i=1}^{n} s_i^{MoL} E_{iB} E_{iA} x \qquad (3)$$

### 3.2 Union-Merged Language Experts with UniCore Module

Our UniLoRA approach establishes union-connections between experts in the MoE structure along with a shared UniCore module as seen in Figure 2(c). In this design, each expert in the MoE structure is assigned to a specific language, while interconnected expert pairs jointly model translation relationships between distinct language pairs.

For an $N$-language multilingual translation task, a dedicated LoRA expert $E_i$ is assigned to each language (e.g., $E_1, E_2, \ldots, E_n$), where $E_i$ corresponds to the $i$-th language. Each language-specific expert is further decomposed into two low-rank matrices:

- $E_{iA} \in R^{d_i \times r}$: Activated when the corresponding language is the source (e.g., $Zh_A$ for Chinese as the source language).

- $E_{iB} \in R^{r \times d_o}$: Activated when the corresponding language is the target (e.g., $En_B$ for English as the target language).

During translation, only the expert modules relevant to the source-target pair and the shared UniCore layer are activated. For example, in the translation direction Chinese $\rightarrow$ English, the $Zh$ and $En$ experts are engaged while all other LoRA experts (e.g., $Cs$) remain inactivated.

To ensure precise activation of the source and target language-specific low-rank matrices during translation, we implement a static language router that routes inputs based on pre-defined language labels, as is shown in Figure 2(c). For a given input sequence $x$ with the source language ($src$) and the target language $tgt$, the output ($y$) is computed as:

$$y = W_0 x + \sum_{i=1}^{n} \sum_{j=1}^{n} f(x;i,j) E_{jB} \Delta W_{uni} E_{iA} x, \qquad (4)$$

where $\Delta W_{uni}$ is the weight matrix of UniCore layer, which can be further decomposed into two trainable low-rank matrices: $Uni_A \in R^{r_{uni} \times r}$ and $Uni_B \in R^{r \times r_{uni}}$, as shown in Figure 3. Besides, $f(x;i,j)$ is the gating function of the static router:

$$f(x;i,j) = \begin{cases} 1 & \text{if } i = src \text{ and } j = tgt \\ 0 & \text{otherwise} \end{cases} \qquad (5)$$

### 3.3 Staged Fine-Tuning on UniLoRA

Our framework implements a two-stage fine-tune to enable comprehensive any-to-any translation with limited parallel data, as formalized below:

**Stage 1: English-Centric Specialization.** The first stage focuses on training the UniLoRA module using English-involved corpora across all languages. The objective is to enhance the model's translation capabilities in both En$\Rightarrow$Any and Any$\Rightarrow$En directions, and facilitate cross-lingual transfer via the UniCore layer by designating English as the pivot language.

Consider a simplified case, where only three language experts are involved: Chinese (Zh), English (En), and Czech (Cs). Let $\mathscr{L} = \{Zh, En, Cs\}$ denote the set of all supported languages, and a translation direction can be defined as $(src, tgt) \in \mathscr{L} \times \mathscr{L}$. In Stage 1, the UniLoRA module is trained on the subset of English-pivoted pairs: $\mathscr{L}_{EN} = \{(En, Zh), (En, Cs), (Zh, En), (Cs, En)\}$. For the subset of non-English pairs $\{(Cs, Zh), (Zh, Cs)\}$, the corresponding routing paths in the UniLoRA module remain inactive. However, all expert weight matrices (in both UniLoRA and UniCore layers) are updated during this stage, enabling the

4

model to learn robust representations for English-involved translation directions.

**Stage 2: Ant-to-Any Activation.** In the second stage, the pre-trained UniLoRA model is further fine-tuned on a limited-size parallel corpus covering all translation directions (e.g., six language pairs in the simplified case). This stage activates all routing paths, allowing each language expert to adapt to both source and target roles across arbitrary language pairs via the UniCore layer. By leveraging the knowledge acquired in Stage 1, the model achieves comprehensive translation capabilities in all directions, ensuring optimal performance regardless of the specific source-target pair.

# 4 Experiments

## 4.1 Dataset and Metrics

Following the ALMA model's configuration (Xu et al., 2024), six languages(i.e., English (En), German (De), Chinese (Zh), Russian (Ru), Czech (Cs), and Icelandic (Is))are selected for evaluation. To comprehensively assess translation performance, we test the model across all 30 possible language directions.

For Stage 1 fine-tuning, the English-centric training set is selected from the OPUS-100 dataset (Tiedemann, 2012). 20k parallel sentence pairs per direction are randomly sampled from 10 English-involved translation directions for Stage 1 training. For Stage 2 activation, all available pairs (with fewer than 1k parallel sentence pairs per direction) are employed from the Flores-200 development set(NLLB Team et al., 2022). Due to the limited availability of non-English-centric test data in OPUS-100, our final evaluation combines the OPUS-100 test sets involving English with the Flores-200 test sets for non-English directions.

We adopt the widely used sentence-level translation prompt template (Hendy et al., 2023), structured as "*Translate the following* $\{src\}$ *sentences into* $\{tgt\}$*:* ". The training loss is not computed for the prompt template or the source sentence itself.

For evaluation metrics, the SacreBLEU (Post, 2018) and COMET-22 (Rei et al., 2022) are selected to evaluate translation quality.

## 4.2 Implementation Details

The UniLoRA framework is implemented on state-of-the-art backbone LLMs, including LLaMA-3-8B-Instruct (Grattafiori et al., 2024) and Qwen2.5-7B-Instruct (Qwen et al., 2025).

For Stage 1 fine-tuning, the training setup includes a batch size of 32, 2 training epochs, and 1 epoch for Stage 2 activation. The initial learning rate is set to 5e-4. Given the six languages involved in the translation task, the number of experts in the MoE structure is fixed at 6. For LoRA configurations, the hyperparameters are set as follows: *lora rank* $r = 16$, *lora alpha* $\alpha = 32$, and *lora dropout p* = 0.1. The UniCore module uses a rank of $r_{uni} = 1024$.

## 4.3 Baselines

To ensure a fair comparison, we evaluate UniLoRA against the following LoRA-based methods under identical staged fine-tuning configurations:

- **LoRA**. We scale up the *lora rank* and *lora alpha* parameters within a single LoRA adapter to match the total number of trainable parameters used in the UniLoRA setup.

- **MoLoRA (Top-k)**. We employ MoLoRA adapter with the same number of experts as UniLoRA alongwith a top-1 router, activating only one expert per translation process.

- **MoLoRA (Static)**. A variant of MoLoRA equipped with a static router, where each translation direction is assigned a dedicated expert, increases the total number of experts to 30. While this ensures deterministic activation and eliminates routing instability, it incurs 3.8× higher training overhead due to redundant expert allocation.

In addition, we benchmark the UniLoRA framework against state-of-the-art multilingual translation models based on LLMs, including **M2M100-12B** (Fan et al., 2021), **BigTranslate-13B** (Yang et al., 2023), and **NLLB-3.3B** (NLLB Team et al., 2022). These models represent leading end-to-end paradigms for multilingual translation. To further contextualize our findings, we include **ALMA-7B** (Xu et al., 2024), an English-centric model that employs a staged fine-tuning strategy similar to our framework. Notably, ALMA-7B's performance in non-English-centric directions is evaluated via an English pivot translation pipeline.

## 4.4 Main Experiments

The comprehensive performance across all translation directions is presented in Table 1. Overall,

| Models | Training Tokens | En-Centric | | non-En-centric | | Average | |
|---|---|---|---|---|---|---|---|
| | | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| ALMA-7B (En-pivot) | > 20B | 24.45 | 78.12 | 17.66 | 80.79 | 19.92 | 79.90 |
| M2M100-12B | > 7.5B | 24.06 | 74.59 | 18.98 | 82.52 | 20.68 | 79.88 |
| BigTranslate-13B | 89.8B | 22.02 | 72.95 | 18.94 | 81.88 | 19.98 | 78.90 |
| NLLB-3.3B | > 21.5B | 27.85 | 77.01 | 20.53 | 82.88 | 22.97 | 80.92 |
| **LLaMA-3-8B-Instruct** | Trainable Parameters | Training Tokens | En-Centric | | non-En-centric | | Average | |
| | | | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| Base | — | — | 18.41 | 68.61 | 13.94 | 77.41 | 15.43 | 74.48 |
| +LoRA | 2.09% | | 25.45 | 75.34 | 15.77 | 79.60 | 19.00 | 78.18 |
| +MoLoRA (Top-k) | 2.14% | 12.4M | 26.03 | 75.25 | 15.97 | 80.02 | 19.32 | 78.43 |
| +MoLoRA (Static) | 10.18% | | 27.46 | 77.10 | 16.39 | 81.27 | 20.08 | 79.88 |
| +UniLoRA | | | | | | | | |
| — Stage 1 | 2.12% | 12.4M | **29.97** | **77.54** | 5.13 | 58.72 | 13.41 | 64.99 |
| — Stage 2 | | | 27.61 | 77.21 | **17.98** | **81.91** | **21.19** | **80.35** |
| **Qwen2.5-7B-Instruct** | Trainable Parameters | Training Tokens | En-Centric | | non-En-centric | | Average | |
| | | | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| Base | — | — | 19.09 | 70.37 | 12.85 | 76.55 | 14.93 | 74.49 |
| +LoRA | 2.08% | | 24.76 | 74.83 | 16.20 | 79.89 | 19.05 | 78.20 |
| +MoLoRA (Top-k) | 2.14% | 12.4M | 25.37 | 74.96 | 16.35 | 80.04 | 19.36 | 78.35 |
| +MoLoRA (Static) | 10.17% | | 26.15 | 76.61 | 17.09 | 80.62 | 20.11 | 79.28 |
| +UniLoRA | | | | | | | | |
| — Stage 1 | 2.12% | 12.4M | **27.97** | **78.51** | 7.44 | 59.40 | 14.28 | 65.77 |
| — Stage 2 | | | 26.92 | 77.63 | **18.35** | **81.70** | **21.21** | **80.34** |

Table 1: The overall results in all directions. We mark the amount of tokens in training data and the proportion of trainable parameters in the table as well. Except for UniLoRA, all LoRA-based fine-tuning approaches report only Stage 2 results. UniLoRA outperforms all other fine-tuning configurations and is comparable to the state-of-the-art translation models on multilingual translation tasks. **Bold results** indicate the highest scores among fine-tuning methods on the same backbone model.

the proposed UniLoRA method demonstrates superior effectiveness after Stage 2 fine-tuning, outperforming other LoRA-based fine-tuning approaches. The optimized model achieves competitive performance relative to state-of-the-art multilingual translation systems, with results closely aligning with the NLLB model.

**Compared with backbone LLMs.** After Stage 1 English-centric fine-tuning, UniLoRA significantly enhances translation quality for all English-involved directions. Following Stage 2 any-to-any activation, UniLoRA exhibits statistically significant improvements across all language pairs, with particularly pronounced gains in non-English translation directions, achieving an average BLEU score increase of +6.28. However, notable performance degradation can be observed in English-centric directions compared to Stage 1 results, attributed to the knowledge forgetting of English-specific patterns during the activation phase.

**Comparison with LoRA-based fine-tuning methods.** The experiments reveal that UniLoRA achieves the most significant performance gains among all LoRA-based approaches under identical or less parameter budgets, with consistent improvements observed across key evaluation metrics. Notably, while MoLoRA's top-k routing strategy outperforms standard LoRA by dynamically selecting experts, its static routing variant, though achieving enhanced accuracy, demands multiplied computational resources due to its fixed expert assignment. After Stage 2 fine-tuning, MoLoRA with top-k routing experiences expert fluctuations when the number of experts is insufficient, leading to a general performance decline. In contrast, UniLoRA consistently outperforms both MoLoRA configurations, demonstrating superior efficiency.

**Compared with prior multilingual translation models.** The UniLoRA-fine-tuned model outperforms most existing multilingual translation systems and achieves performance comparable to the NLLB model in average metrics. Notably, while the ALMA model exhibits strong performance in English-centric translation directions compared to

| Methods | English-Centric | | non-English-centric | | Average | |
|---|---|---|---|---|---|---|
| | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| UniLoRA | | | | | | |
| — Stage 1 | **29.97** | 77.54 | 5.13 | 58.72 | 13.41 | 64.99 |
| — Stage 2 | 27.61 | 77.21 | **17.98** | **81.91** | **21.19** | **80.35** |
| *w/o* Stage 1 | 23.48 | 73.17 | 17.36 | 81.11 | 19.40 | 78.47 |
| *w/o* UniCore | | | | | | |
| — Stage 1 | 29.50 | **78.26** | 13.95 | 77.35 | 19.13 | 77.65 |
| — Stage 2 | 26.59 | 75.22 | 17.22 | 81.30 | 20.34 | 79.27 |

Table 2: The results of the ablation study on UniLoRA with different framework and fine-tuning configurations, based on the LLaMA-3-8B-Instruct model. The best scores are marked in **bold**.

| Methods | English-Centric | | non-English-centric | | Average | | Trainable |
|---|---|---|---|---|---|---|---|
| | BLEU | COMET | BLEU | COMET | BLEU | COMET | Parameters |
| 1 Shared Source Expert | 26.49 | 76.37 | 17.70 | 80.52 | 20.63 | 79.14 | 1.28% |
| 1 Shared Target Expert | 26.18 | 76.10 | 17.81 | 80.55 | 20.60 | 79.07 | |
| 3 Experts | 27.32 | 76.94 | **18.04** | 81.27 | 21.13 | 79.83 | 1.10% |
| 6 Experts | **27.61** | **77.21** | 17.98 | **81.91** | **21.19** | **80.35** | 2.12% |

Table 3: The results of the ablation study on UniLoRA with different configurations of merged language experts, based on the LLaMA-3-8B-Instruct model. The best scores are marked in **bold**.

other baselines, its effectiveness in non-English-centric directions is significantly limited when relying on an English pivot pipeline for any-to-any translation. This performance gap highlights the inherent limitations of English-centric architectures in direct cross-lingual scenarios. UniLoRA's advantages are further underscored by its ability to reduce dependency on large-scale non-English parallel corpora, which have traditionally been considered essential for robust multilingual translation. Despite this reduction in data requirements, the framework maintains competitive performance across diverse language pairs, showcasing its efficiency in parameter utilization without compromising translation quality.

## 5 Ablation Studies

We conduct further research on the UniLoRA framework with diverse configurations to gain a more comprehensive understanding. All experiments for analysis are conducted on the LLaMA-3-8B-Instruct model.

### 5.1 Component and Staged Training Analysis

To validate the necessity of the UniCore module in the UniLoRA framework and the effectiveness of staged fine-tuning, we perform ablation experiments by modifying the UniLoRA framework or training process. As shown in Table 2, we systematically evaluate the impact of individual components and training strategies, with corresponding results presented for comparative analysis. The results reveal two critical findings:

**The UniCore module is essential for any-to-any translation.** After Stage 1 fine-tuning, the model without the UniCore module achieves comparable performance to the full UniLoRA model in English-centric directions (notably with higher COMET scores) and retains the backbone model's translation capability in non-English-centric directions. However, following Stage 2 fine-tuning, the UniCore-free variant underperforms the baseline UniLoRA model in both English-centric and non-English-centric directions, demonstrating its critical role in enabling robust any-to-any translation capabilities.

**Staged fine-tuning is indispensable.** When Stage 1 fine-tuning is omitted and the model is trained directly on a limited any-to-any corpus, performance drops significantly compared to the optimal configuration. This decline is observed across all translation directions, including non-English-centric ones, highlighting the importance of Stage 1 in transferring and preserving cross-lingual capabilities before any-to-any activation.

### 5.2 Merged Language Experts

Exploring the application of expert compression techniques within the UniLoRA framework is critical for advancing parameter efficiency. To further
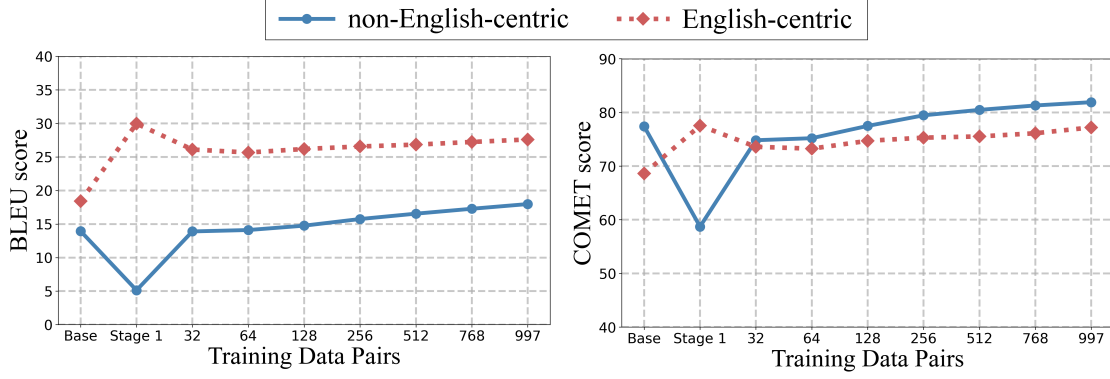
Figure 4: Data Requirements for Stage 2 activation.

optimize resource allocation, we conduct experiments employing two distinct expert compression strategies.

**Shared source-target language experts.** Inspired by the asymmetric architecture in HydraLoRA (Tian et al., 2024), we implement configurations where a unified expert module handles either all source language inputs or all target language outputs. This design enables shared expertise between source and target modalities, thereby reducing redundancy and optimizing parameter usage. By consolidating language-specific adaptations into fewer specialized experts, the framework achieves greater efficiency without compromising directional translation performance.

**Linguistically-informed group experts.** Drawing from the integration of language typology in MoE-based translation systems (Li et al., 2023), we propose linguistically-informed group experts. As detailed in Table 4, selected languages are clustered into typological groups. For instance, English, German and Icelandic are consolidated into a single expert. This approach reduces the total number of language experts, decreasing trainable parameters from 2.12% to 1.10% of the total model parameters while maintaining the core UniLoRA architecture.

Experimental results in Table 3 demonstrate that these strategies achieve significant parameter compression with performance trade-offs. The linguistically-informed grouping configuration even shows improved performance on non-English-centric directions, while maintaining robustness across major language pairs. These findings establish a scalable pathway for expanding the UniLoRA framework to additional languages, demonstrating its capacity for efficient multilingual adaptation without compromising translation quality.

### 5.3 Data Requirements for Any-to-any Activation

We further investigate the data requirements for Stage 2 any-to-any adaptation in the UniLoRA framework. To analyze scalability, we subsample the original training data into subsets with varying sizes, where the number of parallel sentences per translation direction ranges from 32 to 997 (i.e., the Flores-200 development set size). Notably, larger subsets hierarchically include smaller ones to ensure consistent comparisons. Using these subsets for Stage 2 training, we evaluate how translation performance scales with data quantity, with results visualized in Figure 4.

Experimental results reveal two critical phenomena: First, Stage 2 fine-tuning temporarily degrades performance in English-centric directions, but this knowledge degradation diminishes as training data increases; Second, non-English-centric directions require approximately 128 parallel sentence pairs to activate translation capabilities comparable to the backbone model, with further performance gains achieved through additional data scaling.

### 6 Conclusion

We present UniLoRA, a framework that integrates LoRA with an MoE architecture to enable efficient multilingual translation in large language models. By combining language-specific expert modules with a shared unified layer, UniLoRA achieves robust any-to-any translation capabilities through a two-stage training approach that eliminates reliance on extensive non-English parallel corpora. Extensive experiments demonstrate that UniLoRA is a scalable solution for multilingual translation, offering both technical innovation and practical value for resource-constrained deployment scenarios.

## Limitations

This work provides an efficient framework for multilingual neural machine translation via LLM fine-tuning, yet several key limitations remain and warrant further investigation.

**Language Coverage Constraints.** While the UniLoRA framework demonstrates potential to reduce reliance on non-English training corpora, our experiments are limited to six languages (including one low-resource language: Icelandic). Although ablation studies on linguistically-informed expert groups in Section 5.2 suggest language scalability of UniLoRA, empirical validation is required to evaluate performance across diverse language pairs, particularly in low-resource settings.

**Knowledge Forgetting in Staged Training.** The staged fine-tuning process leads to degraded performance in English-centric directions due to knowledge forgetting of Stage 1 capabilities. This highlights the need for architectural innovations to preserve cross-stage knowledge retention.

**Model-Specific Generalization.** Our experiments are conducted on LLaMA-3-8B-Instruct and Qwen2.5-7B-Instruct, which represent strong baselines but limit insights into model diversity and size scalability. Future work should systematically evaluate UniLoRA's effectiveness across models with varying capabilities to ensure broader applicability.

**Diversified Training Process.** Our research primarily explores supervised fine-tuning of LLMs using parallel corpora. However, recent studies indicate that translation capabilities can be further enhanced through techniques such as continual pre-training with monolingual data (Xu et al., 2024) and preference learning (Xu et al., 2025). Further exploration of integrating these methods with UniLoRA is essential for enhancing its versatility.

## References

Michael Adjeisah, Guohua Liu, Douglas Omwenga Nyabuga, Richard Nuetey Nortey, and Jinling Song. 2021. Pseudotext injection and advance filtering of low-resource corpus for neural machine translation. *Computational Intelligence and Neuroscience*.

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *Preprint*, arXiv:1907.05019.

Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuitho Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. Overview of the iwslt 2017 evaluation campaign. In *Proceedings of the 14th International Workshop on Spoken Language Translation, IWSLT 2017*.

Kaidi Chen, Ben Chen, Dehong Gao, Huangyu Dai, Wen Jiang, Wei Ning, Shanqing Yu, Libin Yang, and Xiaoyan Cai. 2024. General2specialized llms translation for e-commerce. In *Companion Proceedings of the ACM on Web Conference 2024, WWW 2024*.

Yun Chen, Yang Liu, Yong Cheng, and Victor OK Li. 2017. A teacher-student framework for zero-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2017*.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, and 1 others. 2021. Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*.

Dehong Gao, Kaidi Chen, Ben Chen, Huangyu Dai, Linbo Jin, Wen Jiang, Wei Ning, Shanqing Yu, Qi Xuan, Xiaoyan Cai, and 1 others. 2024. Llms-based machine translation for e-commerce. *Expert Systems with Applications*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *Preprint*, arXiv:2302.09210.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022*.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, and 1 others. 2022. Findings of the 2022 conference on machine translation (wmt22). In *Proceedings of the Seventh Conference on Machine Translation, WMT 2022*.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, and 1 others. 2022. Quality at a glance: An

audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*.

Sneha Kudugunta, Yanping Huang, Ankur Bapna, Maxim Krikun, Dmitry Lepikhin, Minh-Thang Luong, and Orhan Firat. 2021. Beyond distillation: Task-level mixture-of-experts for efficient inference. In *Findings of the Association for Computational Linguistics: EMNLP 2021*.

Shangjie Li, Xiangpeng Wei, Shaolin Zhu, Jun Xie, Baosong Yang, and Deyi Xiong. 2023. Mmnmt: Modularizing multilingual neural machine translation with flexibly assembled moe and dense blocks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023*.

Xiaobo Liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, and Tie-Yan Liu. 2021. R-drop: Regularized dropout for neural networks. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021*.

Chao-Hong Liu, Catarina Cruz Silva, Longyue Wang, and Andy Way. 2018. Pivot machine translation using chinese as pivot language. In *China Workshop on Machine Translation*.

Junpeng Liu, Kaiyu Huang, Jiuyi Li, Huan Liu, Jinsong Su, and Degen Huang. 2022. Adaptive token-level cross-lingual feature mixing for multilingual neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*.

Xiner Liu, Jianshu He, Mingzhe Liu, Zhengtong Yin, Lirong Yin, and Wenfeng Zheng. 2023. A scenario-generic neural machine translation data augmentation method. *Electronics*.

Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. Llamax: Scaling linguistic horizons of llm by enhancing translation capabilities beyond 100 languages. In *Findings of the Association for Computational Linguistics: EMNLP 2024*.

Benjamin Marie and Atsushi Fujita. 2021. Synthesizing monolingual data for neural machine translation. *Preprint*, arXiv:2101.12462.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, and 20 others. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

Shilong Pan, Zhiliang Tian, Liang Ding, Haoqi Zheng, Zhen Huang, Zhihua Wen, and Dongsheng Li. 2024. Pomp: Probability-driven meta-graph prompter for llms in low-resource unsupervised neural machine translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024*.

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

Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, and 24 others. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation, WMT 2022*.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Édouard Grave, Armand Joulin, and Angela Fan. 2021. Ccmatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*.

Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *5th International Conference on Learning Representations, ICLR 2017*.

Lucia Specia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary, and André FT Martins. 2021. Findings of the wmt 2021 shared task on quality estimation. In *Proceedings of the Sixth Conference on Machine Translation, WMT 2021*.

Xu Tan, Yi Ren, Di He, Tao Qin, Zhou Zhao, and Tie-Yan Liu. 2019. Multilingual neural machine translation with knowledge distillation. In *7th International Conference on Learning Representations, ICLR 2019*.

Chunlin Tian, Zhan Shi, Zhijiang Guo, Li Li, and Cheng-Zhong Xu. 2024. Hydralora: An asymmetric lora architecture for efficient fine-tuning. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024*.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation, LREC 2012*.

Isidora Chara Tourni and Subhajit Naskar. 2024. Direct neural machine translation with task-level mixture of experts models. *Preprint*, arXiv:2310.12236.

Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, and 1 others. 2024. Aya model: An instruction finetuned open-access multilingual language model. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024*.

Haoran Xu, Kenton Murray, Philipp Koehn, Hieu Hoang, Akiko Eriguchi, and Huda Khayrallah. 2025. X-alma: Plug & play modules and adaptive rejection for quality translation at scale. *Preprint*, arXiv:2410.03115.

Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages. *Preprint*, arXiv:2305.18098.

Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermis, Acyr Locatelli, and Sara Hooker. 2024. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. In *The Twelfth International Conference on Learning Representations, ICLR 2024*.

Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. 2024. Teaching large language models to translate with comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence, AAAI 2024*.

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, ACL 2020*.

Jinyi Zhang and Tadahiro Matsumoto. 2019. Corpus augmentation for neural machine translation with chinese-japanese parallel corpora. *Applied sciences*.

Zhirui Zhang, Shujie Liu, Mu Li, Ming Zhou, and Enhong Chen. 2018. Joint training for neural machine translation models with monolingual data. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*.

Yun Zhu, Nevan Wichers, Chu-Cheng Lin, Xinyi Wang, Tianlong Chen, Lei Shu, Han Lu, Canoee Liu, Liangchen Luo, Jindong Chen, and Lei Meng. 2023. Sira: Sparse mixture of low rank adaptation. *Preprint*, arXiv:2311.09179.

11

## A Appendix

### A.1 Training Details

We hereby supplement the model training configuration not mentioned in the main text. For both backbone LLMs, we fine-tune the models using a warm-up ratio of 5e-4, a maximum sequence length of 512 tokens, and a weight decay of 0.02. LoRA adapters are applied to the *gate_proj*, *up_proj*, and *down_proj* modules of the backbone LLMs. Stage 1 fine-tuning requires 3 epochs, while Stage 2 activation requires 1 epoch. Model training process is conducted on 4 NVIDIA A100 GPUs, with each GPU handling batches with *batch_size* of 2 and employing a *gradient_accumulation_step* of 4.

### A.2 Data Settings

For the staged fine-tuning data details:

**Stage 1 English-Centric Specialization:** The pre-divided development subset from OPUS-100 serves as our development set. The training data consists of the randomly sampled subset from the OPUS-100 training dataset.

**Stage 2 Ant-to-Any Activation:** In this stage, non-English-centric directions use the full Flores-200 development subsets for training, with 20% of the randomly sampled training data serving as the development set. For English-centric directions, the training data consists of 5% of the randomly sampled subset from Stage 1's training data. Detailed configurations are summarized in Table 5.

### A.3 Self-Contrastive Semantic Enhancement

To further improve regularization capability, we take R-Drop (Liang et al., 2021) to reduce the inconsistency existing in training and inference. In each training step, the R-Drop method seeks to regularize the model's predictions by minimizing the bidirectional Kullback-Leibler (KL) divergence between the two output distributions for the same sample.

| Training Stage | Directions | Parallel data pairs | | |
|---|---|---|---|---|
| | | train | dev | test |
| Stage 1 | En⇔Any | 20000 | 2000 | 2000 |
| Stage 2 | En⇔Any | 1000 | 200 | 2000 |
| | Others | 997 | 200 | 1012 |

Table 5: The statistics for the data we utilize for the main experiment.

We evaluate its impact via ablation studies on the UniLoRA model based on LLaMA-3-8B-Instruct, comparing fine-tuning with and without R-Drop. Results in Table 6 demonstrate that self-contrastive semantic enhancement significantly boosts the generalization capability of UniLoRA, achieving consistent performance improvements across all translation directions relative to the baseline, while incurring no additional inference cost. This highlights the effectiveness of R-Drop in stabilizing training dynamics without compromising computational efficiency.

### A.4 Full Results of the Main Experiment

We present in Table 7 and Table 8 the specific performance of the UniLoRA model based on the LLaMA-3-8B-Instruct backbone model across all translation directions in the main experiment. The performance metrics include BLEU scores, ROUGE-L, and COMET scores. For comparison, the table also includes the performance of prior studies and the backbone LLM baseline.

| Language | Language Family |
|---|---|
| (En) English | |
| (De) German | Germanic, Indo-European |
| (Is) Icelandic | |
| (Cs) Czech | Balto-Slavic, Indo-European |
| (Ru) Russian | |
| (Zh) Chinese | Sino-Tibetan |

Table 4: The languages selected in the main experiment and their corresponding language families.

| Configurations | English-Centric | | non-English-centric | | Average | |
|---|---|---|---|---|---|---|
| | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| *w/o* R-Drop | 27.33 | 76.79 | 17.47 | 81.46 | 20.76 | 79.90 |
| *with* R-Drop | **27.61** | **77.21** | **17.98** | **81.91** | **21.19** | **80.35** |

Table 6: Results of the ablation study on the effect of R-Drop regularization, based on the LLaMA-3-8B-Instruct backbone model. Higher scores are marked in **bold**. Employing the R-Drop method results in a comprehensive performance improvement.

| Models | Zh⇒En | | | En⇒Zh | | | De⇒En | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 29.75 | 50.11 | 79.51 | 28.07 | 44.35 | 80.58 | 28.08 | 45.06 | 75.24 |
| M2M100-12B | 27.66 | 51.72 | 78.97 | 27.76 | 45.06 | 79.81 | 30.90 | 50.48 | 78.27 |
| LLaMA-3-8B-Instruct | 20.84 | 39.64 | 74.93 | 18.76 | 33.86 | 73.47 | 23.07 | 37.75 | 71.10 |
| UniLoRA Stage 1 | 35.28 | 58.52 | 81.26 | 36.19 | 53.84 | 82.49 | 33.45 | 54.47 | 78.46 |
| UniLoRA Stage 2 | 33.08 | 56.40 | 81.75 | 34.56 | 52.85 | 82.73 | 31.70 | 52.51 | 79.03 |
| **Models** | **En⇒De** | | | **Ru⇒En** | | | **En⇒Ru** | | |
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 27.54 | 43.22 | 78.24 | 29.03 | 47.86 | 76.80 | 28.41 | 43.21 | 82.51 |
| M2M100-12B | 27.29 | 45.93 | 76.48 | 26.65 | 46.34 | 76.97 | 23.39 | 36.81 | 79.54 |
| LLaMA-3-8B-Instruct | 21.26 | 34.14 | 70.36 | 22.54 | 38.50 | 71.10 | 18.74 | 30.84 | 73.84 |
| UniLoRA Stage 1 | 28.34 | 48.77 | 76.45 | 32.69 | 55.36 | 78.25 | 26.44 | 46.54 | 80.92 |
| UniLoRA Stage 2 | 27.61 | 47.87 | 78.31 | 31.25 | 54.02 | 79.07 | 22.81 | 41.53 | 78.63 |
| **Models** | **Cs⇒En** | | | **En⇒Cs** | | | **Is⇒En** | | |
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 31.10 | 47.71 | 76.15 | 28.11 | 41.25 | 81.39 | 25.47 | 43.63 | 72.63 |
| M2M100-12B | 26.12 | 41.56 | 76.31 | 21.19 | 32.69 | 77.70 | 16.41 | 38.40 | 64.67 |
| LLaMA-3-8B-Instruct | 22.58 | 38.45 | 70.06 | 15.38 | 25.64 | 71.02 | 11.89 | 21.46 | 55.51 |
| UniLoRA Stage 1 | 33.86 | 57.03 | 78.82 | 24.89 | 46.60 | 80.19 | 26.55 | 51.01 | 71.08 |
| UniLoRA Stage 2 | 30.97 | 53.90 | 77.96 | 22.82 | 43.61 | 79.67 | 22.90 | 45.78 | 69.62 |
| **Models** | **En⇒Is** | | | **De⇒Zh** | | | **Zh⇒De** | | |
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 22.98 | 38.37 | 67.03 | 25.11 | 43.38 | 79.31 | 18.17 | 41.43 | 80.52 |
| M2M100-12B | 13.21 | 32.84 | 57.15 | 27.24 | 48.11 | 84.06 | 16.47 | 39.34 | 80.09 |
| LLaMA-3-8B-Instruct | 9.05 | 17.29 | 54.71 | 16.81 | 32.70 | 76.69 | 13.26 | 33.07 | 77.25 |
| UniLoRA Stage 1 | 21.98 | 45.42 | 67.45 | 5.39 | 10.34 | 51.08 | 9.74 | 26.48 | 70.88 |
| UniLoRA Stage 2 | 18.36 | 40.46 | 65.35 | 27.61 | 48.59 | 83.56 | 15.28 | 40.58 | 79.95 |
| **Models** | **De⇒Ru** | | | **Ru⇒De** | | | **De⇒Cs** | | |
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 25.29 | 46.45 | 87.12 | 24.17 | 49.12 | 81.89 | 24.13 | 47.40 | 89.48 |
| M2M100-12B | 22.07 | 43.26 | 86.55 | 21.30 | 45.92 | 80.52 | 23.35 | 46.50 | 89.61 |
| LLaMA-3-8B-Instruct | 17.79 | 36.24 | 82.52 | 17.84 | 39.78 | 77.34 | 17.11 | 37.73 | 85.10 |
| UniLoRA Stage 1 | 1.79 | 2.20 | 52.71 | 1.99 | 2.82 | 49.45 | 1.86 | 2.92 | 58.89 |
| UniLoRA Stage 2 | 21.58 | 45.51 | 86.48 | 20.81 | 47.84 | 81.24 | 19.87 | 45.66 | 89.09 |

Table 7: Part 1 of the full results for all translation directions of the main experiment.

13

| Models | Cs⇒De | | | De⇒Is | | | Is⇒De | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 26.02 | 51.21 | 84.59 | 18.19 | 43.15 | 82.37 | 20.63 | 44.89 | 78.80 |
| M2M100-12B | 24.00 | 49.15 | 83.55 | 13.72 | 37.02 | 79.35 | 18.99 | 42.91 | 78.30 |
| LLaMA-3-8B-Instruct | 20.26 | 42.71 | 80.47 | 8.06 | 26.90 | 72.05 | 9.92 | 24.22 | 66.90 |
| UniLoRA Stage 1 | 4.33 | 8.63 | 61.52 | 2.02 | 4.50 | 52.84 | 7.12 | 17.78 | 59.31 |
| UniLoRA Stage 2 | 22.35 | 49.43 | 83.80 | 11.56 | 35.19 | 74.37 | 18.02 | 43.79 | 77.75 |

| Models | Zh⇒Ru | | | Ru⇒Zh | | | Zh⇒Cs | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 17.50 | 36.90 | 85.57 | 24.96 | 42.84 | 80.22 | 15.64 | 36.42 | 86.20 |
| M2M100-12B | 15.76 | 34.68 | 85.39 | 26.10 | 46.57 | 83.60 | 14.87 | 35.39 | 86.59 |
| LLaMA-3-8B-Instruct | 12.25 | 27.87 | 81.79 | 23.02 | 71.13 | 79.93 | 11.29 | 28.05 | 82.84 |
| UniLoRA Stage 1 | 9.03 | 22.22 | 66.63 | 5.12 | 9.02 | 49.33 | 9.77 | 25.83 | 81.13 |
| UniLoRA Stage 2 | 15.55 | 37.14 | 86.15 | 27.87 | 49.56 | 84.03 | 13.71 | 37.34 | 86.08 |

| Models | Cs⇒Zh | | | Zh⇒Is | | | Is⇒Zh | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 24.38 | 42.81 | 79.50 | 12.81 | 34.72 | 79.63 | 20.83 | 38.99 | 77.27 |
| M2M100-12B | 26.96 | 47.80 | 84.38 | 9.89 | 30.84 | 77.44 | 21.14 | 42.30 | 80.73 |
| LLaMA-3-8B-Instruct | 18.49 | 34.99 | 77.71 | 6.34 | 21.53 | 71.12 | 16.20 | 33.11 | 73.70 |
| UniLoRA Stage 1 | 3.99 | 7.15 | 45.90 | 2.52 | 8.54 | 51.92 | 0.93 | 2.64 | 46.18 |
| UniLoRA Stage 2 | 30.86 | 51.94 | 85.04 | 7.26 | 29.21 | 72.07 | 22.12 | 44.14 | 80.06 |

| Models | Cs⇒Ru | | | Ru⇒Cs | | | Cs⇒Is | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 24.26 | 45.55 | 87.82 | 21.10 | 43.93 | 88.36 | 16.43 | 40.52 | 81.39 |
| M2M100-12B | 21.65 | 43.34 | 87.76 | 20.18 | 42.64 | 88.99 | 12.49 | 35.55 | 76.42 |
| LLaMA-3-8B-Instruct | 17.69 | 36.25 | 83.09 | 15.04 | 35.02 | 83.93 | 7.97 | 26.07 | 71.36 |
| UniLoRA Stage 1 | 16.77 | 34.99 | 81.83 | 1.62 | 2.32 | 53.20 | 7.74 | 24.53 | 71.23 |
| UniLoRA Stage 2 | 20.43 | 44.02 | 87.53 | 17.80 | 43.20 | 87.73 | 10.04 | 34.97 | 76.22 |

| Models | Is⇒Cs | | | Ru⇒Is | | | Is⇒Ru | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET | BLEU | ROUGE | COMET |
| NLLB-3.3B | 17.32 | 38.99 | 84.35 | 15.23 | 39.21 | 80.58 | 18.43 | 38.32 | 82.64 |
| M2M100-12B | 16.05 | 36.89 | 82.39 | 11.49 | 33.54 | 74.37 | 15.96 | 35.35 | 80.20 |
| LLaMA-3-8B-Instruct | 10.77 | 27.58 | 77.63 | 7.31 | 25.06 | 70.51 | 11.60 | 27.48 | 76.28 |
| UniLoRA Stage 1 | 4.01 | 9.53 | 58.41 | 3.88 | 13.65 | 60.07 | 2.98 | 6.42 | 51.88 |
| UniLoRA Stage 2 | 13.74 | 36.94 | 83.63 | 8.65 | 32.76 | 71.77 | 14.53 | 35.98 | 81.68 |

Table 8: Part 2 of the full results for all translation directions of the main experiment.