

TRANS-ZERO: Self-Play Incentivizes Large Language Models to Achieve Multilingual Translation Without Parallel Data

Anonymous ACL submission

Abstract

The rise of Large Language Models (LLMs) has reshaped machine translation (MT), but multilingual MT still relies heavily on parallel data for supervised fine-tuning (SFT), facing challenges like data scarcity for low-resource languages and catastrophic forgetting. To address these issues, we propose TRANS-ZERO, a self-play framework that leverages only monolingual data and the intrinsic multilingual knowledge of LLM. TRANS-ZERO combines a novel Monte-Carlo Tree Search, GMCTS, with preference optimization, achieving strong translation performance that rivals supervised methods. Experiments demonstrate that this approach not only matches the performance of models trained on large-scale parallel data but also excels in non-English translation directions. Further analysis reveals that GMCTS itself significantly enhances translation quality by exploring semantically consistent candidates through iterative translations, providing a robust foundation for the framework’s success.

1 Introduction

The advent of Large Language Models (LLMs) witnesses a fundamental shift in machine translation (MT) paradigms from the supervised end-to-end training (Vaswani et al., 2017) to the sophisticated generation of fine-tuned language models (Achiam et al., 2023).

Unlike various downstream tasks that gain impressive proficiency through lightweight instruction tuning, multilingual translations between languages still necessitate sufficient fine-tuning with parallel data for specific translation directions. Besides external human annotations, researchers like Xu et al. (2024b); Li et al. (2024) obtain translation annotations or preferences from external LLM, as long as they prove multilingual capability. Either way, it faces a notable deficit in data scarcity for

less popular languages. Moreover, issues arise as the multilingual translation fine-tuning scales up. The reliance on one-on-one MLE supervision has been criticized for potential biases that clash with natural language’s inherent multilingualism (Zhu et al., 2024a) and poses a risk of catastrophic forgetting. Furthermore, exceeding multilingual annotations inversely dilutes the pre-trained knowledge in supported languages, thus degrading overall cross-lingual performance (Xu et al., 2023; Zhu et al., 2024b). Xu et al. (2024a) proposes a mixture-of-expert with hand-crafted route across language modules. However, the route and distributed overheads increase exponentially as the number of translation directions involved increases.

Whereas traditional fine-tuning scaling approaches a plateau, leveraging LLM’s inherent knowledge rather than external supervision for self-improvement is trending (Chen et al., 2024; Kumar et al., 2024). However, adapting this approach to MT introduces two technical challenges. First, systematic *cross-lingual exploration* requires navigating complex semantic spaces beyond simple prompt engineering. Traditional LLM planning involves delicate prompt-based reasoning, even fine-tuning, which is generally unavailable for most scenarios. Second, *multilingual quality assessment* must overcome the limitations of data-dependent quality estimation (QE) metrics and reward model training complexities.

In this work, we introduce TRANS-ZERO. This innovative self-play framework enables LLMs to bootstrap their multilingualism by strategically exploring their semantic space, achieving self-improvement for multilingual translation given only monolingual data. We start by defining a Multilingual Translation Process (MTP) that generates translations by interweaving the languages supported by LLM, so the inference is scaled up to explore more potential translations. First, we

implement the Genetic Monte-Carlo tree search (G-MCTS) upon MTP, exploring potential translations. Second, we harness the search to assess translation preferences based on intrinsic multilingual consistency. Experiments verify that TRANS-ZERO improves the lesser translations through iterative G-MCTS and preference optimizations given only monolingual data. We summarize our contributions as follows:

- First self-play framework extending to multilingual MT training with monolingual data.
- Novel integration of MCTS that explores improved translation for preference.
- An intrinsic translation preference without additional QE modules enables iterative self-improvement.

2 Preliminary

2.1 Machine Translation via LLM

Traditional machine translation depends on large parallel supervision in specific language pairs, which is limited by insufficient annotation in less popular translation directions. The emergence of multilingual pre-trained language models has revolutionized the field, enabling impressive performance across various downstream tasks with minimal supervision (Wei et al., 2022a). State-of-the-art (SOTA) LLM-based translation systems now achieve competitive results with significantly fewer annotations by leveraging supervised fine-tuning guided by sophisticated instructions (see Appendix A for details). Notably, Zhu et al. (2024b) demonstrated that LLMs exhibit remarkable zero-shot and few-shot machine translation capabilities, even without explicit instruction formatting or exemplars. This highlights the intrinsic potential of LLMs for self-improvement beyond finely calibrated supervision, opening new avenues for resource-efficient translation paradigms.

2.2 Monte-Carlo Tree Search

Monte-Carlo Tree Search (MCTS, Browne et al., 2012; Świechowski et al., 2023) is a heuristic search algorithm proficient for complex decision processes such as the game of Go (Silver et al., 2016). MCTS proceeds through four steps: selection, expansion, simulation, and backpropagation.

- **Selection.** MCTS descends the tree from its root by selecting the top amongst child nodes

by their upper confidence bounds (UCB):

$$\text{UCB}(\alpha) = \nu(\alpha) + 2\sqrt{\frac{\log N(A)}{1 + N(\alpha)}}, \quad (1)$$

where node α is a child of node A , and $N(*)$ is the node’s visit count, with $\nu(*)$ as its current utility. The utility is the cumulated rewards r averaged by the visit count:

$$\nu(\alpha) = \frac{\sum r(\cdot)}{N(\alpha)} \quad (2)$$

The UCB balances the exploration and exploitation in the heuristic search.

- **Expansion.** Upon reaching a selected node, MCTS expands the tree by adding a new child node representing a possible move from the current state.
- **Simulation.** From the newly expanded node, a random simulation is performed until either the decision process concludes or reaches a maximum step limit, estimating the reward r for this decision path.
- **Backpropagation.** The obtained reward is propagated backward, updating the utility values of all nodes along the path from the expanded node to the root. This update process refines the UCB values, progressively enhancing the efficiency of subsequent search iterations.

2.3 Self-Optimization in LLM

Typically, optimizations for LLM rely on external rewards (e.g., critic and revise modules) for tuning (Huang et al., 2023; Tian et al., 2024; Zhang et al., 2024). Recent work explores self-optimization using LLMs’ internal knowledge, leveraging performance gaps across test scenarios. For example, in multilingual tasks, higher-performing languages can optimize lesser ones (Geng et al., 2024). She et al. (2024) use a strong language as a pivot to map and optimize weaker languages. Self-play preference optimization (Chen et al., 2024, SPPO) adopts the gaming theory where the post-update models shall prevail by a win rate to optimize the preference. SPPO suits the human preferences’ non-transitive, unstable nature and enables iterative self-play optimization.

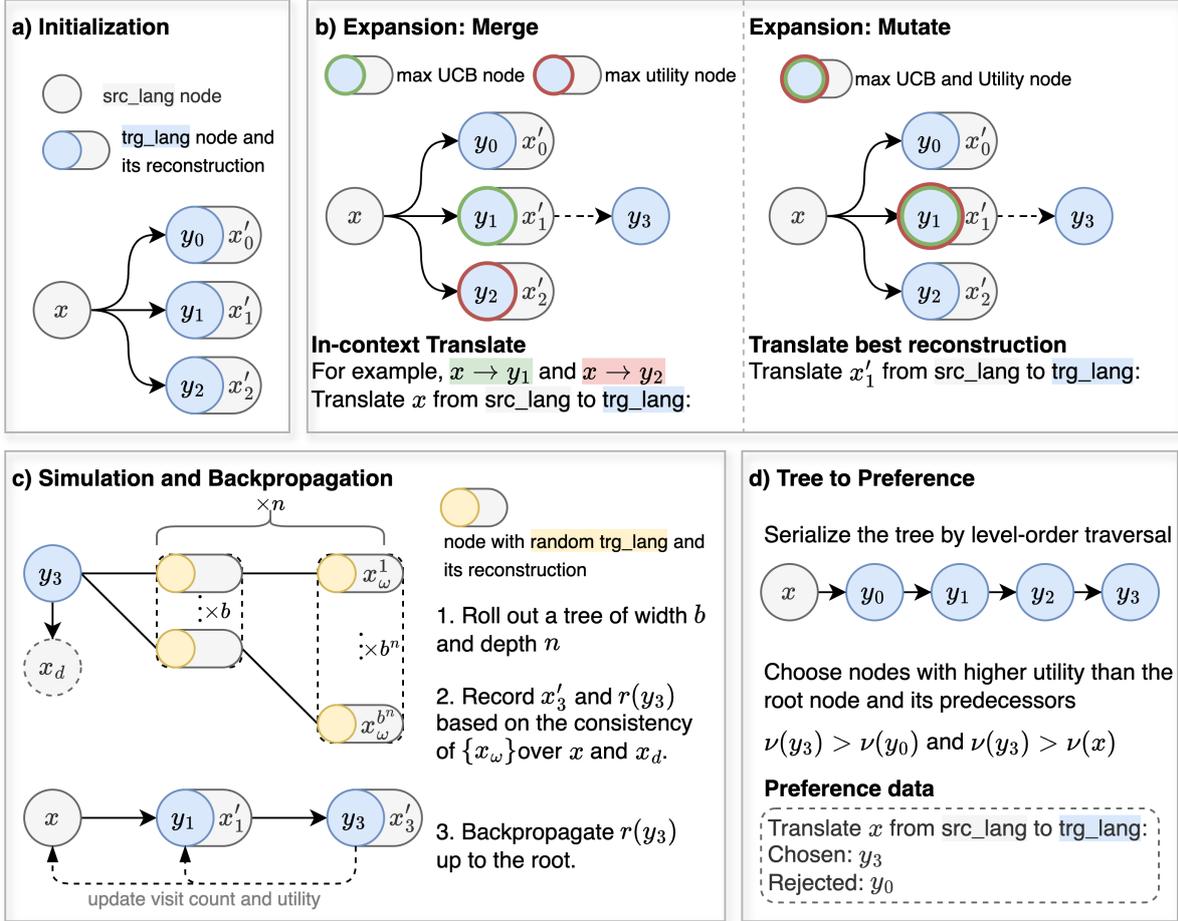


Figure 1: **Overview of TRANS-ZERO.** Once the tree is initiated, the search cycles the selection, expansion, simulation, and back-propagation for new nodes. b) G-MCTS selects the node with maximum UCB to expand a new child node. There are two types of genetic expansion: merge and mutate. c) A mass roll-out simulation of MTP trajectories assesses the semantic consistency. The assessed reward is backpropagated to guide the search. d) Finally, we harvest the search tree into data pairs for preference optimization.

3 TRANS-ZERO

In this section, we present TRANS-ZERO. First, we introduce the Multilingual Translation Process (MTP), a novel framework orchestrating multi-step translation across multiple languages (§3.1). Second, we implement Genetic Monte Carlo Tree Search (G-MCTS) upon MTP to explore promising translations (§3.2), which derive preference from cross-lingual semantic consistency in the search purely through translation prompts, eliminating the need for explicit reasoning training or reward learning. Finally, we utilize the search results for further preference optimization (§3.3), enabling the unsupervised MT training given only monolingual data.

3.1 Multilingual Translation Process

We define the Multilingual Translation Process (MTP) as an iterative translation involving at least two languages, denoted as $\{L_i\}_{|\{L_i\}|>1}$, where each translation step maps the sentence from

one language to another distinct language within the set. An MTP trajectory of length T starts from a source language l :

$$l_T = f(l|l_1, l_2, \dots, l_{T-1}),$$

where the l_i is a sentence in one language in $\{L_i\}$, and $f(\cdot)$ is the translation function. MTP iteratively scales up the translation across languages, enabling similar cross-lingual preferences in work by Geng et al. (2024); She et al. (2024).

Notably, an optimized translation ensures that *the semantics are maintained throughout the multilingual translations*. Such consistency preference intuitively guides the search toward better translation candidates, which can also contribute to preference optimization. E.g., translation optimized via back-translation promotes the bilingual semantic consistency given an MTP $x \rightarrow y \rightarrow x'$.

3.2 Genetic Monte-Carlo Tree Search

As shown in Figure 1, we conduct a genetic Monte-Carlo tree search based on the defined MTP. The search employs genetic expansion coupled with semantic consistency simulation, constructing a search tree in which each node corresponds to a translation candidate in the target language.

Initialization The input language is X with input text x , and the output language is Y . We initialize the search tree with x as the root and perform top-k sampling with a width of b to generate b translation candidates $\{y_i\}_b$ in language Y . Each candidate is assigned as a child node of the root. These nodes are quickly initialized by back-translation to language X , with the corresponding reconstruction recorded, denoted as x' . We define the consistency function $S(a, b)$ over sentence pair (a, b) of **same** language based on a mutual evaluation metric $M(a, b)$, e.g. BLEURT (Sellam et al., 2020). The consistency score is computed as:

$$S(a, b) = \frac{M(a, b) + M(b, a)}{2}, \quad (3)$$

where x' is to compute fast-initiated reward $r(y_i) = S(x, x')$, followed by backpropagations.

Genetic Expansion Given an initialized search tree, each expansion step of the MTCS selects the node with the highest UCB value to generate a new translation in the target language. However, a straightforward generation does not naturally lead to diverse exploration. Inspired by genetic algorithms (Sastry et al., 2005), we propose two strategies for expansion based on the status of the current maximum UCB node:

- **Merge.** We perform a merging when the current maximum UCB node differs from the maximum utility node: Merging is a few-shot translation given the current best translation (i.e., the maximum utility node) as demonstrations:

$$\begin{aligned} y_t &= f(x \mid y_{\text{UCB}}, y_\nu) \\ y_{\text{UCB}} &= \arg \max_{y < t} \{\text{UCB}(y)\} \\ y_\nu &= \arg \max_{y < t} \{\nu(y)\} \end{aligned}$$

Specifically, (x, y_ν) and (x, y_{UCB}) , the pairs from both the maximum utility node and the maximum UCB node, are prepended to the instruction. The LLM then translates the **original input** x to the target language, with the given context.

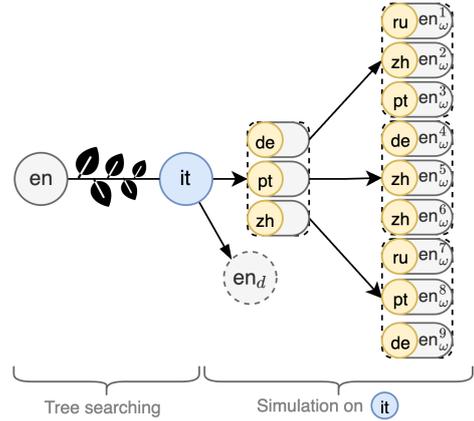


Figure 2: **An example of simulation an English-to-Italian translation candidate using $b = 3$ and $n = 2$.** Through roll-outs of MTP, the Italian candidate (it) is assessed by semantics consistency of b^n English reconstructions $\{\text{en}_\omega^1, \dots, \text{en}_\omega^9\}$ from simulated trajectories.

- **Mutate.** We perform a mutation when the current maximum UCB node is the same as the maximum utility node. Mutation enables creative exploration based on existing translations, which is performed by translating a variant of the original input:

$$y_t = f(x' \mid \arg \max_{y < t} \{\text{UCB}(y)\}).$$

Specifically, the LLMs translate the **best** reconstruction $x' \in \{x_\omega\}$ recorded during simulations of the parent node, instead of the original input x :

Each candidate¹ is generated autoregressively given existing translations. The merge operation reduces the most promising MTP trajectories via a contexted translation of the original x . Meanwhile, the mutation extends the existing trajectory through additional reconstructed x' from simulation, which may come from languages other than source and target (§ 3.2 on simulations). Thus a search tree is constructed, with any path from the root to a node in the tree a valid MTP trajectory.

Simulation with Multilingual Semantic Consistency The simulation generates reward signals through semantic agreement across translations, which directly updates the UCB values in tree search to assess translation preferences without external supervision. The key insight is that a better translation candidate achieves higher semantic consistency for themselves and their descendants, thus guiding an optimized search.

¹The translation instructions within G-MCTS are sampled from Table 4 in Appendix A.

Utility (ν)	Tree Nodes
0.5255	If you always harm others, the chickens gonna come home to roost.
0.5595	① 如果你总是伤害别人, 那么必然会有一天会有报应。(If you always hurt others, there will be retribution one day.)
1.2290	② 你总是伤害别人, 最后伤害的就是自己。(You always hurt others, and in the end you end up hurting yourself.)
0.5801	③ 你总是伤害别人, 最后伤害的还是你自己。(You always hurt others, but in the end you end up hurting yourself.)
2.3230	④ 你总是伤害别人, 最后总会有报应的。(You always hurt others, and you will get your comeuppance in the end.)
0.6117	⑤ 你若总是伤害别人, 鸡就要飞回你的巢了。(If you keep hurting others, the chickens will fly back to your nest.)
0.5275	⑥ 如果你总是伤害别人, 最后鸡蛋就要落回自己的头上。(If you always hurt others, eventually the eggs will fall back on your head.)
0.5626	⑦ 你总是伤害别人, 最后鸡蛋就要落回自己的头上。(If you always hurt others, in the end the eggs will fall back on your head.)
0.5657	⑧ 如果你总是伤害别人, 那么你也会遭到报应。(If you always hurt others, then you will also suffer retribution.)
0.5601	⑨ 如果你总是伤害别人, 那么最后总会有报应的。(If you keep hurting others, you will get punished in the end.)
0.4676	⑩ 如果你总是伤害别人, 那么麻雀总会飞回窝的。(If you always hurt others, the sparrow will always fly back to the nest.)

Extracted Preference Pairs for Self-Play Preference Optimization (SPPO)

Chosen	Rejected	Win rates (softmax)
② 你总是伤害别人, 最后伤害的就是自己。	① 如果你总是伤害别人, 那么必然会有一天会有报应。	1.23 : 0.56 (0.6614)
④ 你总是伤害别人, 最后总会有报应的。	② 你总是伤害别人, 最后伤害的就是自己。	2.32 : 1.23 (0.7491)
④ 你总是伤害别人, 最后总会有报应的。	③ 你总是伤害别人, 最后伤害的还是你自己。	2.32 : 0.58 (0.8511)
⑦ 你总是伤害别人, 最后鸡蛋就要落回自己的头上。	⑥ 如果你总是伤害别人, 最后鸡蛋就要落回自己的头上。	0.56 : 0.52 (0.5088)
⑧ 如果你总是伤害别人, 那么你也会遭到报应。	⑦ 你总是伤害别人, 最后鸡蛋就要落回自己的头上。	0.57 : 0.56 (0.5008)

Table 1: **Example of Tree-to-Preference.** We perform level-order traversal of a search tree for English-to-Chinese translation. The first line presents the source input as the root. The Chinese translations are accompanied by corresponding English explanations enclosed in parentheses. Given that the utility near the root shall be larger, we apply a sorting algorithm to arrange the sequence in descending order, where each swap during the sort makes a preference pair.

As shown in Figure 2, assessing candidate y involves rolling out a temporary sub-tree on y through MTP with a width of b until it reaches a maximum depth of n . Each rollout step translates the parent node into a different language sampled from the supported languages $\{L_i\}$, with corresponding input reconstruction x_ω . Ultimately, this results in b^n MTP trajectories with reconstructions $\{x_\omega\}_{b^n}$. Given Eq. 3, the semantic consistency of these reconstructions $\{x_\omega\}_{b^n}$ with the original input x can be calculated as follows:

$$r(y) = \max(\underbrace{S(x_\omega, x)}_{\text{literal}}, \underbrace{S(x_\omega, x_d)}_{\text{free}}), x_\omega \in \{x_\omega\}_{b^n}$$

where x is the original input, and x_d is a direct reconstruction of the candidate y . Translation can be literal or free, with free translations assessed via straightforward back-translation x_d . Each assessment is *averaged* over all trajectories' $\{x_\omega\}_{b^n}$, with the superior one as the reward $r(y)$ derived from the simulation. The best simulation is recorded as x' for further expansions and simulations:

$$x' = \arg \max_{\{x_\omega\}_{b^n}} (S(x_\omega, x), S(x_\omega, x_d))$$

During the backpropagation, the reward $r(y)$ updates utility ν and visit counts of all nodes on the trajectory from the current node back to the root according to Eq. 2.

3.3 Tree-to-Preference Algorithm

As Table 1 shows, once the G-MCTS is finished, we extract preference data from the translation can-

Algorithm 1 Tree-to-Preference Algorithm

Require: Translation candidates $\{y\}_*$ with utilities $\{\nu\}_*$, search tree \mathcal{T}

Ensure: Preference pairs with win rates for SPPO

1: **Step 1: Serialize and Sort Tree**

2: $\mathcal{S} \leftarrow \text{LevelOrderTraversal}(\mathcal{T}) \triangleright$ Serialize tree and merge duplicates

3: $\mathcal{S} \leftarrow \text{SelectionSort}(\mathcal{S}, \text{descending}) \triangleright$ Sort nodes by utility

4: **Step 2: Generate Preference Pairs**

5: $\mathcal{P} \leftarrow \emptyset$

6: **for** each swap (y_i, y_j) in \mathcal{S} **do**

7: **if** $\nu_i > \nu_{\text{root}}$ **and** $\nu_i > \nu_j$ **then**

8: win rate $\leftarrow \frac{\exp(\nu_i)}{\exp(\nu_i) + \exp(\nu_j)}$

9: $\mathcal{P} \leftarrow \mathcal{P} \cup \{(y_i, y_j, \text{win rate})\}$

10: **end if**

11: **end for**

12: **Return** $\mathcal{P} \triangleright$ Preference pairs with win rates for SPPO

didates $\{y\}_*$ based on their utility $\{\nu\}_*$ by Algorithm 1. The tree is serialized where duplicate nodes are merged to their ancestor, with utilities and visit counts accumulated. Intuitively, nodes away from the root take more MTP steps to generate and thus face more risk of semantic loss as the translation step increases. Consequently, if a node has a higher utility than its ancestors or brothers, the corresponding translation is preferred for optimization. Furthermore, the root node tracks a comprehensive utility, which reflects the expected semantic consistency of the entire search. Thus,

the utility of the preferred node shall also be higher than that of the root node. The preference pairs are utilized in SPPO, where their utilities are transformed into win rates through a softmax function. Note that translations that failed language detection may be generated during MTP, and their utility is halved as a penalty during the sorting and filtering.

4 Experiments

4.1 Settings

Data. We conduct experiments across six widely used languages: English (EN), German (DE), Portuguese (PT), Italian (IT), Chinese (ZH), and Russian (RU). We utilize the latest monolingual data from the WMT datasets. The SFT data for baselines is generated from the combination of the Flores-200 development set, with equal size for all translation directions. To evaluate multilingual translation performance, we employ the Flores-200 benchmark (Costa-jussà et al., 2022), assessing three key translation directions: (1) EN \Rightarrow X (English to other languages), (2) X \Rightarrow EN (other languages to English), and (3) X \Rightarrow X (inter-translations between non-English languages). These evaluations cover all translation directions of the six languages mentioned above.

Metrics. We evaluate translation quality with the reference-oriented metric BLEURT (Sellam et al., 2020) and reference-free metric COMET-KIWI (KIWI, Rei et al., 2022).

Baselines. We compare TRANS-ZERO with the following representative baselines:

- **General Instruct LLMs:** LLMs with off-the-shelf instruct-following ability for MT, e.g., Mixtral-8 \times 7B, Llama3.1-8B-instruct and Qwen2.5-7B-instruct.
- **MT-oriented LLMs:** LLMs supervised by MT annotations, e.g., ALMA (Xu et al., 2023) with parallel annotations, ALMA-R (Xu et al., 2024b) with preference annotations and Tower-Instruct (Colombo et al., 2024) with multi-task MT-related annotations, representing strong supervised baselines.
- **Base model and SFT model:** The base LLM for TRANS-ZERO, and their supervised counterparts.

Implementation. The training is conducted on 32 NVIDIA A100 GPUs (80GB). We utilize two base LLMs without instruction tuning: Llama-3.1-8b and Qwen-2.5-7b. Since some

base LLMs (e.g., Llama-3.1-Base) lack the initial instruction-following for translation, we cold-start the LLM with a snippet of translation instructions. The TRANS-ZERO is parallelized across 32 threads, with each thread a batch of 10 sentences assigned to random translation directions for G-MCTS. Upon completion of the search, we reduce all search threads for filtered preference data and apply Self-Play Preference Optimization (Chen et al., 2024, SPPO). Additional details are in Appendix B.

4.2 Main Results

As shown in Table 2, TRANS-ZERO based on Llama3.1 and Qwen2.5 achieves performance comparable to MT baselines trained on large-scale annotations, despite using only monolingual data for self-play. While TRANS-ZERO matches the English translation performance (EN \Rightarrow X) of ALMA-R, which also utilizes preference optimization, it significantly surpasses ALMA-R in non-English translation directions. Compared to Tower-instruct, an LLM trained on large-scale annotations, TRANS-ZERO exhibits slightly lower performance but remains highly competitive.

We also evaluated the translation performance of the base model and its instruction fine-tuned version for comparison. Through exploration learning, TRANS-ZERO significantly enhances the performance of both base models.

Additionally, we include SFT baselines using 5m parallel data by pairing the Flores200 development set, and those 5k instructions as parallel data. Although TRANS-ZERO does not match the performance by 5m supervision on EN \Rightarrow X and X \Rightarrow EN, it shows significant improvements in the non-English translation direction (X \Rightarrow X), achieving performance comparable to 5m supervision.

We further compared the performance improvement to varying scales of parallel annotations in SFT. As shown in Figure 3, the translation quality saturates after more than 100k samples, especially for the EN \Rightarrow X and X \Rightarrow X directions. This suggests that simply increasing the amount of parallel annotations may not lead to proportional translation improvements.

Increasing the number of languages expands the potential exploration thus improving the TRANS-ZERO’s performance upper bound. We explored with 4 and 6 languages for the G-MCTS. Figure 4 illustrates the performance changes in German-Chinese translation during Llama3.1-8b training.

	EN \Rightarrow X		X \Rightarrow EN		X \Rightarrow X		Average	
	BLEURT	KIWI	BLEURT	KIWI	BLEURT	KIWI	BLEURT	KIWI
Mixtral-8x7B-Instruct	55.42	69.07	75.41	81.63	54.49	71.64	61.77	74.11
Llama3.1-Instruct	62.57	74.28	72.07	77.90	62.52	76.49	65.72	76.22
Qwen2.5-Instruct	<u>72.16</u>	81.99	77.70	84.60	<u>68.59</u>	79.87	<u>72.82</u>	82.15
ALMA	71.98	82.60	<u>78.25</u>	84.34	61.07	80.89	70.43	82.61
ALMA-R	69.38	<u>83.10</u>	77.52	<u>84.87</u>	51.03	<u>82.47</u>	65.98	<u>83.48</u>
Tower-Instruct	76.74	85.26	78.73	85.02	72.98	83.08	76.15	84.45
Llama3.1-Base	33.18	25.53	50.83	55.64	34.38	53.52	39.46	44.89
w/ SFT (40k)	<u>74.46</u>	82.30	77.30	84.03	71.23	80.02	74.33	82.12
w/ SFT (5m)	75.80	84.61	78.47	84.61	73.30	<u>82.33</u>	75.86	83.85
w/ TRANS-ZERO	73.71	<u>83.20</u>	<u>77.60</u>	<u>84.34</u>	<u>73.28</u>	82.71	<u>74.86</u>	<u>83.42</u>
Qwen2.5-Base	62.91	73.98	70.98	80.50	62.70	77.86	65.53	77.45
w/ SFT (5m)	75.32	84.79	78.21	85.42	72.99	82.92	75.49	84.38
w/ TRANS-ZERO	<u>75.05</u>	<u>84.48</u>	78.21	<u>84.88</u>	<u>72.23</u>	<u>82.30</u>	<u>75.16</u>	<u>83.89</u>

Table 2: **TRANS-ZERO achieves comparable and improved translation compared to SFT baselines with only monolingual self-play.** We highlight the best and the second-best performances in **bold** and underlined, respectively.

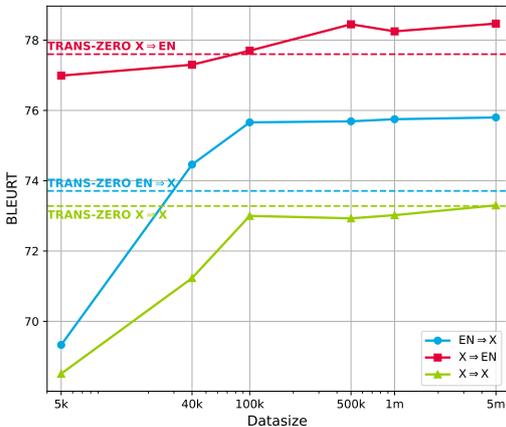


Figure 3: **BLEURT performance for SFT based on the Llama3.1-Base at different data sizes.** We include the performance of TRANS-ZERO in each language direction.

The number of languages used in tree search significantly impacts TRANS-ZERO’s performance upper bound: increasing the number of languages can enhance the overall learning performance. With 6 languages, TRANS-ZERO essentially matches the performance of the open-source baseline system.

4.3 Inference-time Scaling with G-MCTS

Inference time scaling, such as Chain of Thought (CoT) (Wei et al., 2022b), has become popular for improving the performance of LLMs. However, CoT requires additional learning given multiple natural language understanding supervisions. In contrast, G-MCTS enables straightforward inference-time scaling using only translation instructions. During the tree search, merging and mutation continuously explore and integrate rele-

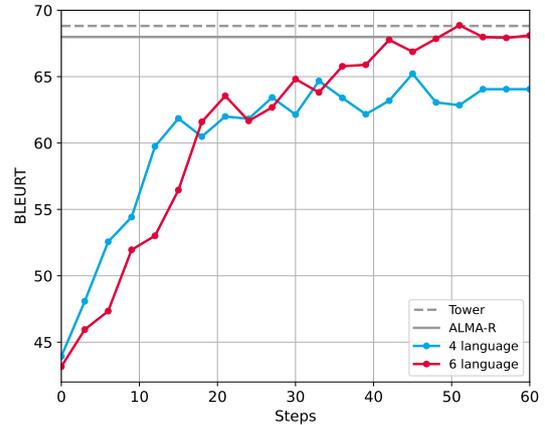


Figure 4: **The learning diagram of TRANS-ZERO on Llama3.1-Base for German-to-Chinese translation demonstrates the search process in 4-language and 6-language settings under G-MCTS.** By incorporating 6 languages, TRANS-ZERO attains BLEURT scores on par with the baseline systems.

vant expressions from the multilingual semantic space, revising translations based on the LLM’s conditional generation capabilities.

We employ G-MCTS with 6 languages to explore translations for the Llama-3.1 baselines, as well as Tower-Instruct and ALMA-R. The candidate of the highest utility makes the final translation. As Table 3 shows, G-MCTS requires a language model with basic instruction-following capabilities. Consequently, the Llama3.1-Base model fails the search due to its lack of instruction-following in various translation directions. Similarly, ALMA-R also fails due to its limited multilingualism, as indicated by its significantly lower X \Rightarrow X performance.

In contrast, Tower-Instruct and Llama3.1-Instruct significantly improve translation perfor-

	EN \Rightarrow X		X \Rightarrow EN		X \Rightarrow X		Average	
	BLEURT	KIWI	BLEURT	KIWI	BLEURT	KIWI	BLEURT	KIWI
ALMA-R	69.38	83.10	77.52	84.87	51.03	82.47	65.98	83.48
+ G-MCTS	Failed	Failed	Failed	Failed	Failed	Failed	Failed	Failed
Tower-Instruct	76.74	85.26	78.73	85.02	72.98	83.08	76.15	84.45
+ G-MCTS	76.44 _{-0.30}	85.33 _{+0.07}	78.28 _{-0.45}	85.12 _{+0.10}	74.42 _{+1.44}	83.57 _{+0.49}	76.38 _{+0.23}	84.67 _{+0.22}
Llama3.1-Base	33.18	25.53	50.83	55.64	34.38	53.52	39.46	44.89
+ G-MCTS	Failed	Failed	Failed	Failed	Failed	Failed	Failed	Failed
Llama3.1-Instruct	62.57	74.28	72.07	77.90	62.52	76.49	65.72	76.22
+ G-MCTS	64.21 _{+1.64}	80.12 _{+5.84}	70.01 _{-2.06}	79.86 _{+1.96}	68.12 _{+5.60}	77.12 _{+0.63}	67.45 _{+1.73}	79.03 _{+2.81}
Llama3.1-SFT (5k)	69.33	80.19	76.99	83.97	68.51	78.38	71.61	80.85
+ G-MCTS	71.55 _{+2.22}	82.23 _{+2.04}	76.89 _{-0.10}	84.00 _{+0.03}	71.92 _{+3.41}	81.23 _{+2.85}	73.45 _{+1.84}	82.49 _{+1.64}
Llama3.1-SFT (5m)	75.80	84.61	78.47	84.61	73.30	82.33	75.86	83.85
+ G-MCTS	76.16 _{+0.36}	84.95 _{+0.34}	78.71 _{+0.24}	84.68 _{+0.07}	73.36 _{+0.06}	82.48 _{+0.15}	76.08 _{+0.22}	84.04 _{+0.19}

Table 3: **G-MCTS enhances translation by scaling up inference, given the models’ own instruction-following capability and multilingualism.** Performance improvements beyond one point are highlighted in bold. The base LLM and ALMA-R exhibit limitations due to their failure to follow instructions in various translation directions. The search particularly enhances X \Rightarrow X translations, where the availability of SFT annotations is significantly limited compared to English-related annotations.

mance in the X \Rightarrow X direction, benefiting from multilingual priors. The base model trained with small-scale supervision also shows notable improvements. However, the improvement upon Llama3.1-SFT (5m) with large-scale supervision, is almost negligible. This suggests that when translations are fully activated, the performance gains from G-MCTS, rooted in LLM’s inherent multilingualism, are not statistically significant.

5 Related Work

The utilization of LLM for machine translation has become popular, aligning with the prevailing trends in LLM applications. Xu et al. (2023) first shifts the machine translation paradigm to fine-tuned LLM with moderate parallel supervision. Though LLM seems more data-efficient, it does not have a big appetite for large-scale supervision due to potential catastrophic forgetting (Xu et al., 2023; Kondo et al., 2024). Directly scaling up multilingual MLE supervision hurts performance on resource-rich languages (Xu et al., 2024a). Therefore, XALMA (Xu et al., 2024a) hand-craft a mixture-of-expert to route different translation directions through separated modules, while ALMA-R (Xu et al., 2024b) turn to scale up more expensive preference tuning.

Preference tuning offers flexibility when fitting LLM with subjective human expectations for open-ended generations. However, the expense of preference annotation has led researchers to seek more cost-effective data sources, e.g., with additional assessments such as critic and revise modules (Huang et al., 2023; Tian et al., 2024; Zhang et al., 2024)

Intuitively, the cross-lingual gaps in LLM offer a more scalable self-improvement preference as the proficient languages improve the lesser ones (Geng et al., 2024), e.g., a straightforward improvement roots in the direct mapping from the dominant linguistic ability as preference (She et al., 2024). Researchers further scale the preference by iterative and competitive gaming theory (Chen et al., 2024, SPPO), making it possible for models to self-improve. Recent work by Deepseek (Guo et al., 2025) has empirically validated its self-improving potential by employing large-scale reinforcement learning with its multilingual reasoning abilities.

6 Conclusion

In this work, we present TRANS-ZERO, a novel framework for multilingual machine translation that leverages multilingual LLMs with monolingual data only. Our experiments demonstrate that the proposed Genetic Monte-Carlo Tree Search (G-MCTS) effectively enhances translation quality by exploiting the LLM’s inherent multilingual and instruction-following capabilities. Furthermore, we show that iterative training of G-MCTS, combined with preference optimization using monolingual data, TRANS-ZERO achieves scalable performance improvements, with the number of supported languages positively correlating with final translation quality. These findings establish a new direction for resource-efficient MT by shifting the paradigm from supervised parallel data to self-supervised monolingual learning.

518 **Limitations**

519 While TRANS-ZERO demonstrates promising re-
520 sults, it has several limitations that warrant discus-
521 sion. First, the framework introduces higher compu-
522 tational overhead than supervised baselines, as
523 the search process requires extensive exploration
524 of the cross-lingual semantic space. Second, its
525 effectiveness is inherently tied to the multilingual
526 capabilities of the underlying LLM, rendering it
527 less suitable for weaker models with limited cross-
528 lingual alignment. Third, due to computational
529 constraints, our experiments were limited in scale:
530 we were unable to explore larger LLMs or extend
531 the search process to a broader range of languages.
532 Finally, the framework’s performance upper bound
533 may be influenced by the quality and diversity of
534 monolingual data used for training, highlighting
535 the need for future research into identifying the
536 most cost-effective data types for self-supervised
537 training.

538 **Ethics Statement**

539 The authors declare no competing interests. The
540 datasets used in the training and evaluation come
541 from publicly available sources and do not contain
542 sensitive content such as personal information.

543 **References**

544 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
545 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
546 Diogo Almeida, Janko Altenschmidt, Sam Altman,
547 Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#).
548 *ArXiv preprint*, abs/2303.08774.

549 Cameron B Browne, Edward Powley, Daniel White-
550 house, Simon M Lucas, Peter I Cowling, Philipp
551 Rohlfshagen, Stephen Tavener, Diego Perez, Spyri-
552 don Samothrakis, and Simon Colton. 2012. A survey
553 of monte carlo tree search methods. *IEEE Transac-
554 tions on Computational Intelligence and AI in games*,
555 4(1):1–43.

556 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji,
557 and Quanquan Gu. 2024. [Self-play fine-tuning con-
558 verts weak language models to strong language mod-
559 els](#). *ArXiv preprint*, abs/2401.01335.

560 Pierre Colombo, Duarte Alves, José Pombal, Nuno
561 Guerreiro, Pedro Martins, Joao Alves, Amin Fara-
562 jian, Ben Peters, Ricardo Rei, Patrick Fernandes, et al.
563 2024. [Tower: An open multilingual large language
564 model for translation-related tasks](#). *ArXiv preprint*,
565 abs/2402.17733.

566 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha
567 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*,
abs/2207.04672. 568 569 570 571

Xiang Geng, Ming Zhu, Jiahuan Li, Zhejian Lai, Wei
Zou, Shuaijie She, Jiaxin Guo, Xiaofeng Zhao,
Yinglu Li, Yuang Li, et al. 2024. [Why not transform
chat large language models to non-english?](#) *ArXiv
preprint*, abs/2405.13923. 572 573 574 575 576

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
Peiyi Wang, Xiao Bi, et al. 2025. [Deepseek-r1: In-
centivizing reasoning capability in llms via reinforce-
ment learning](#). *ArXiv preprint*, abs/2501.12948. 577 578 579 580 581

Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi
Wang, Hongkun Yu, and Jiawei Han. 2023. [Large
language models can self-improve](#). In *Proceedings
of the 2023 Conference on Empirical Methods in Natu-
ral Language Processing*, pages 1051–1068, Singa-
pore. Association for Computational Linguistics. 582 583 584 585 586 587

Minato Kondo, Takehito Utsuro, and Masaaki Nagata.
2024. [Enhancing translation accuracy of large lan-
guage models through continual pre-training on par-
allel data](#). *ArXiv preprint*, abs/2407.03145. 588 589 590 591

Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su,
John D Co-Reyes, Avi Singh, Kate Baumli, Shariq
Iqbal, Colton Bishop, Rebecca Roelofs, et al. 2024.
[Training language models to self-correct via rein-
forcement learning](#). *ArXiv preprint*, abs/2409.12917. 592 593 594 595 596

Jiahuan Li, Shanbo Cheng, Shujian Huang, and Jia-
jun Chen. 2024. [MT-PATCHER: Selective and ex-
tendable knowledge distillation from large language
models for machine translation](#). In *Proceedings of
the 2024 Conference of the North American Chap-
ter of the Association for Computational Linguistics:
Human Language Technologies (Volume 1: Long
Papers)*, pages 6445–6459, Mexico City, Mexico. As-
sociation for Computational Linguistics. 597 598 599 600 601 602 603 604 605

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,
Chrysoula Zerva, Ana C Farinha, Christine Maroti,
José G. C. de Souza, Taisiya Glushkova, Duarte
Alves, Luisa Coheur, Alon Lavie, and André F. T.
Martins. 2022. [CometKiwi: IST-unbabel 2022 sub-
mission for the quality estimation shared task](#). In
*Proceedings of the Seventh Conference on Machine
Translation (WMT)*, pages 634–645, Abu Dhabi,
United Arab Emirates (Hybrid). Association for Com-
putational Linguistics. 606 607 608 609 610 611 612 613 614 615

Kumara Sastry, David Goldberg, and Graham Kendall.
2005. Genetic algorithms. *Search methodologies:
Introductory tutorials in optimization and decision
support techniques*, pages 97–125. 616 617 618 619

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020.
[BLEURT: Learning robust metrics for text genera-
tion](#). In *Proceedings of the 58th Annual Meeting of
the Association for Computational Linguistics*, pages
7881–7892, Online. Association for Computational
Linguistics. 620 621 622 623 624 625

626	Shuaijie She, Wei Zou, Shujian Huang, Wenhao Zhu, Xiang Liu, Xiang Geng, and Jiajun Chen. 2024. Mapo: Advancing multilingual reasoning through multilingual alignment-as-preference optimization . <i>ArXiv preprint</i> , abs/2401.06838.	684
627		685
628		686
629		687
630		
631	David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneshelvam, Marc Lanctot, et al. 2016. Mastering the game of go with deep neural networks and tree search. <i>nature</i> , 529(7587):484–489.	688
632		689
633		690
634		691
635		692
636		693
637		694
638		695
639		696
640		
641		
642	Maciej Świechowski, Konrad Godlewski, Bartosz Sawicki, and Jacek Mańdziuk. 2023. Monte carlo tree search: A review of recent modifications and applications. <i>Artificial Intelligence Review</i> , 56(3):2497–2562.	697
643		698
644		699
645		700
646		701
647		702
648		703
649		704
650		
651		
652		
653	Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. 2024. Toward self-improvement of llms via imagination, searching, and criticizing . <i>ArXiv preprint</i> , abs/2404.12253.	705
654		706
655		707
656		708
657		709
658		710
659		711
660		712
661		713
662		714
663		715
664		716
665		717
666		718
667		719
668		720
669		721
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		

Model	Translation Instruction
ALMA & ALMA-R	Translate this from {src_lang} to {trg_lang}: \n{src_lang}: {src_sent}\n{trg_lang}:
Tower-Instruct	Translate the following text from {src_lang} into {trg_lang}.\n{src_lang}: {src_sent} \n{trg_lang}:
Others	Please translate the {src_lang} into {trg_lang}: {src_sent} {src_lang}: {src_sent} = {trg_lang}: {src_sent} in {src_lang} can be translated to {trg_lang} as: {src_lang}: {src_sent} \n {trg_lang}: Explain the following {src_lang} sentence in {trg_lang}: {src_sent}

Table 4: Commonly adopted translation instructions for LLM, {src_lang} and {trg_lang} indicates the corresponding languages for source and target, and {src_sent} presents the input sentence of the source language.

sampling. Direct inference with G-MCTS searches
by the width of $b = 10$ given 6 languages.

736
737