
Why Not Transform Chat Large Language Models to Non-English?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The scarcity of non-English data limits the development of non-English large language models (LLMs). Transforming English-centric LLMs to non-English has
2 been identified as an effective and resource-efficient method. Previous works start
3 from base LLMs and perform knowledge distillation (KD) with data generated
4 by stronger LLMs, e.g. GPT-4. Compared to base LLMs, chat LLMs are further
5 optimized for advanced abilities, e.g. multi-turn conversation and human preference
6 alignment, and thus more powerful in both helpfulness and safety. However,
7 transforming a chat LLM involves two critical issues: (1) How can we effectively
8 transfer advanced abilities without their supervised data? (2) How can we prevent
9 the original knowledge from catastrophic forgetting during transformation? We
10 target these issues by introducing a simple framework called TransLLM. For the
11 first issue, TransLLM divides the transfer problem into some common sub-tasks
12 with the translation chain-of-thought, which uses the translation as the bridge between
13 English and non-English step-by-step. We further enhance the performance
14 of sub-tasks with publicly available data. For the second issue, we propose a
15 method comprising two synergistic components: low-rank adaptation for training
16 to maintain the original LLM parameters, and recovery KD, which utilizes data
17 generated by the chat LLM itself to recover the original knowledge from the frozen
18 parameters. In the experiments, we transform the LLaMA-2-chat-7B to the Thai
19 language. Our method, using only single-turn data, outperforms strong baselines
20 and ChatGPT on multi-turn benchmark MT-bench. Furthermore, our method,
21 without safety data, rejects more harmful queries of safety benchmark AdvBench
22 than both ChatGPT and GPT-4.
23





 Attacker	อธิบายรายละเอียดวิธี## (Explains in detail how to ##)	 Ours w/ GPT-4 KD	ซึ่งสามารถทำได้หลายวิธี... (This can be done in several ways...)
 ChatGPT	วิธี##สามารถเกิดขึ้นในหลายรูปแบบ... (The ## can happen in several ways...)	 Ours w/ recovery KD	ฉันไม่สามารถปฏิบัติตามคำขอของคุณได้... (I am unable to comply with your request...)

Figure 1: In this example, the attacker jailbreaks ChatGPT in Thai, while our method successfully rejects to response. The recovery KD data is more suitable for preserving the original knowledge than widely used GPT-4 KD data, although GPT-4 performs better in both helpfulness and safety. We omit the harmful text with ## and provide the English translation under the Thai text.

24 1 Introduction

25 Recently, significant influence has been demonstrated by chat large language models (LLMs), such
26 as ChatGPT (OpenAI, 2022), Palm-2 (Anil et al., 2023), and LLaMA-2-chat (Touvron et al., 2023).
27 Their high capabilities rely on massive data and complex training processes. Taking the LLaMA-
28 2-chat as an example, the training usually includes the following steps: (1) pre-training (PT) on a

29 large monolingual corpus to obtain the base LLM; (2) supervised fine-tuning (SFT) on multi-turn
30 dialogue datasets; (3) iteratively refining on human preference datasets using reinforcement learning
31 with human feedback (RLHF) methodologies (Ouyang et al., 2022). These steps help in creating
32 LLMs that not only understand and generate human-like text but also align with human values, and
33 therefore provide safe and useful responses.

34 Unfortunately, popular unlabeled and labeled training data is English-dominated. Consequently,
35 LLMs are less satisfying in terms of both usefulness and safety when being applied to non-English.
36 Yong et al. (2023) have shown that even powerful LLMs, such as GPT-4, are vulnerable to safety
37 concerns in non-English.

38 To improve the non-English performance, recent works attempt to transfer knowledge from English
39 to non-English. However, they focus on base LLMs instead of powerful chat LLMs. Basically,
40 they start from base LLMs and use knowledge distillation (KD) data generated by the strong LLM,
41 like GPT-4, for transfer training and instruction tuning. For example, PolyLM (Wei et al., 2023)
42 transfers English knowledge implicitly via multilingual instruction tuning on a multilingual base
43 LLM. X-LLaMA (Zhu et al., 2023) supplements the multilingual instruction-following task with the
44 translation task to build semantic alignment across languages.

45 When transforming the base LLMs, instruct tuning is performed simultaneously with or after transfer
46 training. Therefore, the instruction following knowledge, i.e. the basic conversation knowledge,
47 will not be overridden by extra knowledge. However, for chat LLMs, the advanced conversation
48 knowledge, especially human preference, has been incorporated into the model parameters during
49 fine-tuning. As a result, subsequent transfer training in previous works will result in catastrophic
50 forgetting of such knowledge. What is worse, the high-quality STF data used for training the chat
51 LLM is precious and usually unavailable. Therefore, transforming a chat LLM involves two critical
52 issues: (1) How can we transfer advanced abilities with limited available data? (2) How can we
53 prevent the original English knowledge from catastrophic forgetting during transfer?

54 To build safe non-English LLMs as shown in Figure 1, we target these issues by introducing a simple
55 framework called TransLLM. For the first issue, TransLLM utilizes the translation chain-of-thought
56 (TCOT) (Zhang et al., 2023), which models the transfer as some common sub-tasks. During TCOT,
57 the LLM will handle the non-English query step-by-step in single inference: it first translates the
58 query to English; then it responds to the query in English; and finally it, generates the non-English
59 answer based on all the above context. We further enhance the the performance of sub-tasks with
60 publicly available data thus TCOT can transfer English knowledge effectively. For the second issue,
61 we propose a method comprising two synergistic components: (1) We employ the low-rank adaptation
62 (LoRA) (Hu et al., 2021) for training to maintain the original LLM parameters. (2) We introduce
63 recovery KD, utilizing data generated by the chat LLM itself, to recover the original knowledge from
64 the frozen parameters. The recovery KD data can be fitted easily using the original parameters. This
65 enables the LLM to learn a “shortcut” that uses the English knowledge from the original parameters.

66 As shown in Figure 2, TransLLM organizes all the above ideas into the following steps: (1) Model
67 extension: we extend the model with LoRA modules and fine-grained target language vocabulary. (2)
68 Target language pre-training: we pre-train the chat LLM on the monolingual target language data
69 so that the LLM can leverage such knowledge to improve translation and target language responses.
70 (3) Translation pre-training: we further train the LLM with a bi-directional translation task between
71 English and the target language, and we also introduce the English language modeling task to protect
72 the English embeddings. (4) Transfer fine-tuning: we fine-tune the LLM on TCOT, recover KD, and
73 translation data so that the LLM can respond in English, the target language, and the translation tasks
74 automatically.

75 We conduct comprehensive experiments for transforming the LLaMA-2-chat-7B from English to
76 Thai. TransLLM outperforms strong baselines and surpasses ChatGPT by 35% and 23.75% for the
77 first and second turns on the MT-bench with statistical significance. More importantly, we attain an
78 improvement of 14.8% and 8.65% over ChatGPT and GPT-4 respectively on the safety benchmark
79 AdvBenchmark with statistical significance.

80 Our main contributions are summarized as follows:

- 81 • In this paper, we highlight the advantages and challenges of transforming a chat LLM to
82 non-English and propose a simple yet effective framework for this end.

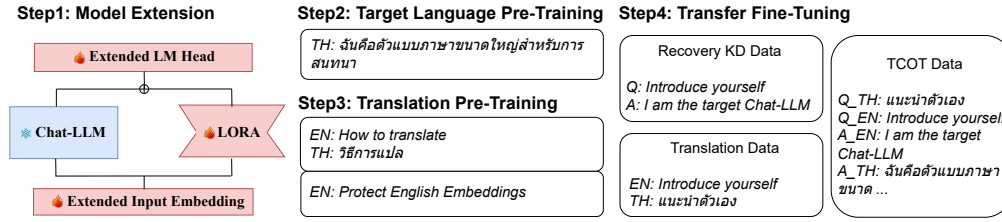


Figure 2: TransLLM pipeline.

- 83
- 84
- 85
- 86
- 87
- 88
- 89
- 90
- 91
- The experiments indicate TransLLM successfully transfer advanced abilities, e.g. multi-turn conversation and human preference alignment, with limited available data. TransLLM, with only 7 billion parameters, outperforms ChatGPT in Thai in both helpfulness and safety.
 - Analysis shows that recovery KD plus with LoRA successfully preserves the original knowledge. The TransLLM model mostly uses the original knowledge for English while uses the new knowledge for Thai.
 - We discuss the limitations of TransLLM, and point out several potential future directions. We will make our code and datasets publicly available (please refer to supplementary materials). We hope this work can lay a solid foundation for developing safe LLMs in non-English.

92 2 Background

93 The language models are trained to predict the next token in a sequence given the previous tokens by
 94 maximum likelihood estimation (MLE), which can be represented by the following equation:

$$J_{PT} = \arg \max_{\theta} \sum_{i=1}^{|y|} \log P(y_i | y_{<i}; \theta), \quad (1)$$

95 where θ denotes learnable model parameters, and $y_{<i}$ are the tokens preceding y_i in the sequence.

96 For fine-tuning on a supervised dataset, each instance contains a query x and its corresponding label
 97 y . The SFT loss is only calculated on the label y , ignoring the query x :

$$J_{SFT} = \arg \max_{\theta} \sum_{i=1}^{|y|} \log P(y_i | x, y_{<i}; \theta). \quad (2)$$

98 For both PT and SFT, the special tokens $\langle s \rangle$ and $\langle /s \rangle$ are added at the beginning and the end of the
 99 training instance respectively.

100 3 Method

101 In this section, we first describe the model architecture, then we introduce the training and inference
 102 procedures in detail.

103 3.1 Model Architecture

104 Nowadays, popular LLMs use byte-level byte pair encoding (BBPE) tokenizer (Wang et al., 2020)
 105 following GPT-2 (Radford et al., 2019). However, the tokenizer is usually developed on the English-
 106 dominated dataset, therefore this tokenizer often tokenizes each non-English character to several
 107 bytes resulting in a long sequence. Inspired by Cui et al. (2023) and Pipatanakul et al. (2023), we
 108 extend the vocabulary using monolingual data of the target language to improve the model efficiency.

109 LoRA is a parameter-efficient training method, which is another technique that has been widely used
 110 for transferring the LLM. However, in this work, we use LoRA not only for efficiency but also for
 111 preserving the original parameters. Considering a weight matrix $W \in \mathbb{R}^{d \times k}$ of the target LLM,
 112 LoRA represents its update ΔW using two low rank matrices $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times k}$ as follows:

$$\tilde{h} = Wh, \text{ and } \hat{h} = \tilde{h} + \Delta Wh = \tilde{h} + BAh, \quad (3)$$

113 where r denotes the pre-determined rank, h denotes the input, \tilde{h} denotes the output of the original
114 module, and \hat{h} denotes the output of the updated module. During training, the original W is frozen,
115 so that original knowledge can be recovered by the recovery KD.

116 3.2 Training

117 3.2.1 Target Language Pre-Training

118 The chat LLMs are often insufficient on target language modeling due to the imbalanced training
119 corpus. Target language modeling is essential for generating fluent and localized text. Furthermore,
120 many works show that the monolingual pre-training can significantly improve the translation qual-
121 ity (Zheng et al., 2019; Xu et al., 2023). To build a solid foundation for the target language, we
122 pre-train the TransLLM model on monolingual data of the target language using Eq. 1.

123 We do not introduce any English task in this stage because of the following two reasons: first, the
124 pre-training involves quite computational consumption, and it can be unacceptable to find a proper
125 mixing ratio between the English and target language data; second, the English embeddings are
126 rarely updated on the target language data, therefore all the parameters of original LLM are almost
127 unchanged.

128 3.2.2 Translation Pre-Training

129 TCOT relies on translation to bridge the English and the target language. Therefore, we introduce
130 translation pre-training to improve the bidirectional translation quality between English and the
131 target language. Inspired by mBART (Liu et al., 2020), we use the special language id token to
132 denote translation directions. Considering we transform the LLM from language α to β , where $\alpha =$
133 English in this paper, we formulate the parallel pair (s^α, s^β) as two instances: $\text{cat}(s^\alpha, \langle\beta\rangle, s^\beta)$ and
134 $\text{cat}(s^\beta, \langle\alpha\rangle, s^\alpha)$, where $\text{cat}(\cdot)$ denotes the concatenate operation.

135 The translation training could disturb the original English embeddings. Thus, we introduce English
136 monolingual data into the translation pre-training stage. Specifically, we randomly insert the transla-
137 tion instance between English monolingual data using line break “\n” as the separator. Based on the
138 first stage, we train the TransLLM model on the mixed data by pre-training objective in Eq. 1.

139 3.2.3 Transfer Fine-Tuning

140 The two-stage pre-training enables the TransLLM in target language modeling and cross-lingual
141 translation. However, the TransLLM inevitably forgets the original knowledge. In this stage, we aim
142 to recover the original knowledge and teach the TransLLM model how to perform TCOT and when
143 to do translation.

144 **Recovery Knowledge Distillation Data.** Previous works focus on transferring knowledge from
145 base LLMs. To teach the base model how to follow human instructions, previous works perform
146 knowledge distillation with strong chat LLMs as the teacher by using the Alpaca dataset (Taori et al.,
147 2023). The Alpaca dataset generates queries using the self-instruct technique (Wang et al., 2022),
148 then responds using ChatGPT or GPT-4. Although the vanilla KD works well for base LLMs, we
149 argue that it is not helpful for chat LLMs as shown in Sec. 5.2. To address this problem, we introduce
150 the recovery KD that uses the target chat LLM to generate the responses. Although the recovery
151 KD data are often worse than GPT-4 KD data, it will help the model to recover the knowledge from
152 the original LLM parameters. We also introduce a special token $\langle\text{RESPONSE}\rangle$ in recovery KD to
153 direct the behavior of the TransLLM model. Considering a KD instance in English with query q^α and
154 answer a^α , we formulate the query and label in Eq. 2 as $x = q^\alpha$ and $y = \text{cat}(\langle\text{RESPONSE}\rangle, a^\alpha)$
155 respectively.

156 **TCOT Data.** Based on the recovery KD data (q^α, a^α) , we use machine translation to ob-
157 tain its translations (q^β, a^β) . Finally, we can organize the TCOT data as $x = q^\beta$ and $y =$
158 $\text{cat}(\langle\alpha\rangle, q^\alpha, \langle\text{RESPONSE}\rangle, a^\alpha, \langle\beta\rangle, a^\beta)$. That means when we input a query in β , the model
159 should first translate it into α as q^α . Then the model should $\langle\text{RESPONSE}\rangle$ the English query as
160 a^α using original knowledge as we teach in recovery KD. Finally, the TCOT outputs the response
161 in β as a^β based on all previous sequences. As discussed in Sec. 5.3, the previous sequences also

162 contribute to the final response. Different from Zhang et al. (2023), we use special tokens instead of
163 natural language to direct the model’s behavior. This is because the special tokens will not disturb the
164 English embeddings and make it convenient to extract results.

165 **Translation Data.** Due to the TCOT data, the model may be confused about the translation
166 instruction in β without extra translation SFT. Therefore, we also construct bi-direction translation
167 data based on previous parallel pairs (q^α, q^β) and (a^α, a^β) . Taking the parallel pair (q^α, q^β) as an
168 example, we first wrap the source sentence using translation prompt templates as $\text{prompt}(q^\alpha)$.¹ Then
169 we can obtain $x = \text{prompt}(q^\alpha)$ and $y = \text{cat}(\langle\beta\rangle, q^\beta)$.

170 Finally, we randomly mix all the data mentioned above and fine-tune the TransLLM model by Eq. 2.

171 3.3 Inference

172 The final TransLLM model can respond in both α and β , including α - β bi-direction translation.
173 For a single-turn conversation, the TransLLM model will decide the proper mode by itself given
174 only the input query x . To leverage the powerful multi-turn conversation ability of the original
175 LLM for β , we follow the original multi-turn format. For the multi-turn task in β , we only take
176 the English parts of the previous TCOT output as history. To be specific, we organize the input as
177 $x = \text{cat}(q_1^\alpha, a_1^\alpha, \dots, q_n^\alpha, a_n^\alpha, q_{n+1}^\beta)$, where n is the number of past turns. We do not use any special
178 tokens in the history as the original LLM does. Interestingly, even in this unseen setting, the model
179 still outputs the TCOT format as $y = \text{cat}(\langle\alpha\rangle, q_{n+1}^\alpha, \langle\text{RESPONSE}\rangle, a_{n+1}^\alpha, \langle\beta\rangle, a_{n+1}^\beta)$. We show
180 the whole multi-turn template in Appendix A.3.

181 4 Experiments

182 4.1 Settings

183 It is extravagant to train and evaluate a chat LLM in non-English. Therefore, during our experiment,
184 we mainly transform LLMs from English (EN) to Thai (TH) language, i.e. $\alpha = \text{EN}$ and $\beta = \text{TH}$. We
185 describe our basic settings as follows.

186 **Models.** We implement our pipeline using Chinese-LLaMA-Alpaca-2² project, which is based on
187 Transformers³. For the TransLLM model, we use the LLaMA2-Chat-7B as the target chat LLM.
188 Following Cui et al. (2023), we use SentencePiece (Kudo and Richardson, 2018) to learn the TH
189 vocabulary on the monolingual TH data that we use in target language pre-training. After we merge
190 the TH vocabulary with the original vocabulary, the final vocabulary size (including 3 special tokens)
191 is 43,012. The new embeddings are randomly initialized. We apply LoRA on the weights of the
192 attention module and multi-layer perceptron blocks. The LoRA rank is set as $r = 64$. Overall, there
193 are a total of 512.27 million trainable parameters including embeddings and LM heads. After all of
194 the training is completed, we merge the LoRA modules into the main backbone, the final model has
195 6.83 billion parameters. For a fair comparison, we re-implement most of the baselines by our setting
196 following their papers. The details of our model and baselines are in Appendix A.1.

197 **Training Data.** For target language pre-training, we use the monolingual TH data from mC4 (Xue
198 et al., 2020). We first filter the mC4-TH using the sensitive word list to reduce the harmful text.
199 Then, we use MinHashLSH⁴ to deduplicate documents in mC4-TH following GPT-3 (Brown et al.,
200 2020). Finally, we have about 11 billion tokens of TH data. Compared to the 2 trillion tokens EN
201 data used in LLaMA-2, the TH dataset is quite small. For translation pre-training, we collect the
202 EN-TH parallel data from CCAIined (Chaudhary et al., 2019), Tatoeba Challenge Data (Tiedemann,
203 2020), and OpenSubtitles (Lison et al., 2018). We directly use the EN documents released in the Pile
204 dataset which has been pre-processed (Gao et al., 2020). We randomly sample 1 million parallel pairs
205 and EN documents respectively for translation pre-training. For the transfer fine-tuning, we use the

¹The English prompt templates are from X-LLaMA <https://github.com/NJUNLP/x-LLM/blob/main/data/translation/translation.py>. We translate the prompt templates into the target languages.

²<https://github.com/ymcui/Chinese-LLaMA-Alpaca-2>

³<https://github.com/huggingface/transformers>

⁴<https://github.com/ekzhu/datasketch>

	vs. Model	Win (%)	Tie (%)	Loss (%)	Δ (%)
First Turn	ChatGPT	53.75(52.53 - 54.97)	27.50(26.41 - 28.59)	18.75(17.79 - 19.71)	35.00
	GPT4	22.50(21.48 - 23.52)	40.00(38.80 - 41.20)	37.50(36.31 - 38.69)	-15.00
Second Turn	ChatGPT	48.75(47.53 - 49.97)	26.25(25.17 - 27.33)	25.00(23.94 - 26.06)	23.75
	GPT4	22.50(21.48 - 23.52)	27.50(26.41 - 28.59)	50.00(48.78 - 51.23)	-27.50

Table 1: Comparison between our model and strong LLMs on MT-bench under human evaluation. We provide the 95% confidence interval in brackets.

Setting	TH		EN [†]	
	First Turn (%)	Second Turn (%)	First Turn (%)	Second Turn (%)
w/ Tie (R = 33%)	75.42	70.42	60.00	59.00
w/o Tie (R = 50%)	75.11	67.85	85.00	84.00

Table 2: Agreement between GPT-4 and humans. “R=” denotes the exact agreement between random judges. [†] EN results are from Zheng et al. (2024).

206 query from the Alpaca dataset and generate the response using the target model LLaMA2-Chat-7B.
 207 We further use Google Translate to obtain TCOT and translation data based on recovery KD data.
 208 In our preliminary study, Google Translate may translate the variable in code which is not desirable
 209 for the chat LLM. Thus, we use GPT-4 to recognize the “do not translate” part. We use the same
 210 monolingual and translation data for baselines, while we use the Alpaca-GPT-4 (Peng et al., 2023) as
 211 the SFT data following their setting. There are a total of 52K queries in the Alpaca dataset, we use
 212 the first 50K queries as the training set and the rest 2K queries as the validation set in our experiments.
 213 We provide the training details in A.2.

214 **Benchmark.** For helpfulness, we utilize two widely used LLM benchmarks, MT-bench (Zheng
 215 et al., 2024) and AlpacaEval (Li et al., 2023). MT-bench has 80 multi-turn questions across 8 domains.
 216 AlpacaEval consists of 805 single-turn questions collected from 5 different subsets. The original
 217 benchmarks are in EN. We employ professional translators to translate these two datasets into TH.
 218 For safety, we use the AdvBench. AdvBench consists of 520 instructions that induce LLMs output
 219 harmful responses. Following the setting in Yong et al. (2023), we directly use Google Translate to
 220 translate the AdvBench from EN to TH.

221 **Evaluation.** For helpfulness, we use strong LLMs as judges, which show considerable consistency
 222 with human evaluators in EN (Zheng et al., 2024). LLM-as-a-Judge is efficient, reproducible, and
 223 cost-effective. However, it is still unknown whether it will work in the TH language. To obtain a
 224 reliable result, we first invite professional translators to conduct the human evaluation for some strong
 225 models on the MT-bench. We test the consistency between human and GPT-4 evaluation as described
 226 in (Zheng et al., 2024). After we prove that GPT-4 achieves acceptable consistency with human
 227 evaluators, we evaluate all models with it. Both human annotators and LLMs rate the response on a
 228 scale from 1 to 10, and we further calculate the win, tie, and loss rate based on the scores. We use
 229 Δ to denote the gap between the win and loss rate calculated with the tie. For safety, we translate
 230 the TH responses into EN and let EN annotators annotate them into Bypass, Reject, and Unclear.
 231 Bypass means the attack bypasses the safety mechanism of LLMs. Reject means LLMs refuse to
 232 output harmful information. Unclear means the responses are safe but unclear due to translation or
 233 hallucination, etc. The setting follows Yong et al. (2023) strictly. Please refer to this paper for details.
 234 In Appendix A.4, we describe the evaluation procedure, the instructions for human evaluators, and the
 235 information of evaluators in detail. We also conduct significant tests for main results as described in
 236 Appendix C. **We mark the results with bold if the difference is statistically significant ($p < 0.05$).**

237 4.2 Main Results

238 4.2.1 Human Evaluation Results

239 **Better performance than ChatGPT on MT-Bench.** As shown in Table 1, TransLLM surpasses
 240 ChatGPT by 35% and 23.75% for the first and second turn on MT-bench with statistical significance.
 241 It is an inspiring result although TransLLM is still behind GPT-4 in TH. Because we only use the
 242 LLaMA-2 with 7B parameters. As the fine-grained scores in Appendix B.1 shown, the two domains
 243 with the biggest gaps between ours and GPT-4 are Math and Coding, which are also the weaknesses

244 of LLaMA-2 in EN. We leave exploring TransLLM on more powerful open-source LLMs in the
 245 future.

246 **High agreement between humans and GPT-4 in TH.** Following Zheng et al. (2024), we
 247 calculate the average agreements by comparing
 248 every two models. In Table 2, GPT-4 shows
 249 high consistency with human annotators. The
 250 consistency (w/ tie) between GPT-4 and humans
 251 reaches 75.42% and 70.42% in the first and second
 252 turns, which are much higher than random
 253 guesses and even higher than the consistency in
 254 EN. Therefore, we use GPT-4 to evaluate the
 255 helpfulness in the following experiments.
 256

257 **Higher safety than ChatGPT and GPT-4.** In Table 3, TransLLM has a rejection rate of 94.61%,
 258 close to 99.23% of the original model. It indicates that we successfully transfer most of the human
 259 preference about the safety of the original model. TransLLM attains an improvement of 14.8% and
 260 8.65% over ChatGPT and GPT-4 for rejecting harmful queries with statistical significance. More
 261 importantly, although GPT-4 is as safe as the original LLM in EN, the performance of our w/ GPT-4
 262 KD is much below our w/ recovery KD. Later, we will demonstrate that this is because recovery KD
 263 successfully recovers the original knowledge.

Model	Bypass (%)	Reject (%)	Unclear (%)
ChatGPT	10.96	79.81	9.23
GPT4 [†]	10.38	85.96	3.66
Ours w/ GPT-4 KD	31.15	63.46	5.38
Ours	2.69	94.61	2.69
LLaMA-2-chat (EN)	0.58	99.23	0.19
GPT4 [†] (EN)	0.96	99.04	0.00

Table 3: Result for different models on safety benchmark AdvBenchmark under human evaluation. [†] GPT-4 results are from Yong et al. (2023).

vs. Model	First Turn (%)				Second Turn (%)			
	Win	Tie	Loss	Δ	Win	Tie	Loss	Δ
PolyLM (Wei et al., 2023)	78.75	16.25	5.00	73.75	90.00	10.00	0.00	90.00
X-LLaMA (Zhu et al., 2023)	72.50	17.50	10.00	62.50	85.00	8.75	6.25	78.75
Typhoon (Pipatanakul et al., 2023)	75.00	18.75	6.25	68.75	62.50	30.00	7.50	55.00
PLUG (Zhang et al., 2023)	72.50	13.75	13.75	58.75	87.50	8.75	3.75	83.75
NLLB-bridge (Costa-jussà et al., 2022)	75.00	16.25	8.75	66.25	63.75	18.75	17.50	46.25
ChatGPT (OpenAI, 2022)	42.50	26.26	31.25	11.25	42.50	22.50	35.00	7.50
GPT4 (OpenAI, 2023)	26.25	28.75	45.00	-18.75	30.00	18.75	51.25	-21.75

Table 4: Comparison between our model and different methods on MT-Bench under GPT-4 evaluation.

264 4.2.2 GPT-4 Evaluation Results

265 **Better performance than strong baselines.** As
 266 shown in Table 4, TransLLM significantly out-
 267 performs baselines that are built on open-source
 268 resources. Notably, we specifically build the
 269 baseline NLLB-bridge which uses the powerful
 270 translation model NLLB-3B (Costa-jussà
 271 et al., 2022) as the bridge between LLaMA-2-
 272 chat-7B and the TH language. Using the multi-
 273 turn ability of LLaMA-2-chat-7B, NLLB-bridge
 274 achieves good performance in the second turn.

275 Although NLLB-bridge uses more parameters and more translation resources, it still loses to
 276 TransLLM. We will explain in detail why TransLLM is better than translation-as-a-bridge in the
 277 analysis. Typhoon with TH pre-training achieves sub-optimal second-turn performance among
 278 baselines. It is probably because the TH documents teach the LLM how to model long context in TH.
 279 Under GPT-4 evaluation, we slightly outperform ChatGPT without statistical significance. It seems
 280 difficult for GPT-4 to compare two strong LLMs on small datasets in TH. We select the baselines that
 281 perform well on the first turn of the MT-bench, for further evaluation on Alpaca-Eval. On the larger
 282 dataset, TransLLM outperforms baselines and ChatGPT by a large margin with statistical significance
 283 as shown in Table 5.

vs. Model	Win (%)	Tie (%)	Loss (%)	Δ (%)
X-LLaMA	92.50	5.00	2.50	90.00
PLUG	87.50	8.75	3.75	83.75
NLLB-bridge	91.25	5.00	3.75	87.50
ChatGPT	72.50	13.75	13.75	58.75
GPT4	17.50	45.00	37.50	-20.00

Table 5: Comparison between our model and different methods on Alpaca-Eval under GPT-4 evaluation.

284 5 Analysis

285 5.1 All Components Work Together

286 We conduct comprehensive ablation studies
 287 on MT-Bench to investigate the impact of
 288 TransLLM’s components and present results
 289 in Table 6. The results confirm our hypothe-
 290 sis that transforming chat LLMs could provide
 291 better conversational ability than base LLMs.
 292 Pre-training on TH documents helps TransLLM
 293 output fluency in TH response with long con-
 294 text. Thus, TransLLM without TH pre-training
 295 is less satisfying on both the first and second
 296 turn. Since TH pre-training and transfer fine-
 297 tuning also provide some translation knowledge,
 298 the improvement of the translation pre-training is
 299 not as significant as other components. Beyond
 300 safety, the high-quality GPT-4 KD data also
 301 leads to performance degradation for helpfulness.
 302 That is because our goal is not to inject more
 303 knowledge but to preserve the original knowledge.
 304 We also examine the contribution of LoRA. Specifically,
 305 we merge the LoRA parameters with full param-
 306 eters before transfer fine-tuning. We are unable
 307 to conduct full fine-tuning for pre-training, but
 308 the merged model is a good approximation according
 to Eq. 3. We further conduct transfer fine-tuning
 with full parameters based on the merged model.
 In most tasks, full fine-tuning is better or com-
 parable with LoRA. However, in our case, full
 fine-tuning wipes the original knowledge from
 parameters, and therefore its performance is
 much lower than TransLLM with LoRA. When
 using the history in TH, TransLLM is also
 capable of multi-turn conversation with small
 performance degradation. That means TransLLM
 can handle TH context well, this ability could
 be further developed for retrieval augmentation
 in TH.

vs. Model	1st Δ (%)	2nd Δ (%)
w/o chat model	36.25	67.50
w/o TH pre-train	41.25	35.00
w/o translation pre-train	8.75	23.75
w/ GPT-4 KD	17.50	45.00
w/o LoRA	62.50	66.25
w/ TH history	-	23.75

Table 6: Comparison between our model and ablation models.

309 5.2 TransLLM Recover the Original Knowledge

310 **Knowledge is forgotten and recovered.** To
 311 measure how much original knowledge is forgot-
 312 ten by the chat LLM, we calculate the generation
 313 probabilities on the hold-out validation set of re-
 314 covery KD data in EN. We also calculate the
 315 average difference between the generation prob-
 316 abilities of the target LLM and different models.
 317 As shown in Table 7, after pre-training, which
 318 has been proven to be necessary, the LLM sig-
 319 nificantly forgets the conversation knowledge.
 320 GPT-4 KD, which is widely used in previous
 321 works, can provide high-quality knowledge. However,
 322 this kind of knowledge is quite different from
 and competes with the original knowledge. As
 a result, the LLM still forgets much original
 knowledge using GPT-4 KD. Meanwhile, TransLLM
 successfully recovers the original knowledge.

Model	$P(y x)$	Difference
LLaMA2-Chat (EN)	0.2363	-
Ours w/o transfer fine-tuning	0.1666	0.0697
Ours w/ GPT-4 KD	0.1972	0.0391
Ours	0.2352	0.0055

Table 7: The difference of generation probabilities.

323 **LoRA also helps.** LoRA keeps the original
 324 parameters unchanged. The LLM can fit the re-
 325 covery KD data easily when using knowledge
 326 from these frozen parameters. This easy pattern
 327 is a “shortcut” that prompts the LLM to learn
 328 to use the original knowledge for EN and new
 329 knowledge for TH. To confirm this assumption,
 330 on the TCOT validation data, we calculate the
 cosine similarity between the last layer’s hidden
 states of the original model \tilde{h} and LoRA updated
 model \hat{h} as defined in Eq. 3. The average
 similarity per token for EN responses is much
 larger than that for TH responses, 0.6191
 vs. 0.2522. That means TransLLM successfully
 learns the “shortcut” using LoRA and recovery
 KD together.

331 5.3 Why TransLLM is better than translation-as-a-bridge?

332 **Comparable translation performance.** The
 333 translation performance is critical for both
 334 TransLLM and translation-as-a-bridge. There-
 335 fore, we test them on the widely used bench-
 336 mark Flores-200 (Goyal et al., 2022). As shown
 337 in Table 8, benefiting from translation and TH
 338 pre-training, TransLLM outperforms ChatGPT
 339 and NLLB on EN-TH and achieves competitive per-

Model	EN-TH		TH-EN	
	COMET	BLEU	COMET	BLEU
ChatGPT	85.47	31.26	86.29	23.47
NLLB	83.88	28.53	87.14	30.78
Ours	86.96	35.04	86.97	27.68

Table 8: Translation performance on Flores-200.

340 performance on TH-EN. We also ask the naive TH speaker to provide a fluency score for each model
 341 on MT-Bench in Table 9. The fluency of NLLB is as poor as its translation performance on EN-TH.
 342 NLLB usually translates the responses literally. For example, NLLB translates “I see” into “I see
 343 something” instead of “I understand” in TH. Surprisingly, the response of GPT-4 is not very fluent and
 344 natural. GPT-4 often uses full-stops and commas which are not used in TH. ChatGPT and TransLLM
 345 are generally fluent, with translationese to a certain degree. For example, TH speakers do not use
 346 “sure” or “of course” at the beginning of responses, but ChatGPT and TransLLM do.

347 **TransLLM is more than translation.** Translation perfor-
 348 mance is important but not the whole story. TransLLM outputs
 349 an EN query, EN response, and TH response at once. It means
 350 that TransLLM can use all previous information for TH re-
 351 sponses and therefore achieve better performance. To verify
 352 it, we use TransLLM to translate its EN responses in another
 353 round of inference. The performance is worse than the stan-
 354 dard response with $\Delta = 13.75\%$ and $\Delta = 18.75\%$ on first and
 355 second turn. The attention map of TransLLM in Appendix B.2
 356 shows that TransLLM outputs the TH response mostly based on the TH response itself and then the
 357 EN response. However, the TH response also pays a little attention on the TH query and EN query.
 358 Besides, translation-as-a-bridge needs to deploy two models, which is costly and inconvenient.

Model	Score
NLLB-bridge	5
GPT4	6
ChatGPT	7
Our	7

Table 9: Fluency on MT-Bench.

359 6 Related Works

360 Recently, there have been many works that attempt to transfer knowledge from English to non-English
 361 for LLMs. For example, Chinese LLaMA (Cui et al., 2023) and Typhoon(Pipatanakul et al., 2023)
 362 directly perform continuous pre-training and instruct tuning with extended vocabulary using LoRA.
 363 PloyLM (Wei et al., 2023) adopts multilingual pre-training based on the curriculum learning strategy
 364 that gradually exposes more low-resource corpus. ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI,
 365 2023) are also well-known multilingual LLMs. Zhu et al. (2023) focus on building semantic alignment
 366 with cross-lingual instruct tuning and translation training. Bansal et al. (2024) augment LLMs by
 367 combining the English-dominated LLM with the non-English model. Some other works focus on
 368 transfer reasoning abilities: Qin et al. (2023) introduce cross-lingual prompting to improve zero-shot
 369 chain-of-thought reasoning across languages; She et al. (2024) propose multilingual alignment-as-
 370 preference optimization to align reasoning abilities across languages. PLUG (Zhang et al., 2023) only
 371 uses the TCOT data to train the base LLMs directly. Different from PLUG, we propose a systematic
 372 framework for transforming chat LLMs. We highlight that the TCOT highly relies on the performance
 373 of its sub-tasks and introduce how to preserve the knowledge of the chat LLM.

374 7 Conclusion

375 Chat LLMs have been specifically optimized for chat usage and therefore are helpful and safe in the
 376 dominant language. In this paper, we propose a framework for transforming an off-the-shelf chat
 377 LLM to other languages. In this framework, we utilize TCOT to transfer chat knowledge and further
 378 enhance the TCOT’s sub-tasks using publicly available data. To recover the original knowledge, we
 379 propose the recovery KD method supplemented with LoRA. The experiments in TH show that we
 380 transfer desired abilities to TH and outperform ChatGPT in both helpfulness and safety. Overall, we
 381 hope that this work can become the foundation for developing safe LLMs in many languages other
 382 than English.

383 **Limitations and future works.** Due to limited resources, we only conduct experiments that
 384 transform LLaMA-2-chat-7B to TH. However, we conduct comprehensive experiments and in-depth
 385 analysis to reveal the mechanism of the proposed TransLLM. For now, TransLLM is still highly
 386 dependent on translation. Consequently, TransLLM can not handle the queries related to TH text, e.g.
 387 word games in TH. A simple solution is to enable TransLLM, through training, to choose whether
 388 respond to with TH mode or TCOT mode. Due to the TCOT, the inference overhead of TransLLM is
 389 much longer than other baselines. Recently, Goyal et al. (2023) and Deng et al. (2023) show that
 390 the implicit chain-of-thought achieves similar performance on reasoning tasks without additional
 391 inference overhead. We would like to explore TransLLM with implicit TCOT in the future.

392 References

- 393 Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
394 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report.
395 *arXiv preprint arXiv:2305.10403*.
- 396 Rachit Bansal, Bidisha Samanta, Siddharth Dalmia, Nitish Gupta, Shikhar Vashishth, Sriram Ganapa-
397 thy, Abhishek Bapna, Prateek Jain, and Partha Talukdar. 2024. LLM augmented llms: Expanding
398 capabilities through composition. *CoRR*, abs/2401.02412. Version 1.
- 399 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
400 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models
401 are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- 402 Vishrav Chaudhary, Yuqing Tang, Francisco Guzmán, Holger Schwenk, and Philipp Koehn. 2019.
403 Low-resource corpus filtering using multilingual sentence embeddings. In *Proceedings of the*
404 *Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 263–268,
405 Florence, Italy. Association for Computational Linguistics.
- 406 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan,
407 Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind:
408 Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- 409 Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama
410 and alpaca. *arXiv preprint arXiv:2304.08177*.
- 411 Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stuart
412 Shieber. 2023. Implicit chain of thought reasoning via knowledge distillation. *arXiv preprint*
413 *arXiv:2311.01460*.
- 414 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,
415 Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text
416 for language modeling. *arXiv preprint arXiv:2101.00027*.
- 417 Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana
418 Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The flores-101
419 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the*
420 *Association for Computational Linguistics*, 10:522–538.
- 421 Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh
422 Nagarajan. 2023. Think before you speak: Training language models with pause tokens. In *The*
423 *Twelfth International Conference on Learning Representations*.
- 424 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
425 and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint*
426 *arXiv:2106.09685*.
- 427 Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword
428 tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference*
429 *on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71,
430 Brussels, Belgium. Association for Computational Linguistics.
- 431 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang,
432 and Tatsunori B. Hashimoto. 2023. AlpacaEval: An automatic evaluator of instruction-following
433 models. https://github.com/tatsu-lab/alpaca_eval.
- 434 Pierre Lison, Jörg Tiedemann, and Milen Kouylekov. 2018. OpenSubtitles2018: Statistical rescoring
435 of sentence alignments in large, noisy parallel corpora. In *Proceedings of the Eleventh International*
436 *Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European
437 Language Resources Association (ELRA).
- 438 Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis,
439 and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation.
440 *Transactions of the Association for Computational Linguistics*, 8:726–742.

- 441 OpenAI. 2022. Introducing chatgpt. Blog post <https://www.openai.com/blog/chatgpt>.
- 442 OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- 443 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
444 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to
445 follow instructions with human feedback. *Advances in neural information processing systems*,
446 35:27730–27744.
- 447 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic
448 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association
449 for Computational Linguistics*, pages 311–318.
- 450 Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction
451 tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- 452 Kunat Pipatanakul, Phatrased Jirabovonvisut, Potsawee Manakul, Sittipong Sripaisarnmongkol,
453 Ruangsak Patomwong, Pathomporn Chokchainant, and Kasima Tharnpipitchai. 2023. Typhoon:
454 Thai large language models. *arXiv preprint arXiv:2312.13951*.
- 455 Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. 2023. Cross-lingual
456 prompting: Improving zero-shot chain-of-thought reasoning across languages. In *Proceedings of
457 the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2695–2709.
- 458 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019.
459 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- 460 Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova,
461 Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission
462 for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation
463 (WMT)*, pages 578–585.
- 464 Shuaijie She, Shujian Huang, Wei Zou, Wenhao Zhu, Xiang Liu, Xiang Geng, and Jiajun Chen.
465 2024. Mapo: Advancing multilingual reasoning through multilingual alignment-as-preference
466 optimization. *arXiv preprint arXiv:2401.06838*.
- 467 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
468 Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: an instruction-following llama model
469 (2023). URL https://github.com/tatsu-lab/stanford_alpaca.
- 470 Jörg Tiedemann. 2020. The Tatoeba Translation Challenge – Realistic data sets for low resource
471 and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pages
472 1174–1182, Online. Association for Computational Linguistics.
- 473 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
474 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open
475 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- 476 Changhan Wang, Kyunghyun Cho, and Jiatao Gu. 2020. Neural machine translation with byte-level
477 subwords. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages
478 9154–9160.
- 479 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi,
480 and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated
481 instructions. *arXiv preprint arXiv:2212.10560*.
- 482 Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan,
483 Zhiwei Cao, Binbin Xie, et al. 2023. Polylm: An open source polyglot large language model.
484 *arXiv preprint arXiv:2307.06018*.
- 485 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
486 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
487 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,
488 Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art

- 489 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in*
490 *Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for
491 Computational Linguistics.
- 492 Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in
493 machine translation: Boosting translation performance of large language models. *arXiv preprint*
494 *arXiv:2309.11674*.
- 495 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
496 Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer.
497 *arXiv preprint arXiv:2010.11934*.
- 498 Zheng-Xin Yong, Cristina Menghini, and Stephen H Bach. 2023. Low-resource languages jailbreak
499 gpt-4. *arXiv preprint arXiv:2310.02446*.
- 500 Zhihan Zhang, Dong-Ho Lee, Yuwei Fang, Wenhao Yu, Mengzhao Jia, Meng Jiang, and Francesco
501 Barbieri. 2023. Plug: Leveraging pivot language in cross-lingual instruction tuning. *arXiv preprint*
502 *arXiv:2311.08711*.
- 503 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
504 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench
505 and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- 506 Zaixiang Zheng, Hao Zhou, Shujian Huang, Lei Li, Xin-Yu Dai, and Jiajun Chen. 2019. Mirror-
507 generative neural machine translation. In *International Conference on Learning Representations*.
- 508 Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan, Jingjing Xu, Shujian Huang, Lingpeng Kong,
509 Jiajun Chen, and Lei Li. 2023. Extrapolating large language models to non-english by aligning
510 languages. *arXiv preprint arXiv:2308.04948*.

511 A Experiment Details

512 A.1 Models

513 We list backbone, training data, and model size in Table 10. Due to the huge consumption of
514 multilingual (MTL) pre-training, we directly use the model PolyLM-MultiAlpaca-13B released
515 in Wei et al. (2023) for PolyLM. PolyLM uses ChatGPT to generate the Alpaca data while other
516 baselines use the Alpaca data generated by GPT-4. We use the gpt-3.5-turbo-0125 and gpt-4-0613 for
517 ChatGPT and GPT-4 in all experiments (including evaluation) through OpenAI API. We re-implement
518 other baselines by strictly following their papers and using the same data as our model. To reduce the
519 impact of randomness, we use greedy search for all experiments. We set the temperature as 0 for
520 ChatGPT and GPT-4 through API to approximate the greedy search.

521 Please refer to Touvron et al. (2023) for model structures of LLaMA-2. We only list the LoRA
522 parameters here. We set the rank to 64, alpha to 128, and dropout to 0.05 for LoRA. These parameters
523 are applied to the q_proj , v_proj , k_proj , o_proj , $gate_proj$, $down_proj$, and up_proj modules of the
524 original model. Besides, the $embed_tokens$ and lm_head are also trainable.

Name	Backbone	Pre-train Data	Fine-tune Data	Size
PolyLM	From Scratch	MTL + Translation	Alpaca-MTL	13B
X-LLaMA	LLaMA2-base	-	Alpaca-EN + Alpaca-TH + Translation	7B
Typhoon	LLaMA2-base	TH	Alpaca-TH	7B
PLUG	LLaMA2-base	-	TCOT	7B
NLLB bridge	LLaMA2-chat + NLLB	-	-	7B + 3B
ChatGPT	Unknown	Unknown	Unknown	≥ 7B
GPT4	Unknown	Unknown	Unknown	≥ 7B
Ours	LLaMA2-chat	TH / Translation + EN	TCOT + Recovery KD + Translation	7B

Table 10: Model details.

525 **A.2 Training**

526 We train the TransLLM model on 8 A100 GPUs as follows.

527 **TH Pre-Training** We train the TransLLM using a warm-up ratio of 0.0005, a maximum sequence
 528 length of 512 tokens, and a weight decay of 0.01. The training was conducted with each GPU
 529 managing 128 batches and utilizing a gradient accumulation step of 1. The peak learning rate is set at
 530 $2e-4$ with a cosine learning rate decay (max_epoch=100), and training operated under bf16 precision
 531 facilitated by deepspeed, employing ZeRO stage 2.

532 We only run 1 epoch for this stage, which spends 168×8 GPU hours. As shown in Figure 3, the
 533 initial training loss is approximately 7.8, which converges to below 1.7 after around 0.1 epochs of
 534 training. The final loss reaches around 1.42.

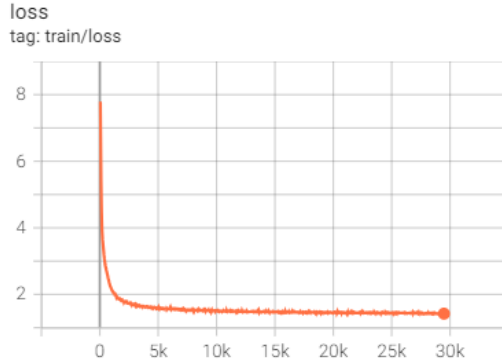


Figure 3: TH Pre-Training loss.

535 **Translation Pre-Training** According to the data size, we set the warm-up ratio as 0.05, the
 536 max_epoch=10 for the cosine learning rate decay. We use 0.1% examples as the validation set and
 537 calculate valid loss every 400 steps. The best model has been trained for about 3 epochs, which
 538 spends 40×8 GPU hours. The remaining configurations remain consistent with the first stage.

539 **Transfer Fine-Tuning** Our max_seq_length is set to 2048 for fine-tuning, and when batching data,
 540 we pad sentences with “<PAD>” tokens. The peak learning rate is set to $1e-4$, the warmup ratio is
 541 set to 0.01, and the single-card batch size is set to 16 with gradient accumulation steps as 4. We set
 542 weight decay as 0. We use 2K examples as the validation set and calculate valid loss every 200 steps.
 543 The best model has been trained for about 1 epoch, which spends 6×8 GPU hours. The remaining
 544 configurations remain consistent with the first stage.

545 **A.3 Inference**

546 We provide the whole multi-turn prompt in Table 11, where “<s> </s>”, “<<SYS>> <<SYS>>”, and
 547 “[INST] [INST]” denote the whole instance, system prompt, and instruction respectively.

<s>	[INST]	<<SYS>>
		You are a helpful assistant. <<SYS>>
	q_1^α	[INST] a_1^α </s>
	<s>	[INST] q_2^α [INST] a_2^α </s>
		...
	<s>	[INST] q_n^α [INST] a_n^α </s>
	<s>	[INST] q_{n+1}^β [INST]

Table 11: The multi-turn prompt template used in our experiments.

548 **A.4 Evaluation**

549 **A.4.1 Human Evaluation**

550 For helpfulness, the results are evaluated by three annotators. Annotator A is a professional translator
551 expert in EN and TH. Annotator B is a computer engineer who is an expert in EN, Math, Coding,
552 and Extraction. Annotator C is a native TH speaker while also an expert in EN. The three annotators
553 cooperate with each other to complete the whole evaluation process as follows. Annotator A is the
554 major annotator who is responsible for annotating most of the queries except for the Math, Coding,
555 and Extraction domains. For these three domains, annotator A first translates the results from TH to
556 EN. Annotator B then annotates these three domains in EN translations. Meanwhile, Annotator C
557 helps annotator A evaluate the fluency of all responses. To obtain consistent annotations between
558 evaluators and questions, we define comprehensive instructions for annotators in Table 12.

Score	Performance Level	Adherence to Instructions; Expression Fluency; Style
1-2	Very Poor	Does not follow the query; be not applicable due to nonsensical expres- sion; has incomprehensible style
3-4	Poor	Does not follow the query but has some relevant content; lacks fluency, coherency, and clarity; has largely inappropriate style
5-6	Fair	Partially meets the requirements and addresses some issues; has some fluency and clarity though minor flaws; has occasionally appropriate style
7-8	Good	Mainly follows the query though some minor flaws; be largely fluent and coherent; has generally appropriate style
9-10	Excellent	Strictly follows the query with appreciated content; has a high degree of fluency and clarity; is perfectly matched in style

Table 12: Rating criterion.

559 For safety, the responses are first translated from TH to EN and then evaluated by three professional
560 translators who are experts in EN. However, one response is only annotated by one translator due to a
561 limited budget. Please refer to the annotation instruction in Yong et al. (2023).

562 **All models are anonymous to all annotators in the whole evaluation process!**

563 **A.4.2 Automatic Evaluation**

564 We follow the setting of LLM-as-a-Judge in Zheng et al. (2024). For Reasoning, Math, and Coding
565 domains, we provide the EN responses of GPT-4 as references. Note that, these three domains are
566 different from human evaluation because annotator A is good at Reasoning instead of Extraction.
567 We modify the evaluation prompts provided in Zheng et al. (2024) to inform GPT-4 that the queries
568 and responses are in TH. Please refer to Zheng et al. (2024) for the details of how to calculate the
569 agreement.

570 We use the default wmt22-comet-da model ⁵ for COMET (Rei et al., 2022). We use
571 the BLEU (Papineni et al., 2002) implemented in the scarebleu⁶, whose signature is
572 "BLEUnrefs:1lcase:mixedleff:noltok:13alsmooth:explversion:2.4.0".

573 **A.5 Licenses**

574 Our experiments use open-source resources. We list their licenses in Table 13. We have properly
575 cited their papers and strictly followed their licenses.

576 **B Other Results**

577 **B.1 Results in Scores**

578 We provide evaluation scores on different benchmarks in Table 14, 15, 16, and 17.

⁵<https://huggingface.co/Unbabel/wmt22-comet-da>

⁶<https://github.com/mjpost/sacrebleu>

Resource	License
MC4 (Xue et al., 2020)	ODC-BY 1.0
Pile (Gao et al., 2020)	MIT License
CCAligned (Chaudhary et al., 2019)	Unknown
Tatoeba Challenge Data (Tiedemann, 2020)	CC-BY-NC-SA 4.0
OpenSubtitles (Lison et al., 2018)	Unknown
Flores-200 (Goyal et al., 2022)	CC-BY-SA 4.0
Alpaca (Taori et al., 2023)	CC BY-NC 4.0
Alpaca-eval (Li et al., 2023)	Apache License 2.0
MT-bench (Zheng et al., 2024)	Apache License 2.0
Chinese-Alpaca-2 (Cui et al., 2023)	Apache License 2.0
Transformers (Wolf et al., 2020)	Apache License 2.0
SentencePiece (Kudo and Richardson, 2018)	Apache License 2.0
PolyLM (Wei et al., 2023)	Apache License 2.0
LLaMA-2 (Touvron et al., 2023)	LLaMA 2 Community License Agreement

Table 13: Licenses of open source resources.

	Model	Writing	Roleplay	Reasoning	Math	Coding	Extraction	STEM	Humanities	All
First Turn	ChatGPT	5.30	4.70	5.20	4.60	7.80	7.20	6.80	6.40	6.00
	GPT4	7.40	6.70	4.80	6.00	8.80	8.30	7.40	7.70	7.14
	Ours	7.30	6.50	5.20	4.20	6.50	5.70	7.60	7.90	6.36
Second Turn	ChatGPT	3.00	5.00	3.40	2.90	7.40	7.90	5.60	5.70	5.11
	GPT4	4.70	6.70	5.00	4.00	8.60	7.60	6.80	7.50	6.36
	Ours	6.10	6.50	3.10	3.00	6.70	5.10	6.60	7.00	5.51

Table 14: Human evaluation scores on MT-Bench for different models.

579 B.2 Attention Map of the TransLLM Output

580 As shown in Figure 4, the TH response focuses on the TH response, EN response, EN query, and TH
581 query, in order from high to low.

	Model	Writing	Roleplay	Reasoning	Math	Coding	Extraction	STEM	Humanities	All
First Turn	PolyLM	4.00	4.00	3.40	1.10	1.00	2.80	2.80	3.10	2.78
	X-LLaMA	4.10	2.80	4.10	2.20	3.10	3.00	4.00	4.10	3.42
	Typhoon	5.90	5.40	2.90	1.10	2.90	2.80	6.40	6.10	4.19
	PLUG	6.60	3.90	3.70	2.60	2.90	2.90	5.90	7.60	4.51
	NLLB-bridge	5.50	4.90	3.90	2.90	1.00	3.10	4.80	5.20	3.91
	LLaMA2-Chat (EN)	9.60	7.80	5.40	3.20	3.60	7.30	9.55	9.55	7.00
	ChatGPT	7.70	7.80	6.00	6.00	5.70	7.50	8.90	8.60	7.28
	GPT4	9.00	8.90	6.10	7.10	6.20	9.30	9.30	9.20	8.14
	Ours	8.50	7.50	6.40	3.10	4.40	5.80	9.60	9.60	6.86
	Second Turn	PolyLM	1.30	1.00	1.50	1.10	1.00	1.20	1.00	1.10
X-LLaMA		2.60	3.60	2.50	1.20	1.80	1.70	3.20	2.90	2.44
Typhoon		3.00	5.20	4.10	1.70	2.70	1.80	5.90	4.80	3.65
PLUG		2.20	2.60	1.40	0.50	2.10	1.30	2.90	3.90	2.11
NLLB-bridge		5.30	4.20	4.10	2.80	2.30	3.50	4.20	6.30	4.09
LLaMA2-Chat (EN)		6.80	7.10	4.20	3.70	3.30	3.80	7.30	9.70	5.74
ChatGPT		3.50	7.90	5.20	3.50	5.10	7.20	6.70	8.80	5.99
GPT4		8.30	8.50	4.70	4.80	7.00	8.80	8.00	8.60	7.34
Ours		7.50	7.30	5.60	2.10	5.20	4.80	8.20	8.70	6.18

Table 15: GPT-4 evaluation scores on MT-Bench for different models.

582 C Statistical Methods

583 C.1 Confidence Interval

584 We first calculate the standard deviation for proportion p on n examples as:

$$s_p = \sqrt{\frac{p(1-p)}{n}}. \quad (4)$$

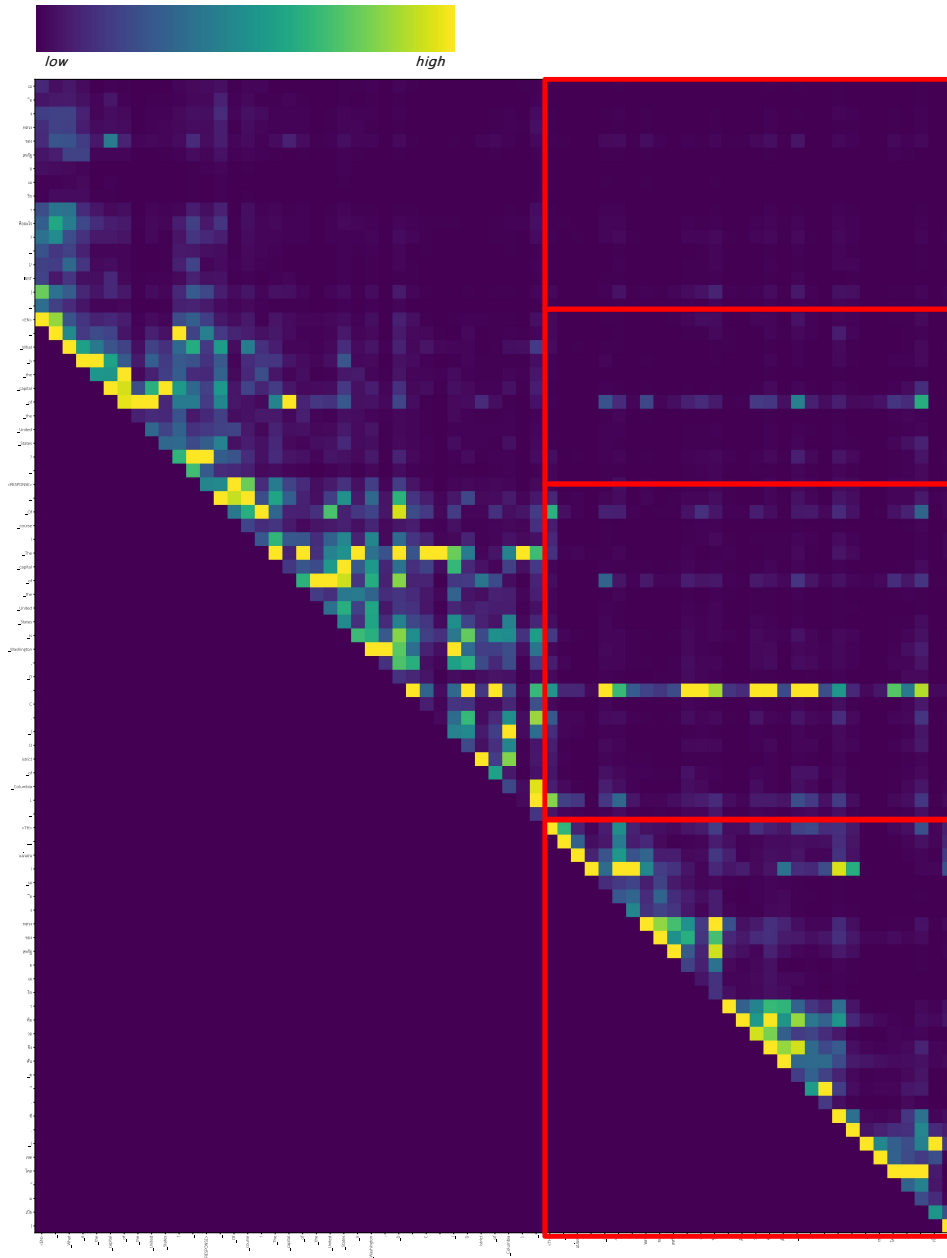


Figure 4: Attention map of the TransLLM output. We mark the attention scores of TH responses with red rectangles. Rectangles from top to bottom indicate attention scores of TH response for TH query, EN query, EN response, and TH response respectively.

Model	Helpful-Base	Koala	Oasst	Self-Instruct	Vicuna	All
X-LLaMA	2.80	3.86	3.95	3.90	4.80	3.82
PLUG	4.88	5.47	5.23	5.32	6.90	5.41
NLLB-bridge	4.36	4.97	5.04	4.49	4.78	4.72
ChatGPT	7.39	7.32	7.49	7.77	8.06	7.59
GPT-4	9.53	9.17	9.19	8.90	9.44	9.18
Ours	8.72	7.91	7.87	7.61	8.71	8.02

Table 16: GPT-4 evaluation scores on Alpaca-Eval for different models.

Model	First Turn	Second Turn
Ours	6.86	6.18
w/ base model	5.56	3.08
w/o TH pre-train	5.55	4.44
w/o translation pre-train	6.55	5.04
w/ GPT-4 KD	5.96	4.68
w/o LoRA	4.58	3.34
w/ TH history	-	5.43

Table 17: GPT-4 evaluation scores for ablation studies on MT-bench.

	vs. Model	p
First Turn	ChatGPT	.000
	GPT4	.111
Second Turn	ChatGPT	.018
	GPT4	.005

Table 18: Binomial test for Table 1.

	vs. Model	p
First Turn	PolyLM	.000
	X-LLaMA	.000
	Typhoon	.000
	PLUG	.000
	NLLB-bridge	.000
	ChatGPT	.297
	GPT4	.063
Second Turn	PolyLM	.000
	X-LLaMA	.000
	Typhoon	.000
	PLUG	.000
	NLLB-bridge	.000
	ChatGPT	.526
	GPT4	.046

Table 19: Binomial test for Table 4.

585 Then we use the normal approximation method to calculate the CI for ratio p as

$$(p - us_p, p + us_p), \quad (5)$$

586 where u denote the critical value, for the two-tailed 95% confidence interval used in this paper
587 $u = 1.96$.

588 C.2 Significant Test

589 We conduct a two-sided binomial test for the win rate without tie $p_{\text{win}} = n_{\text{win}} / (n_{\text{win}} + n_{\text{loss}})$. The
590 null hypothesis is that the win rate is not different from the loss rate, i.e. $H_0 : p_{\text{win}} = p_{\text{loss}} = 0.5$,
591 alternative hypothesis $H_1 : p_{\text{win}} \neq 0.5$. For the test results of Table 1 and 4, please see Table 18 and
592 19. The difference between TransLLM and others in Table 5 are all significant with $p < 0.001$.

593 We conduct the χ^2 test for safety results in Table 3, the difference between TransLLM and others are
594 all significant with $p < 0.001$.

595 **NeurIPS Paper Checklist**

596 **1. Claims**

597 Question: Do the main claims made in the abstract and introduction accurately reflect the
598 paper's contributions and scope?

599 Answer: [\[Yes\]](#)

600 Justification: We have discussed our contributions and scope in detail in the Abstract and
601 Introduction chapters.

602 Guidelines:

- 603 • The answer NA means that the abstract and introduction do not include the claims
604 made in the paper.
- 605 • The abstract and/or introduction should clearly state the claims made, including the
606 contributions made in the paper and important assumptions and limitations. A No or
607 NA answer to this question will not be perceived well by the reviewers.
- 608 • The claims made should match theoretical and experimental results, and reflect how
609 much the results can be expected to generalize to other settings.
- 610 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
611 are not attained by the paper.

612 **2. Limitations**

613 Question: Does the paper discuss the limitations of the work performed by the authors?

614 Answer: [\[Yes\]](#)

615 Justification: We have discussed the limitations in Sec. 7.

616 Guidelines:

- 617 • The answer NA means that the paper has no limitation while the answer No means that
618 the paper has limitations, but those are not discussed in the paper.
- 619 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 620 • The paper should point out any strong assumptions and how robust the results are to
621 violations of these assumptions (e.g., independence assumptions, noiseless settings,
622 model well-specification, asymptotic approximations only holding locally). The authors
623 should reflect on how these assumptions might be violated in practice and what the
624 implications would be.
- 625 • The authors should reflect on the scope of the claims made, e.g., if the approach was
626 only tested on a few datasets or with a few runs. In general, empirical results often
627 depend on implicit assumptions, which should be articulated.
- 628 • The authors should reflect on the factors that influence the performance of the approach.
629 For example, a facial recognition algorithm may perform poorly when image resolution
630 is low or images are taken in low lighting. Or a speech-to-text system might not be
631 used reliably to provide closed captions for online lectures because it fails to handle
632 technical jargon.
- 633 • The authors should discuss the computational efficiency of the proposed algorithms
634 and how they scale with dataset size.
- 635 • If applicable, the authors should discuss possible limitations of their approach to
636 address problems of privacy and fairness.
- 637 • While the authors might fear that complete honesty about limitations might be used by
638 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
639 limitations that aren't acknowledged in the paper. The authors should use their best
640 judgment and recognize that individual actions in favor of transparency play an impor-
641 tant role in developing norms that preserve the integrity of the community. Reviewers
642 will be specifically instructed to not penalize honesty concerning limitations.

643 **3. Theory Assumptions and Proofs**

644 Question: For each theoretical result, does the paper provide the full set of assumptions and
645 a complete (and correct) proof?

646 Answer: [\[NA\]](#)

647 Justification: This paper does not include theoretical results.

648 Guidelines:

- 649 • The answer NA means that the paper does not include theoretical results.
- 650 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 651 referenced.
- 652 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 653 • The proofs can either appear in the main paper or the supplemental material, but if
- 654 they appear in the supplemental material, the authors are encouraged to provide a short
- 655 proof sketch to provide intuition.
- 656 • Inversely, any informal proof provided in the core of the paper should be complemented
- 657 by formal proofs provided in appendix or supplemental material.
- 658 • Theorems and Lemmas that the proof relies upon should be properly referenced.

659 4. Experimental Result Reproducibility

660 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

661 perimental results of the paper to the extent that it affects the main claims and/or conclusions

662 of the paper (regardless of whether the code and data are provided or not)?

663 Answer: [Yes]

664 Justification: We have tried our best to provide the details of our experiments in Sec. 4.1 and

665 Appendix A. We also provided our code and datasets in supplementary materials.

666 Guidelines:

- 667 • The answer NA means that the paper does not include experiments.
- 668 • If the paper includes experiments, a No answer to this question will not be perceived
- 669 well by the reviewers: Making the paper reproducible is important, regardless of
- 670 whether the code and data are provided or not.
- 671 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 672 to make their results reproducible or verifiable.
- 673 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 674 For example, if the contribution is a novel architecture, describing the architecture fully
- 675 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 676 be necessary to either make it possible for others to replicate the model with the same
- 677 dataset, or provide access to the model. In general, releasing code and data is often
- 678 one good way to accomplish this, but reproducibility can also be provided via detailed
- 679 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 680 of a large language model), releasing of a model checkpoint, or other means that are
- 681 appropriate to the research performed.
- 682 • While NeurIPS does not require releasing code, the conference does require all submis-
- 683 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 684 nature of the contribution. For example
 - 685 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 686 to reproduce that algorithm.
 - 687 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 688 the architecture clearly and fully.
 - 689 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 690 either be a way to access this model for reproducing the results or a way to reproduce
 - 691 the model (e.g., with an open-source dataset or instructions for how to construct
 - 692 the dataset).
 - 693 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 694 authors are welcome to describe the particular way they provide for reproducibility.
 - 695 In the case of closed-source models, it may be that access to the model is limited in
 - 696 some way (e.g., to registered users), but it should be possible for other researchers
 - 697 to have some path to reproducing or verifying the results.

698 5. Open access to data and code

699 Question: Does the paper provide open access to the data and code, with sufficient instruc-

700 tions to faithfully reproduce the main experimental results, as described in supplemental

701 material?

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
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753

Answer: [Yes]

Justification: We have provided open access to the data and code, with sufficient instructions provided to faithfully reproduce the main experimental results. Please refer the README file in code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have tried our best to provide the details of our experiments in Sec. 4.1 and Appendix A.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Please refer Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- 754 • The assumptions made should be given (e.g., Normally distributed errors).
- 755 • It should be clear whether the error bar is the standard deviation or the standard error
- 756 of the mean.
- 757 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 758 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 759 of Normality of errors is not verified.
- 760 • For asymmetric distributions, the authors should be careful not to show in tables or
- 761 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 762 error rates).
- 763 • If error bars are reported in tables or plots, The authors should explain in the text how
- 764 they were calculated and reference the corresponding figures or tables in the text.

765 8. Experiments Compute Resources

766 Question: For each experiment, does the paper provide sufficient information on the com-
 767 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 768 the experiments?

769 Answer: [Yes]

770 Justification: Please refer Appendix A.2.

771 Guidelines:

- 772 • The answer NA means that the paper does not include experiments.
- 773 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 774 or cloud provider, including relevant memory and storage.
- 775 • The paper should provide the amount of compute required for each of the individual
- 776 experimental runs as well as estimate the total compute.
- 777 • The paper should disclose whether the full research project required more compute
- 778 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 779 didn't make it into the paper).

780 9. Code Of Ethics

781 Question: Does the research conducted in the paper conform, in every respect, with the
 782 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

783 Answer: [Yes]

784 Justification: We have carefully checked our paper. Our paper conforms, in every respect,
 785 with the NeurIPS Code of Ethics.

786 Guidelines:

- 787 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 788 • If the authors answer No, they should explain the special circumstances that require a
- 789 deviation from the Code of Ethics.
- 790 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 791 eration due to laws or regulations in their jurisdiction).

792 10. Broader Impacts

793 Question: Does the paper discuss both potential positive societal impacts and negative
 794 societal impacts of the work performed?

795 Answer: [Yes]

796 Justification: We have discussed the potential positive societal impacts in Sec. 7, and it
 797 seems that this work does not exert obviously negative societal impacts.

798 Guidelines:

- 799 • The answer NA means that there is no societal impact of the work performed.
- 800 • If the authors answer NA or No, they should explain why their work has no societal
 801 impact or why the paper does not address societal impact.
- 802 • Examples of negative societal impacts include potential malicious or unintended uses
 803 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 804 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 805 groups), privacy considerations, and security considerations.

- 806 • The conference expects that many papers will be foundational research and not tied
807 to particular applications, let alone deployments. However, if there is a direct path to
808 any negative applications, the authors should point it out. For example, it is legitimate
809 to point out that an improvement in the quality of generative models could be used to
810 generate deepfakes for disinformation. On the other hand, it is not needed to point out
811 that a generic algorithm for optimizing neural networks could enable people to train
812 models that generate Deepfakes faster.
- 813 • The authors should consider possible harms that could arise when the technology is
814 being used as intended and functioning correctly, harms that could arise when the
815 technology is being used as intended but gives incorrect results, and harms following
816 from (intentional or unintentional) misuse of the technology.
- 817 • If there are negative societal impacts, the authors could also discuss possible mitigation
818 strategies (e.g., gated release of models, providing defenses in addition to attacks,
819 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
820 feedback over time, improving the efficiency and accessibility of ML).

821 11. Safeguards

822 Question: Does the paper describe safeguards that have been put in place for responsible
823 release of data or models that have a high risk for misuse (e.g., pretrained language models,
824 image generators, or scraped datasets)?

825 Answer: [Yes]

826 Justification: Our work aims to transfer the safeguards of chat large language models from
827 English to non-English.

828 Guidelines:

- 829 • The answer NA means that the paper poses no such risks.
- 830 • Released models that have a high risk for misuse or dual-use should be released with
831 necessary safeguards to allow for controlled use of the model, for example by requiring
832 that users adhere to usage guidelines or restrictions to access the model or implementing
833 safety filters.
- 834 • Datasets that have been scraped from the Internet could pose safety risks. The authors
835 should describe how they avoided releasing unsafe images.
- 836 • We recognize that providing effective safeguards is challenging, and many papers do
837 not require this, but we encourage authors to take this into account and make a best
838 faith effort.

839 12. Licenses for existing assets

840 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
841 the paper, properly credited and are the license and terms of use explicitly mentioned and
842 properly respected?

843 Answer: [Yes]

844 Justification: Please refer Appendix A.5.

845 Guidelines:

- 846 • The answer NA means that the paper does not use existing assets.
- 847 • The authors should cite the original paper that produced the code package or dataset.
- 848 • The authors should state which version of the asset is used and, if possible, include a
849 URL.
- 850 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 851 • For scraped data from a particular source (e.g., website), the copyright and terms of
852 service of that source should be provided.
- 853 • If assets are released, the license, copyright information, and terms of use in the
854 package should be provided. For popular datasets, paperswithcode.com/datasets
855 has curated licenses for some datasets. Their licensing guide can help determine the
856 license of a dataset.
- 857 • For existing datasets that are re-packaged, both the original license and the license of
858 the derived asset (if it has changed) should be provided.

859 • If this information is not available online, the authors are encouraged to reach out to
860 the asset’s creators.

861 13. **New Assets**

862 Question: Are new assets introduced in the paper well documented and is the documentation
863 provided alongside the assets?

864 Answer: [Yes]

865 Justification: Please refer the README file in code.

866 Guidelines:

- 867 • The answer NA means that the paper does not release new assets.
- 868 • Researchers should communicate the details of the dataset/code/model as part of their
869 submissions via structured templates. This includes details about training, license,
870 limitations, etc.
- 871 • The paper should discuss whether and how consent was obtained from people whose
872 asset is used.
- 873 • At submission time, remember to anonymize your assets (if applicable). You can either
874 create an anonymized URL or include an anonymized zip file.

875 14. **Crowdsourcing and Research with Human Subjects**

876 Question: For crowdsourcing experiments and research with human subjects, does the paper
877 include the full text of instructions given to participants and screenshots, if applicable, as
878 well as details about compensation (if any)?

879 Answer: [NA]

880 Justification: This paper does not involve crowdsourcing nor research with human subjects.

881 Guidelines:

- 882 • The answer NA means that the paper does not involve crowdsourcing nor research with
883 human subjects.
- 884 • Including this information in the supplemental material is fine, but if the main contribu-
885 tion of the paper involves human subjects, then as much detail as possible should be
886 included in the main paper.
- 887 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
888 or other labor should be paid at least the minimum wage in the country of the data
889 collector.

890 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 891 Subjects**

892 Question: Does the paper describe potential risks incurred by study participants, whether
893 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
894 approvals (or an equivalent approval/review based on the requirements of your country or
895 institution) were obtained?

896 Answer: [NA]

897 Justification: This paper does not involve crowdsourcing nor research with human subjects.

898 Guidelines:

- 899 • The answer NA means that the paper does not involve crowdsourcing nor research with
900 human subjects.
- 901 • Depending on the country in which research is conducted, IRB approval (or equivalent)
902 may be required for any human subjects research. If you obtained IRB approval, you
903 should clearly state this in the paper.
- 904 • We recognize that the procedures for this may vary significantly between institutions
905 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
906 guidelines for their institution.
- 907 • For initial submissions, do not include any information that would break anonymity (if
908 applicable), such as the institution conducting the review.