# Less Is More: Distilling Large-Scale Data with LLMs for Chinese-Centric **Low-Resource Multilingual Machine Translation**

**Anonymous ACL submission** 

### Abstract

Neural Machine Translation (NMT) between Chinese and low-resource languages (LRLs) faces significant challenges due to limited, noisy training data. We introduce MERIT, a unified translation framework that transforms the traditional English-centric ALT benchmark into a Chinese-centric evaluation suite for five Southeast Asian LRLs. Our approach 010 integrates Language-specific Token Prefixing (LTP) for effective language conditioning and supervised fine-tuning (SFT). A key innovation, Group Relative Policy Optimization (GRPO) guided by the Score Accuracy Reward (SAR) function, strategically filters training data and optimizes model performance. Experiments with models up to 3 billion parameters (MERIT-3B) confirm the efficacy of our method. Abla-019 tion studies demonstrate substantial improvements from SFT-LTP over zero-shot baselines, while GRPO-SAR achieves further significant gains using only 22.8% of the original data, increasing BLEU-chrF scores by 17.4%. MERIT-3B notably surpasses open-source models such as NLLB-200 3.3B by 9.5 BLEU-4 points on Chinese-Indonesian translation and outperforms M2M-100 by 5.1 BLEU-4 points on Chinese-Lao. These findings highlight the pivotal role of targeted data curation and reward-guided training over mere model scaling, advancing multilingual translation in lowresource settings. Code and data are available at https://anonymous.4open.science/r/MERIT-864/.

#### 1 Introduction

011

012

The vision of Neural Machine Translation (NMT) is to provide equitable access to information for speakers of over 7,000 languages worldwide. While English-French translation has achieved near-human BLEU scores (Papineni et al., 2002), 040 many low-resource languages-including Chinese-LRL directions such as Tibetan, Lao, or Tagalog-remain virtually untranslatable due to the lack 043

of parallel corpora, standard orthographies, and annotated data (Costa-jussà et al., 2022).

044

045

047

050

051

059

060

061

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

079

081

Multilingual pretraining has shown impressive zero-shot gains: mBART-50, mT5, and DeltaLM achieve broad coverage but continue to overlook or underperform on Southeast Asian and Chinese domestic LRLs. NLLB-200 (Costa-jussà et al., 2022) improves on this by expanding coverage to 200 languages and introducing the XSTS metric. However, Chinese–LRL performance still trails behind English-pivoted directions.

Moreover, the lack of publicly available and high-quality evaluation benchmarks hinders objective progress measurement. An ideal benchmark should: (i) cover multiple Chinese -> LRL directions, (ii) maintain sufficient balance in scale and domain coverage, and (iii) avoid dependence on English as a pivot. Without these features, model improvements are difficult to reproduce or attribute reliably.

To address the long-standing lack of evaluation benchmarks for Chinese-low-resource language (LRL) directions in Neural Machine Translation (NMT), this paper makes the following three contributions:

- We introduce the first Chinese-centric multilingual benchmark for low-resource languages, derived from the ASEAN Languages Treebank (ALT) through various data filtering (EPDS-DIV) and distillation (QE Agent) techniques. This dataset targets low-resource language scenarios and ensures balanced domain coverage and semantic consistency (Thu et al., 2016).
- We compare multiple metrics through Strict-Overlap and Semantic-Friendly measures on five series of large language models (LLMs). The evaluation is conducted using three approaches: zero-shot, supervised fine-tuning

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

(SFT), and Group Relative Policy Optimization (GRPO), with the latter utilizing the expert-rated validation set for training.

• We explore the impact of model scale, data adaptation, and reward-based quality filtering on Chinese–LRL translation performance. Our results show that smaller open-source models can match proprietary models when trained with high-quality, strictly filtered data and guided reinforcement. The proposed MERIT framework effectively integrates and leverages these components.

## 2 Related Work

090

100

103

104

**Early Chinese–LRL Corpora.** The CCMT shared tasks released fewer than 200k sentence pairs for ZH–UG and ZH–MN (Liu et al., 2021), while Wiki-based mining typically yields only a few thousand pairs for ZH–LO and ZH–FIL (Artetxe and Schwenk, 2019). The ALT corpus (Thu et al., 2016) extends coverage to 13 ASEAN languages but remains English-centric and lacks direct Chinese–LRL alignment.

Multilingual Pretraining. Models like mBART-105 50 (Liu et al., 2020), mT5 (Xue et al., 2021), and 107 DeltaLM (Ma et al., 2021) cover 50–101 languages, but still overlook or underperform on languages 108 such as Tibetan and Uyghur. NLLB-200 (Costa-109 jussà et al., 2022) improves BLEU by 44% on 110 FLORES-200 and adds the XSTS metric, but still 111 underperforms on Chinese->LRL due to data qual-112 ity gaps. 113

LLM-based Machine Translation. Instruction-114 tuned LLMs such as GPT-40 (Huang et al., 2025), 115 Claude-3.5 (Enis and Hopkins, 2024), and Gemini-116 2.5 (DeepMind, 2025) show strong generalization 117 across many languages. Open-source models like 118 Qwen-2.5 (Cui and et al., 2025) and DeepSeek-v3 119 (Huang et al., 2025) serve as strong multilingual 120 baselines, though most published evaluations fo-121 cus on FLORES-200 and overlook Chinese-LRL 122 directions. 123

Evaluation Limitations. Most multilingual
benchmarks pivot through English, mix domain
content, or lack human validation, reducing their
diagnostic value for Chinesee-LRL MT. Our
benchmark addresses these gaps with direct
Chinese–LRL alignment, domain balance, and
human-validated samples.

## 3 Methodology

### 3.1 Dataset

We construct a new test suite based on the ASEAN Languages Treebank (ALT) corpus (Thu et al., 2016). ALT is an English-centric multilingual corpus that already provides sentence-level alignment for several Southeast-Asian languages–Vietnamese (vi), Burmese (my), Lao (10), Tagalog (fil) and Indonesian (id). Although Chinese is included as a target aligned with English, no direct Chinese–LRL alignment exists.

We therefore re-index sentences sharing the same alt\_id and semantic source to form direct Chinese–LRL sentence pairs. In the resulting test set, Chinese can serve either as the source or as the reference language. The benchmark is clean, stylistically consistent and typologically diverse, enabling more controlled and fair evaluation of multilingual NMT systems from a Chinese-centric perspective.

**Language Selection.** We deliberately focus on five Southeast-Asian low-resource languages (vi, my, lo, fil, id) for four data-driven reasons:

All five appear in ALT with reliable English alignments, so that high-quality Chinese–LRL realignment is feasible.

Indonesian rather than Malay is kept because the two belong to the same Malayic subgroup and share 90 % lexical overlap, to the extent that many international surveys treat them as a single "Malay macrolanguage" (Adelaar, 2012; Eberhard et al., 2023b), if both are included simultaneously, it will result in redundancy of the experimental languages. Malay already has over 1 million clean En–Ms

Languages	<b>Speaker Population</b> <sup>1</sup> (Eberhard et al., 2023a)	Filtered Subset <sup>2</sup>	LRL
Chinese (zh)	1180M	X	×
Hindi (hi)	345M	×	X
Bengali (bn)	234M	×	X
Japanese (ja)	121M	×	×
Vietnamese (vi)	86M	10K	1
Indonesian (id)	43M	10K	1
Burmese (my)	33M	10K	1
Tagalog (fil)	24M	10K	<ul> <li>Image: A second s</li></ul>
Thai (th)	20M	×	×
Malay (ms)	18M	×	×
Khmer (km) <sup>3</sup>	16M	_	1
Lao (lo)	4.3M	10K	1

Table 1: Language Statistics from the ALT Corpus for Chinese-Centric Multilingual Translation. ALT corpus statistics sorted by L1 speaker population. All counts refer to L1 speakers and are rounded to the nearest million (M). LRL: Low-resource Language.

210

211

212

213

214

sentence pairs (e.g., MT-Wiki and multiple OPUS
sub-corpora), pushing it into the mid-resource tier
(Duong et al., 2017), whereas Indonesian still lacks
sizeable Chinese parallel data (<50 k pairs in total,</li>
ALT contributes only 20 k) and is classified as lowresource by FLORES-200 and NLLB benchmarks
(Goyal and et al., 2022; Team and AI, 2022).

172

173

174

175

176

178

179

183

184

185

188

189

190

191

192

193

194

196

198

199

201

207

Thai aggregated corpora exceed one million sentence pairs and the language enjoys dedicated WMT/IWSLT tracks (Lowphansirikul and Chuangsuwanich, 2020), so it no longer fits a strict LRL definition.

Khmer parallel resources are both small and highly noisy, the WMT20 corpus-filtering task emphasised that extensive cleaning is required (Koehn et al., 2020); moreover, unlike the other languages considered, Ethnologue reports virtually no L2 speaker community for Khmer (Simons and Fennig, 2023). Including Khmer would therefore demand a language-specific filtering pipeline and would undermine comparability with the other languages.

This reconstructed benchmark complements existing resources such as FLORES-200 (Costa-jussà et al., 2022), particularly for Chinese–LRL directions in mainland and maritime Southeast Asia. Unlike pivot-based benchmarks, our test set avoids semantic distortion introduced by intermediate English, thus enabling more realistic, stable, and reproducible evaluation for Chinese-centric multilingual translation systems.

### 3.2 Model Overview

We evaluate five series representative LLMs, spanning both proprietary and open-source systems:

**Qwen-2.5** (Ghosal et al., 2024; Cui and et al., 2025): Chinese-English bilingual models finetuned for multilingual transfer, evaluated on several LRLs.

**GPT-40** (Huang et al., 2025): OpenAI's flagship model tested on 16 languages, including several low-resource directions such as En–Te and En–Sw.

**Claude-3.5** (Enis and Hopkins, 2024): A multilingual LLM from Anthropic, evaluated via MQM metrics on pairs like En–Yoruba and En–Amharic.

**Gemini-2.5** (DeepMind, 2025): While lacking peer-reviewed benchmarks on Chinese–LRL tasks,

its predecessor covers ultra-low-resource translation (e.g., En–Kalamang).

**DeepSeek** (Huang et al., 2025; Jiang et al., 2025): A competitive open-source model evaluated in the BenchMAX suite alongside GPT-40.

Zero-shot prompting was applied exclusively to closed-source LLMs. For the Qwen-2.5 models, both SFT and GRPO-enhanced SFT were applied. Additional tests were conducted with SFT and GRPO-enhanced SFT using enhanced data derived from closed-source LLMs in the zero-shot regime. The GRPO-enhanced SFT regime utilizes a scoring agent trained on expert-rated development sets, optimized with Score Accuracy Reward (SAR) to ensure the selection of only high-quality translations. Model performance is assessed using the following metrics: BLEU-4, sacreBLEU, chrF, ROUGE-L, METEOR, and BERTScore.

### 3.3 Scoring and Selection

To construct high-quality parallel corpora for lowresource translation, we design a two-stage scoring and filtering pipeline that integrates interpretable statistical features with semantic evaluation, followed by reference-free quality estimation and threshold-based selection.

**Stage I: Statistical and Semantic Scoring.** We extract key surface-level features such as sentence length ratio and digit proportion difference, along with aggregated indicators for token balance, punctuation consistency, and lexical diversity. These help identify formatting mismatches and alignment noise (Munteanu and Marcu, 2005; Sánchez-Cartagena et al., 2018).

To assess deeper semantic alignment, we incorporate two additional signals: (i) conditional perplexity for fluency estimation, and (ii) instructionfollowing discrepancy to capture semantic fidelity, inspired by instruction-tuning objectives (Li et al., 2023). All features are normalized and combined through weighted scoring to penalize semantically misaligned pairs (Esplà-Gomis et al., 2020).

**Stage II: Quality Estimation and Filtering.** A reference-free Quality Estimation (QE) model, trained on human-annotated validation sets, further evaluates translation adequacy and fluency (Rei et al., 2020; Freitag et al., 2021). Sentence pairs surpassing a calibrated threshold are retained. This final stage ensures that the resulting dataset is both scalable and high-quality, suitable for fine-tuning compact models in low-resource scenarios.

<sup>&</sup>lt;sup>1</sup>Speaker numbers derive from the most recent national censuses or *Ethnologue* reports (2023–2025) and are expressed in millions (M).

<sup>&</sup>lt;sup>2</sup>Each ALT language contains approximately 20k aligned sentence pairs from a shared English source. See https:// www2.nict.go.jp/astrec-att/member/mutiyama/ALT/.

<sup>&</sup>lt;sup>3</sup>No L2 speaker community (Simons and Fennig, 2023).



Figure 1: The overall workflow of the MERIT framework for Chinese-centric multilingual translation. (a) Data Selection: Sentences from the ALT corpus are filtered and scored using perplexity (ppl) and inverse frequency diversity (ifd) to construct and refine a high-quality multilingual dataset, yielding training and testing subsets for five low-resource Southeast Asian languages. (b) Translation Scoring: A Translation Quality Estimation (QE) dataset is created by language experts, followed by training a QE agent to assess the translation quality of additional data. (c) Data Distilling: QE Agent trained by Group Relative Policy Optimization (GRPO) with Score Accuracy Reward (SAR) leveraged to further evaluate more training data, yielding more high-quality training samples. (d) Model Training: The MERIT-3B model is trained using Supervised Fine-Tuning (SFT), Language-specific Token Prefixing (LTP) with LLM-Enhanced data argmentation to achieve optimal multilingual translation performance.

### 3.4 Supervised Fine-Tuning

We fine-tune open-source models (Qwen-2.5 0.5B and 3B) on the filtered Chinese–LRL data using supervised fine-tuning (SFT) (Fan et al., 2021). The training objective is to maximize the conditional likelihood of the target sequence  $Y = (y_1, \ldots, y_M)$  given the source sequence  $X = (x_1, \ldots, x_N)$ , using standard sequence-to-sequence formulation with teacher forcing (Sutskever et al., 2014; Williams and Zipser, 1989).

Label smoothing with  $\varepsilon = 0.1$  is applied to avoid over-confidence (Szegedy et al., 2016). Fine-tuning is conducted independently for each Chinese–LRL pair to account for language-specific morphology. This strategy, paired with prior quality filtering, yields strong gains over zero-shot performance, echoing findings on targeted adaptation for multilingual NMT (Arivazhagan et al., 2019).

### 279 **3.5** Language-specific Token Prefixing

To improve language discrimination in one-tomany generation, we adopt Language-specific Token Prefixing (LTP), which prepends a target language token (e.g., [10]) to both the source input and prompt instruction. This token is added to the tokenizer vocabulary and embedded as part of the model input. 283

284

285

287

290

291

292

293

294

296

297

298

299

301

For each training sample, the source input is modified as  $X' = [lang] \oplus X$ , and the final model input becomes a concatenation of the instruction prompt and the language-tagged source sequence:

 $Input = [Prompt \oplus [lang] \oplus x_1, \dots, x_n] \quad (1)$ 

The training objective minimizes the negative log-likelihood of the target sequence:

$$\mathcal{L}_{\text{MLE}} = -\sum_{t=1}^{m} \log P(y_t \mid y_{< t}, \text{Prompt}, X; \theta)$$
(2)

This extends target-language prefixing ideas (Johnson et al., 2017) by combining symbolic and prompt-based conditioning for unified multilingual fine-tuning.

## 3.6 Group Relative Policy Optimization

Group Relative Policy Optimization (GRPO) is a reinforcement learning strategy that refines model

260

261

262

263

264

265

271

272

276

278

outputs using reward feedback. Inspired by Reinforcement Learning with Human Feedback (RLHF) techniques (Ouyang et al., 2022; Lu et al., 2022), GRPO operates on mini-batches of candidate translations in this task, assigning scalar rewards based on Score Accuracy Reward (SAR) scores. Subsequently, the model learns to maximize the expected reward via policy gradient updates.

302

303

304

307

310

311

313

314

315

316

321

323

327

329

332

333

334

336

339

341

343

345

347

348

Unlike conventional pointwise objective functions, GRPO introduces an intra-batch comparison mechanism and normalizes rewards using a moving baseline. This approach helps to reduce the variance of gradient estimates, thereby enhancing the stability of the training process. Our experiments indicate that GRPO is effective for translation evaluation, enabling the selection and improvement of dataset translation quality.

#### **Score Accuracy Reward Function** 3.7

We define the Score Accuracy Reward (SAR) to evaluate model completions based on their ability to accurately reproduce specific numerical scores present in ground-truth answers. This reward function is designed for tasks where precision in extracting or generating key numerical information is critical. SAR rewards model outputs that closely match these target numerical values (representing a form of preference or correctness) while penalizing deviations.

To extract the salient numerical score  $s_i$  from the completion content  $c_i$ , we first delineate a matchset,  $M(c_i)$ . This set comprises all integer values identified within  $c_i$  through the application of a predefined regular expression, denoted as R. We define a predicate  $P_R(m, c_i)$  to be true if and only if m is an integer yielded by matching the regular expression R against the string  $c_i$ . The match-set  $M(c_i)$  is then formally defined as:

$$M(c_i) = \{ m \in \mathbb{Z} \mid P_R(m, c_i) \}$$
(3)

The extractor function  $E(c_i)$  then determines the score  $s_i$  from this match-set. As per the reference, if  $M(c_i)$  is not empty,  $s_i$  is the minimum integer found; otherwise,  $s_i$  is set to -1:

$$s_i = E(c_i) = \begin{cases} \min M(c_i), & \text{if } M(c_i) \neq \emptyset \\ -1, & \text{if } M(c_i) = \emptyset \end{cases}$$
(4)

Given a ground-truth integer answer vector a = $(a_1, a_2, \ldots, a_N)$ , we define a piecewise rewardmapping function  $\phi(d)$  based on the absolute difference  $d = |s_i - a_i|$  between the extracted score

 $s_i$  and the ground-truth answer  $a_i$ . This function, detailed in source, is:

$$\phi(d) = \begin{cases} 2.0, & \text{if } d = 0\\ 1.0, & \text{if } 1 \le d \le 10\\ 0.0, & \text{otherwise} \end{cases}$$
(5)

This mapping assigns the highest reward for an exact match, a partial reward for close matches (difference up to 10), and zero reward for larger deviations or mismatches.

Finally, the reward  $r_i$  for each instance *i* is computed based on  $s_i$  and  $\phi(d)$ . If a valid score  $s_i \ge 0$ was extracted, the reward is  $\phi(|s_i - a_i|)$ . If no score was extracted ( $s_i < 0$ ), the reward is 0.0:

$$r_i = \begin{cases} \phi(|s_i - a_i|), & \text{if } s_i \ge 0\\ 0.0, & \text{if } s_i < 0 \end{cases}$$
(6)

The overall outcome is a reward vector r = $(r_1, r_2, \ldots, r_N)$ . This SAR mechanism, by focusing on the accuracy of extracted numerical scores against ground-truth values, provides a clear signal for tasks requiring numerical precision. Similar score-alignment objectives, where models are rewarded for matching target scores or preferences, have been successfully adopted in alignment training for various generation tasks (Wu et al., 2023).

#### 4 **Experiments and Analysis**

#### **Evaluation Method** 4.1

We evaluate translation performance using both overlap-based and semantic-aware metrics:

Strict-Overlap: BLEU-4 (Papineni et al., 2002), sacreBLEU (Post, 2018), and ROUGE-L (Lin, 2004) assess lexical match and n-gram precision, which are crucial for evaluating surface-level accuracy and fluency.

Semantic-Friendly: chrF (Popovic, 2015), ME-TEOR (Banerjee and Lavie, 2005), and BERTScore (Zhang et al., 2020) measure semantic similarity and fluency robustness, capturing aspects that ngram overlap alone might miss.

Each metric is computed on the reconstructed ALT test suite for five Chinese-LRL pairs. We report averages across directions, comparing zero-shot prompting, SFT, and GRPO-enhanced regimes.

To provide a balanced evaluation that captures both lexical precision and semantic adequacy, we propose a composite metric, BLEU-chrF. This metric integrates insights from both the Strict-Overlap

351

353

354

355

349

350

356 357

359

- 361
- 363
- 366 367 368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

387

388

389

390

- 364



467

468

469

470

471

472

473

474

475



Figure 2: Performance–Scale Trade-offs of MERIT-3B and Baseline Models on Chinese-Centric Multilingual Translation. Comparison of BLEU-chrF scores against model size (log-scale) across MERIT-3B, open-source, and estimated closed-source models.

and Semantic-Friendly categories of evaluation measures by taking the arithmetic mean of the BLEU-4 score and the chrF score:

$$BLEU-chrF = \frac{BLEU-4 + chrF}{2}$$
(7)

By averaging these two widely-used metrics, one emphasizing n-gram precision and the other character n-gram recall and F-score. We aim to achieve a more holistic assessment of translation quality, particularly for tasks where both lexical fidelity and semantic resemblance are important.

### 4.2 Main Result

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419 420

421

422

423

424

Table 2 presents our evaluation results across five Chinese–LRL directions, categorizing metrics into Strict-Overlap (Papineni et al., 2002; Post, 2018,?; Lin, 2004) and Semantic-Friendly (Popovic, 2015; Banerjee and Lavie, 2005; Zhang et al., 2020).

Among leading closed-source models, Gemini-2.5 Flash consistently achieves top scores in BLEU-4 and chrF across multiple languages, such as Filipino (BLEU-4: 49.26, chrF: 42.68) and Indonesian (BLEU-4: 48.73, chrF: 41.93). Claude-3.5 Sonnet excels in ROUGE-L for Lao (14.00) and Burmese (24.44).

The proposed Multilingual Expert-Reward Informed Tuning (MERIT) framework, demonstrates notable strengths. Specifically, MERIT-3B significantly outperforms the similarly sized opensource NLLB-200 3.3B model across several metrics for Filipino, Indonesian, and Vietnamese. For instance, on Filipino, MERIT-3B achieves 29.71 BLEU-4 and 49.88 METEOR compared to NLLB-200's 25.05 and 46.10, respectively. Furthermore, MERIT-3B shows substantial gains over smaller open-source baselines on particularly challenging low-resource language pairs.

Notably, our MERIT-3B model demonstrates substantial advantages over the DeepSeek-r1 7B. MERIT-3B consistently outperforms DeepSeek-r1 7B on lexical similarity metrics such as BLEU-4 and chrF across all five evaluated languages. Furthermore, when benchmarked against Qwen-2.5 7B, MERIT-3B, with only approximately 42.9% of its parameters, achieves highly competitive translation quality. For instance, on Filipino, MERIT-3B reaches 99.7% of Qwen-2.5 7B's ROUGE-L score (31.91 vs. 31.99) and over 98.3% of its BLEU-4 score (29.71 vs. 30.21). Similar competitiveness is observed for Indonesian, where MERIT-3B attains approximately 95.7% of Qwen-2.5 7B's BLEU-4 score (34.73 vs. 36.28) and 93.2% of its ROUGE-L score (26.63 vs. 28.56).

These results underscore the efficacy of our reward-informed filtering and specialized finetuning approach, particularly in enhancing performance for low-resource languages and achieving competitive results within the open-source landscape relative to model scale.

### 4.3 Module Comparison

To investigate the contribution of each component, we conduct an ablation study on Qwen-2.5 0.5B and Qwen-2.5 3B across four distinct setups: zeroshot (serving as our baseline), Supervised Fine-Tuning with Language-Token Prefixing (SFT-LTP), SFT-LTP followed by reward-enhanced tuning using Group Relative Policy Optimization with Score Accuracy Reward (GRPO-SAR), and finally our full SFT-LTP + GRPO-SAR with an additional LLMs-Enhanced (LLME) stage.

As detailed in Table 3, initial SFT-LTP yields substantial improvements in both BLEU-4 and chrF scores over the zero-shot baselines across all languages for both model sizes. For instance, Qwen-2.5 3B sees its overall BLEU-chrF score increase from 9.10 to 15.53 after SFT-LTP. Introducing GRPO-SAR provides further consistent gains. Notably, for the Qwen-2.5 3B model, GRPO-SAR significantly boosts performance on low-resource pairs like Chinese–Lao, improving BLEU-4 from 1.17 (SFT-LTP) to 4.39 and chrF from 2.22 (SFT-LTP) to 5.17. Even with its limited capacity, the Qwen-2.5 0.5B model benefits remarkably from GRPO-SAR, achieving an overall BLEU-chrF score of 4.39, which is nearly a 40-fold increase

Strict-Overlap	BLEU-4				sacreBLEU				ROUGE-L								
Model	fil	id	lo	my	vi	fil	id	lo	my	vi	fil	id	lo	my	vi	FT	OS
GPT-40 Claude-3.5 Sonnet Gemini-2.5 Flash DeepSeek-v3	45.26 42.97 <b>49.26</b> <u>46.19</u>	47.77 47.35 <b>48.73</b> 41.83	28.10 <u>32.20</u> <b>35.79</b> 25.57	30.48 <u>30.92</u> <b>36.91</b> 26.68	$   \begin{array}{r}     45.03 \\     \underline{45.14} \\     \overline{39.24} \\     41.29   \end{array} $	43.60 43.38 <b>47.48</b> <u>44.12</u>	45.81 46.54 <b>47.20</b> 41.38	27.38 <u>32.65</u> <b>35.03</b> 25.25	29.53 <u>30.14</u> <b>36.07</b> 26.29	$ \begin{array}{r}     44.50 \\     \underline{44.55} \\     \overline{38.63} \\     40.90 \end{array} $	33.62 35.03 <b>35.89</b> <u>35.15</u>	<u>31.40</u> 31.17 30.83 30.06	12.33 <b>14.00</b> <u>13.53</u> <u>12.75</u>	23.57 24.44 23.92 22.95	28.25 28.45 24.62 28.27	× × × ×	× × × ✓
Qwen-2.5 32B DeepSeek-r1 32B Qwen 2.5-7B DeepSeek-r1 7B NLLB-200 3.3B DeepSeek-r1 1.5B M2M-100 1.2B <b>MERIT-3B (Ours)</b>	43.56 37.98 30.21 14.77 25.05 0.07 2.13 29.71	$\frac{47.87}{43.29}$ 36.28 20.94 25.27 0.08 9.53 34.73	23.27 12.97 6.17 0.79 15.86 0.05 0.05 5.15	20.43 8.87 6.43 0.37 20.83 1.06 0.00 4.56	<b>46.23</b> 42.57 35.22 12.53 25.30 0.05 3.33 31.20	42.01 37.14 29.75 15.31 24.21 0.03 1.32 27.20	$\frac{46.61}{42.27}$ 36.41 24.17 23.64 0.06 9.47 33.16	22.48 12.71 5.67 0.86 15.19 0.01 0.01 4.28	19.50 8.43 5.42 0.46 20.13 0.11 0.00 3.54	<b>44.63</b> 41.45 35.07 16.16 22.97 0.02 3.19 29.25	34.97 32.40 31.99 24.58 31.45 0.77 9.88 31.91	<b>31.85</b> 29.92 28.56 23.32 25.73 0.90 21.65 26.63	11.84 10.04 9.86 8.67 11.49 1.80 4.69 8.40	20.61 16.31 16.41 5.67 18.68 3.80 0.00 13.68	<b>29.08</b> 28.11 27.81 18.52 25.51 0.55 13.01 21.18	* * * * * * *	******
Semantic-Friendly	chrF			METEOR				BERTScore									
Model	fil	id	lo	my	vi	fil	id	lo	my	vi	fil	id	lo	my	vi	FT	OS
GPT-40 Claude-3.5 Sonnet Gemini-2.5 Flash DeepSeek-v3	39.36 38.71 <b>42.68</b> <u>39.80</u>	41.13 41.30 <b>41.93</b> 36.95	24.20 28.38 <b>30.61</b> 22.37	26.76 28.79 <b>32.09</b> 24.43	39.62 <u>39.65</u> 34.66 36.86	67.80 67.52 <b>70.14</b> <u>68.04</u>	70.34 70.29 <b>70.87</b> 67.10	53.66 58.53 <b>60.21</b> 51.41	56.59 <u>58.88</u> <b>61.85</b> 53.66	<u>69.44</u> 69.23 60.31 67.43	68.59 68.28 <b>71.67</b> <u>70.02</u>	70.84 71.55 <b>72.77</b> 68.81	56.27 60.72 63.93 54.82	57.04 <u>57.68</u> <b>62.39</b> 54.21	$\frac{70.41}{70.40}\\60.88\\68.87$	× × × ×	× × × ✓
Qwen-2.5 32B DeepSeek-r1 32B Qwen-2.5 7B DeepSeek-r1 7B NLLB-200 3.3B DeepSeek-r1 1.5B M2M-100 1.2B <b>MERIT 3B (Ours)</b>	37.77           33.62           27.88           15.43           22.48           0.33           5.16           25.52	41.35 37.62 33.68 21.60 23.57 0.38 18.21 30.22	20.36 12.77 7.39 2.20 14.92 0.36 0.26 5.81	18.13 9.81 7.95 1.35 18.69 1.87 0.01 5.53	<b>39.80</b> 36.97 33.28 14.87 22.43 0.24 6.49 26.98	66.61 61.60 54.51 36.02 46.10 1.99 4.65 49.88	$\frac{70.57}{66.75}$ $\frac{62.56}{49.35}$ $\frac{45.28}{2.34}$ $\frac{14.00}{56.46}$	46.72 34.07 20.84 8.45 36.27 2.17 0.57 16.55	43.29 27.37 23.25 6.46 42.28 3.80 0.11 16.93	<b>69.76</b> 66.95 63.00 37.61 45.73 1.43 6.10 50.50	67.73 63.32 54.95 34.22 48.61 -27.97 -16.71 46.70	71.73 68.54 62.82 49.07 54.34 -26.07 -1.68 45.58	50.50 36.09 22.24 1.46 44.80 -20.48 -10.04 11.45	45.09 27.36 23.35 -4.10 48.93 -19.68 -16.34 16.12	<b>71.13</b> 68.65 63.28 36.54 52.59 -29.48 -9.51 32.87	× × × × × × × × × ×	~~~~~~~~

Table 2: Evaluation on five Southeast Asian languages. Strict-Overlap metrics include BLEU-4, sacreBLEU, and ROUGE-L. Semantic-Friendly metrics include chrF, METEOR, and BERTScore. For each metric column: **Bold** values indicate the highest score, and <u>Underlined</u> values indicate the second highest score across all models. FT: Fine-tuned; OS: Open Source.

			BLEU-4					chrF	chrF			
Model	fil	id	lo	my	vi	fil fil	id	lo	my	vi	(BLEU-chrF)	
Qwen2.5-0.5b	0.03	0.03	0.02	0.01	0.01	0.16	0.12	0.40	0.25	0.06	0.11	
+ SFT-LTP	1.86	$4.02_{\uparrow 4.00}$	$0.25_{\uparrow 0.22}$	$0.15_{\uparrow 0.14}$	$3.12_{\uparrow 3.11}$	4.85	$10.38_{\uparrow 10.26}$	$1.41_{\uparrow 1.00}$	$1.07_{\uparrow 0.82}$	$8.55_{18.50}$	3.57	
+ GRPO-SAR	2.31 12.28	$4.32_{\uparrow 4.29}$	$0.26_{\uparrow 0.24}$	$0.16_{\uparrow 0.16}$	6.29 <sub>16.27</sub>	$5.00_{\uparrow 4.84}$	$10.64_{\uparrow 10.52}$	$1.38_{\uparrow 0.98}$	$1.22_{\uparrow 0.96}$	12.36 12.30	4.39	
+ LLME	0.3210.29	$0.45_{\uparrow 0.42}$	$0.08_{\uparrow 0.06}$	$0.07_{\uparrow 0.06}$	$0.79_{\uparrow 0.78}$	$1.20_{\uparrow 1.04}$	$1.66_{\uparrow 1.54}$	$0.45_{\uparrow 0.05}$	$1.25_{\uparrow 1.00}$	$2.82_{\uparrow 2.76}$	0.91	
Avg.	1.13	2.21	0.15	0.10	2.55	2.80	5.70	0.91	0.95	5.95	2.25	
Qwen2.5-3b	5.80	14.25	1.11	1.83	17.06	8.83	17.71	2.05	3.03	19.35	9.10	
+ SFT-LTP	$23.01_{\uparrow 17.21}$	$26.00_{\uparrow 11.75}$	$1.17_{\uparrow 0.05}$	$2.53_{\uparrow 0.69}$	$27.94_{\uparrow 10.87}$	$20.14_{\uparrow 11.31}$	$24.25_{\uparrow 6.54}$	$2.22_{\uparrow 0.17}$	$3.53_{\pm 0.50}$	$24.55_{\uparrow 5.20}$	15.53 <sub>16.43</sub>	
+ GRPO-SAR	25.58	$29.11_{\uparrow 14.86}$	4.39 <sub>13.28</sub>	$2.77_{\uparrow 0.93}$	$32.54_{\uparrow 15.48}$	23.62	$27.83_{\uparrow 10.12}$	5.17	$3.45_{\uparrow 0.42}$	$27.88_{\pm 8.53}$	18.23	
+ LLME	29.71 23.91	$34.73_{\uparrow 20.48}$	$5.15_{\uparrow 4.04}$	4.56	$31.20_{\uparrow 14.14}$	25.52	$30.22_{\uparrow 12.51}$	5.81 <sub>13.76</sub>	$5.53_{12.50}$	$26.98_{\uparrow 7.63}$	19.94	
Avg.	21.03	26.02	2.96	2.92	27.19	19.53	25.00	3.81	3.89	24.69	15.70	

Table 3: Ablation Study of Qwen-2.5 0.5B and Qwen-2.5 3B on five Southeast Asian languages. All values are rounded to two decimal places. Improvements over the Zero-shot baseline (<u>underlined rows</u>).

(a 3890% relative improvement) over its zero-shot baseline score of 0.11. This underscores the efficacy of reward modeling, consistent with findings in instruction tuning (Ouyang et al., 2022; Wu et al., 2023).

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

Our proposed LLMs-Enhanced (LLME) stage demonstrates further advancements, particularly for the larger Qwen-2.5 3B model. With LLME, the Qwen-2.5 3B model achieves the highest overall BLEU-chrF score of 19.94, representing a 10.84 absolute point improvement (a 119% relative increase) over its zero-shot baseline. This highlights the synergistic benefits of our full pipeline. While the LLME stage yields more modest gains for the Qwen-2.5 0.5B model in the current setup (overall BLEU-chrF of 0.91), the substantial cumulative improvements from SFT-LTP and GRPO-SAR on this smaller model, and the peak performance achieved by the 3B model with LLME, collectively validate the effectiveness and scalability of our modular tuning strategy in significantly enhancing translation quality.

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

### 4.4 Effect of Data Distillation on Performance

We assess the impact of our quality filtering approach by comparing full-scale Supervised Fine-Tuning with Language-Token Prefixing (SFT-LTP) against subsequent reward-informed filtering and tuning via GRPO-SAR, using our MERIT-3B model. Table 4 details the number of retained training instances per language and the corresponding

513

506

521 522 523

- 524 525
- 526 527
- 528 529

530 531 532

533

534 535

53

538 539

540

541 542

> 543 544 545

546

548

551 552

553

5

554 555

556

Second, the current reward model underlying

overall BLEU-chrF scores for these configurations.

training instances. In contrast, the GRPO-SAR

stage strategically curates this data, drastically re-

ducing the volume to only 9,126 instances. This

constitutes an average data reduction of 77.2%,

with the most significant reduction observed for

Vietnamese, where the training data was cut by

87.8% (from 8,000 to 976 instances). Remarkably,

despite this substantial data pruning, the overall

BLEU-chrF score not only signifies the efficient

retention of highly informative samples but actu-

ally improves from 15.53 (achieved with SFT-LTP

on 40,000 instances) to 18.23 with GRPO-SAR

on the reduced dataset. This represents a relative

These findings underscore the efficacy of our

reward-based filtering (GRPO-SAR) as a data-

efficient strategy that can simultaneously reduce

training data requirements and enhance model per-

formance. This offers a compelling alternative to

training on larger, potentially noisier, unfiltered

datasets. The benefits of leveraging reward signals

for targeted data curation align with effective strate-

gies observed in other generative AI tasks, such as

summarization and dialogue tuning (Lu et al., 2022;

Our work, culminating in the MERIT framework,

demonstrates the significant potential of combining

data filtering techniques, such as the Score Accu-

racy Reward (SAR) driven GRPO, with efficient

fine-tuning strategies like Language-Token Prefix-

ing (LTP) for multilingual translation, especially

into low-resource languages (LRLs). The proposed

BLEU-chrF composite metric has also provided a

balanced view of lexical and semantic performance.

While MERIT-3B exhibits strong performance rela-

tive to its scale and against comparable open-source

models, several limitations persist and pave the way

Lao and Burmese, can introduce encoding incon-

sistencies. These not only affect the performance of QE agents used in SAR but also potentially skew

standard evaluation metrics. Future iterations could

incorporate more robust character normalization or

transliteration techniques at the data preprocessing

stage, or develop QE models less sensitive to such

First, script-related challenges, particularly for

Ouyang et al., 2022).

for future exploration.

variations.

Further Discussion

4.5

performance increase of approximately 17.4%.

The SFT-LTP stage utilizes the full set of 40,000

GRPO-SAR, while effective, may inadvertently prioritize adequacy (accuracy of content, as captured by our specific SAR function focusing on numerical or key information matching) sometimes at the expense of optimal fluency. This can occasionally lead to subtle grammatical artifacts in some translations. Future work could investigate multiobjective reward functions that explicitly balance adequacy, fluency, and even other aspects like style or register, potentially drawing on more diverse human feedback signals beyond simple ratings. 557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

Moreover, as noted by Zhang et al. (2022), many LLMs, including some baselines we compared against, are often evaluated in zero-shot or few-shot settings for translation. This might not fully reveal their capabilities, which could be significantly enhanced with more sophisticated prompting strategies or in-context learning techniques specifically tailored for translation. Exploring how our data filtering and fine-tuning methods can synergize with advanced prompting for even larger LLMs is a promising direction.

# 5 Conclusion

This work introduces MERIT, a unified framework combining a reconstructed benchmark and modular training strategies for Chinese–Low-Resource Language (LRL) neural machine translation. We reconstruct a clean and balanced evaluation suite from the ALT corpus, enabling reliable assessment across five Southeast Asian languages. On the training side, MERIT incorporates languageconditioned fine-tuning and reward-guided data selection to improve translation quality efficiently.

Experiments demonstrate that our MERIT-3B model outperforms comparable open-source models and approaches the performance of significantly larger proprietary systems, particularly in LRL scenarios. Ablation results confirm that reward-informed filtering with GRPO-SAR is especially effective, achieving better performance with less data.

Overall, our findings reinforce the importance of strategic data selection and modular fine-tuning over sheer model scale in low-resource settings. Future work will extend MERIT to more languages and explore scaling with stronger reward functions and larger base models.

### References

604

612

613

614

615

616

618

619

623

624

625

636

641

642

646

647

- 2023. Emnlp 2023 ethics faq. https://2023.emnlp. org/ethics/faq. EMNLP 2023 Conference Website.
- Alexander Adelaar. 2012. Indonesian and malay: Are they different languages? *Annual Review of Linguistics*, pages 1–23.
- Naveen Arivazhagan and 1 others. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. In *Proceedings of ACL*.
- Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zeroshot cross-lingual transfer and beyond. *Transactions* of the Association for Computational Linguistics, 7:597–610.
  - Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings* of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT.
- Emily M. Bender and Batya Friedman. 2018. Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6:587–604.
  - Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, and 1 others. 2022. No language left behind: Scaling human-centered machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 11167– 11191.
  - Yiming Cui and et al. 2025. Multilingual machine translation with open large language models at practical scale. *arXiv preprint arXiv:2502.02481*.
- Google DeepMind. 2025. Gemini 2.5 technical report: Unlocking multimodal understanding across millions of tokens. Google AI Blog. URL: https://ai. googleblog.com/gemini-25-report.
- Long Duong, Thanh-Le Ha, and Le-Minh Phuong. 2017. English–malay parallel corpus construction from wikipedia. In *Proceedings of the 21st Workshop on Asian Language Resources*, pages 29–38.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2023a. *Ethnologue: Languages of the World*, 26 edition. SIL International, Dallas, Texas. L1 speaker estimates for 12 ALT languages; language-specific URLs available in Appendix/Table 1.
- David M. Eberhard, Gary F. Simons, and Charles D. (eds.) Fennig. 2023b. Malay Macrolanguage. *Ethnologue: Languages of the World*, 26th ed. https: //www.ethnologue.com/macrolanguages.

Can Enis and Mark Hopkins. 2024. From llm to nmt: Advancing low-resource machine translation with claude. *arXiv preprint arXiv:2404.13813*. 656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

- Miquel Esplà-Gomis, Víctor M. Sánchez-Cartagena, Jeff Zaragoza-Bernabeu, and Felipe Sánchez-Martínez. 2020. Bicleaner at WMT 2020: Universitat d'Alacant–Prompsit's submission to the parallel corpus filtering shared task. In *Proceedings of the 5th Conference on Machine Translation (WMT 2020)*, Online. Association for Computational Linguistics.
- Angela Fan, Shruti Bhosale, and Yihyung C. Aharoni. 2021. Beyond english-centric multilingual machine translation. In *Proc. of ACL*.
- Markus Freitag, Yaser Al-Onaizan, and 1 others. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. In *Proceedings of ACL*.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*.
- Sarthak S Ghosal, Ritam Roy, Sumanth Doddapaneni, and 1 others. 2024. Promptrefine: Enhancing fewshot performance on low-resource indic languages with example selection from related example banks. *arXiv preprint arXiv:2412.05710*.
- Naman Goyal and et al. 2022. The flores-200 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, pages 855–872.
- Yutong Huang, Minghao Wang, Yujie Xie, and 1 others. 2025. Benchmax: A comprehensive multilingual evaluation suite for large language models. *arXiv* preprint arXiv:2502.07346.
- Zhengyong Jiang, Lingfan Zhang, Ruobing Guo, and 1 others. 2025. Deepseek v3 vs o3-mini: How well can reasoning llms evaluate mt and summarization? *arXiv preprint arXiv:2504.08120*.
- Melvin Johnson, Mike Schuster, Quoc V Le, and 1 others. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. In *Transactions of the Association for Computational Linguistics*.
- Philipp Koehn and 1 others. 2020. Findings of the 2020 wmt shared task on parallel corpus filtering and alignment. In *Proceedings of the Fifth Conference on Machine Translation*, pages 726–746.
- Mingxu Li, Yuxian Zhang, Shuai He, Zhipeng Li, Haoran Zhao, Junkai Wang, Ning Cheng, and Tianxiang Zhou. 2023. From quantity to quality: Boosting LLM performance with self-guided data selection for instruction tuning. *arXiv preprint arXiv:2308.12032*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*.

816

817

765

Ming Liu, Jingjing Xu, Jiqiang Zhang, and Yidong Chen. 2021. Research on uyghur–chinese neural machine translation based on ccmt shared tasks. In *Proceedings of the China Conference on Machine Translation* (*CCMT*), pages 120–130.

711

713

714

715

716 717

718

719

721

722

724

727

729

736

737

738

739

740

741

742

743

744

745

746

747

749

751

753

756

757

759

760

761

762

- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Lalida Lowphansirikul and Ekaterina Chuangsuwanich. 2020. Scb-mt en-th: A large, clean parallel corpus for english-thai machine translation. In *Proceedings* of the 28th International Conference on Computational Linguistics, pages 3612–3622.
  - Zichao Lu and 1 others. 2022. Quark: Controllable and compositional generation of text using reinforcement learning. In *Proceedings of ACL*.
- Shuming Ma, Li Dong, Shaohan Huang, Saksham Singhal, Shuming Li, Long Zhou, Yaru Wang, Kaitao Song, Yuxiang Zhang, and Furu Wei. 2021. Deltalm: Encoder-decoder pre-training for language generation and translation by augmenting pretrained multilingual encoders. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 9064–9078.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2018. Model cards for model reporting. *arXiv preprint arXiv:1810.03993*.
- Dragos Ştefan Munteanu and Daniel Marcu. 2005. Improving machine translation performance by exploiting non-parallel corpora. *Computational Linguistics*, 31(4):477–504.
- Long Ouyang, Jeffrey Wu, Xu Jiang, and 1 others. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318.
- Maja Popovic. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation (WMT)*, pages 392–395.
- Matt Post. 2018. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation*.
- Ricardo Rei, Andre F. T. Martins, Ana B. Farinha, and Alon Lavie. 2020. Comet: A neural framework for

mt evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2685–2702.

- ACL Rolling Review. 2023. Responsible nlp research checklist. https://aclrollingreview. org/responsibleNLPresearch/.
- Víctor M. Sánchez-Cartagena, Marta Bañón, Sergio Ortiz-Rojas, and Gema Ramírez-Sánchez. 2018. Prompsit's submission to WMT 2018 parallel corpus filtering shared task. In *Proceedings of the 3rd Conference on Machine Translation (WMT 2018)*, Brussels, Belgium. Association for Computational Linguistics.
- Gary F. Simons and Charles D. (eds.) Fennig. 2023. Central Khmer. *Ethnologue: Languages of the World*, 26th ed. https://www.ethnologue.com/ language/khm.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proc. of NIPS*.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proc. of CVPR*.
- NLLB Team and Meta AI. 2022. No language left behind: Building a complete world-wide machine translation system. Technical Report arXiv:2207.04672, Meta AI Technical Report.
- Ye Kyaw Thu, Pa Pa Win, Masao Utiyama, Andrew Finch, and Eiichiro Sumita. 2016. Introducing the asian language treebank (alt). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC)*, pages 1574–1578.
- Ronald J. Williams and David Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280.
- Yuntao Wu and 1 others. 2023. Fine-grained alignment of language models with preferences. In *arXiv preprint arXiv:2305.18290*.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings* of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 483–498.
- Susan Zhang, Stephen Roller, Naman Goyal, and 1 others. 2022. Opt: Open pre-trained transformer language models. In *arXiv preprint arXiv:2205.01068*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

821

822

823

# A Appendix

# A.1 Ethics Statements

This work presents a Chinese-centric multilingual translation benchmark targeting five Southeast Asian low-resource languages (LRLs), constructed from publicly available corpora and evaluated under reproducible protocols. We aim to support responsible research in multilingual NLP by releasing rigorous evaluation resources while proactively addressing ethical concerns related to data provenance, model fairness, environmental impact, and potential misuse.

**Data Privacy and Consent** All data is derived 830 from the publicly available ASEAN Languages Treebank (ALT), which includes multilingual translations of government and news texts. While the dataset is openly licensed, the original collection did not explicitly document consent procedures or 835 personal identifiable information (PII) removal. To 836 mitigate this, we apply a multi-stage filtering process to exclude named entities, explicit language, and potentially sensitive content. Nonetheless, due to the limitations of automated and manual filtering, some residual risk may remain. We follow the 841 data statements framework (Bender and Friedman, 2018) and document licensing, provenance, and usage constraints in the appendix. 844

Bias and Fairness Despite the use of a three-845 stage filtering pipeline and expert-rated supervision, the training data may still encode latent cultural, linguistic, or regional bias-particularly due to its English-pivoted design and limited coverage 849 of dialectal variations or non-standard orthographies. Annotators are bilingual graduate students, and while they are experienced, demographic diversity is limited. Future work will prioritize the 853 inclusion of more diverse annotators and typolog-854 ically broader sources to mitigate such represen-855 tational imbalances. Our work aligns with global AI ethics principles of fairness, transparency, and 857 non-maleficence (Gebru et al., 2018).

Environmental Impact Model training and inference were conducted on a single-node NVIDIA A100 80GB GPU. We log training FLOPs and wallclock runtime for both the SFT and GRPO stages. While the GRPO procedure improves data efficiency through reward-based filtering, it introduces additional computational cost. We estimate that the total training corresponds to a typical singlenode compute workload and plan to explore more lightweight reward models or compute-efficient alternatives to reduce carbon impact in future iterations (emn, 2023).

867

868

869

870

871

872

873

874

875

876

877

878

879

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

Intended Use and Misuse Risks The benchmark is designed to support objective evaluation and supervised training for Chinese–LRL translation tasks. It is intended for academic research and language technology development, particularly in regions underrepresented in NLP. However, misuse is possible—such as generating misinformation or content targeting marginalized communities. We explicitly discourage such applications and recommend that any downstream use include fairness auditing, risk controls, and human oversight (Mitchell et al., 2018).

**Transparency and Reproducibility** We adhere to the ACL Responsible NLP Research Checklist (Review, 2023) and release all code, data, and model checkpoints under a CC BY-NC 4.0 license. All filtering procedures, model configurations, and hyperparameter settings are fully documented. Some language-specific heuristics (e.g., token ratio thresholds) were empirically selected and may not generalize across domains; future validation is necessary to ensure robustness.

# A.2 Limitations

Despite the encouraging results achieved by our proposed framework and the MERIT-3B model, this work has several limitations that warrant discussion and offer avenues for future improvement.

First, while our evaluation benchmark improves upon existing English-pivoted resources by constructing direct Chinese–LRL sentence pairs, its current scope is confined to five Southeast Asian languages. Other significant low-resource languages, including domestic Chinese minority languages such as Tibetan, Uyghur, and Kazakh, remain unaddressed due to the scarcity of highquality aligned corpora. Expanding the linguistic diversity of our benchmark is crucial for assessing broader generalizability.

Second, the ALT-based test suite, although semantically aligned through shared alt\_id indexing, is fundamentally constrained by its original English-centric design. While our realignment efforts aim to mitigate semantic drift when adapting it for Chinese–LRL evaluation, some residual domain-specific or stylistic artifacts originating

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

967

968

from the English-centric source may persist and subtly influence translation assessment.

916

917

918

919

921

922

923

925

927

928

931

932

934

935

939

942

947

948

951

952

954

957

963

964

966

Third, although our methodology employs a QE agent and the statistical-semantic Score Accuracy Reward (SAR) function for automatic data filtering, the scale of human validation for these components is currently limited. The expert-rated set used for developing or validating the reward model is modest in size. This might restrict the overall robustness and generalizability of the SAR model, particularly its alignment with nuanced human preferences across diverse linguistic phenomena. Future work should prioritize the integration of more extensive and varied human annotations.

Fourth, while we have evaluated a range of LLMs under zero-shot, SFT, and GRPO regimes, the decoding strategies (e.g., beam size, sampling temperature) and specific prompt formats were kept fixed across these models for controlled comparison. These settings can significantly influence translation behavior and perceived quality, especially for proprietary models whose internal mechanisms are opaque. A more exhaustive exploration of model-specific optimal decoding parameters and prompt engineering could reveal further performance variations.

Finally, due to computational resource constraints, our current experiments, including the development and evaluation of the MERIT-3B model and its associated fine-tuning framework (SFT-LTP, GRPO-SAR), have been conducted on models up to the 3B parameter scale. We have not yet extended this framework to significantly larger parameter models (e.g., 7B+, or state-of-the-art models in the tens or hundreds of billions of parameters). Applying and evaluating our data filtering and rewardinformed tuning strategies on such larger-scale models is an important next step to ascertain their scalability and potential for even greater performance gains, though this would require substantial additional computational resources. Furthermore, a detailed efficiency analysis, including training time, inference latency, and computational cost of the filtering and fine-tuning stages, was not conducted and would be valuable for assessing practical deployability.

### A.3 Experimental Setup

All experiments were conducted on a local workstation equipped with two NVIDIA RTX 3090 GPUs (24 GB). Under a 2×2 parallel configuration, the per-GPU batch size was set to 8 with a gradient accumulation step of 2, resulting in an effective total batch size of 32. The maximum input sequence length was set to 1024 tokens, and the initial learning rate was configured as 2e-4. The system environment included Ubuntu 20.04, CUDA 12.1, and Python 3.10, with PyTorch 2.1 and Transformers v4.49 as the core libraries.

All training was performed using standard mixed-precision (fp16) computation via custom training scripts. Due to hardware limitations, the batch size was carefully adjusted to fit within the available GPU memory, and no experiments were conducted using larger-parameter models. To ensure reproducibility, all random seeds were fixed, and detailed runtime logs were maintained for each experiment.

### A.4 Reward Function

In this study, we introduce the Score Accuracy Reward (SAR) function as a key component of the reward mechanism for evaluating the numerical accuracy of generated outputs. Specifically, we design a step-wise reward strategy, which operates as follows:

- A reward of 2.0 is assigned if the model's output exactly matches the reference answer.
- A reward of 1.0 is assigned if the output deviates from the reference by a small margin.
- A reward of 0.0 is assigned if the deviation exceeds the acceptable threshold.

This step-wise reward formulation differs from conventional binary reward functions commonly used in correctness or format-checking tasks. It is motivated by two main considerations:

(1) Task-specific suitability. SAR is designed for numerical question answering and tasks that require precise arithmetic reasoning. In such settings, the reference answer is often a numeric value with an acceptable tolerance range rather than a single exact string. Therefore, outputs that are numerically close to the reference can be considered approximately correct. Compared to rigid binary matching, the step-wise reward better aligns with the intrinsic characteristics of these tasks.

(2) More informative training signal. Unlike traditional 0/1 rewards, step-wise rewards provide finer-grained feedback, allowing the model to receive gradient signals that vary with the degree of deviation. This facilitates smoother optimization

and more stable convergence during training, en-1015 abling the model to gradually improve its numeric 1016 prediction capabilities. In contrast, overly rigid 1017 reward mechanisms may lead to sparse or unin-1018 formative training signals and hinder early-stage learning. 1020

#### A.5 **Recruitment And Payment**

1021

1023

1024

1025

1028

1029 1030

1031

1032

1033

1035

1036

1037

1038

1040

1041

1042

1043

1044

To ensure the accuracy and objectivity of human evaluation, we recruited ten annotators with academic backgrounds in the target Southeast Asian languages. All annotators were either language instructors or graduate students from relevant universities. For each target language, two annotators were assigned, and a cross-review protocol was adopted to enhance annotation quality and consistency.

All participants had formal training in translation or linguistics and possessed strong language comprehension and evaluative capabilities. Annotators were compensated at a rate of 1 RMB per evaluated sample. Before the evaluation began, all participants received detailed instructions and train-13 ing on annotation guidelines. Participation was voluntary, and compensation was provided proportionally based on the amount of completed work.

Since the dataset contains no personally identifiable information (PII) and the task involves only linguistic quality assessment, the annotation process entails no ethical risks and does not require institutional ethics approval.

Algorithm 1 Elite Parallel Data Sampler

**Input:** XML dataset  $\mathcal{D}_{xml}$ ;

Target sizes  $\{T_{\text{train}}, T_{\text{dev}}, T_{\text{test}}\};$ Domain set  $\mathcal{D}$ 

7

10

11

12

1

- **Output:** Datasets  $\{D_{\text{train}}, D_{\text{dev}}, D_{\text{test}}\}$
- **Stage 1: Domain Balancing** Global pool  $M^* \leftarrow$ 1 Ø
- <sup>2</sup> foreach domain  $d \in \mathcal{D}$  do
- $M_d \leftarrow \{ f \in \mathcal{D}_{\text{xml}} \mid \text{domain}(f) = d \} M^* \leftarrow$ 3  $M^* \cup M_d$
- 4  $N \leftarrow \sum_{d} |M_d|; P_d \leftarrow |M_d|/N \; (\forall d)$
- 5 Stage 2: Proportional Split Generation foreach split  $t \in \{train, dev, test\}$  do

$$\left| \begin{array}{ccc} S_t \leftarrow \emptyset \text{ foreach } domain \ d \in \mathcal{D} \text{ do} \\ \left\lfloor \begin{array}{c} n_t(d) \ \leftarrow \ \lfloor T_t \ \cdot \ P_d \rfloor & S_t \ \leftarrow \ S_t \ \cup \\ & \mathsf{SHUFFLE}(M_d)[:n_t(d)] \end{array} \right. \\ \left[ \begin{array}{c} \Delta \leftarrow T_t - |S_t| \ \text{ if } \Delta > 0 \ \text{ then} \\ \left\lfloor \begin{array}{c} S_t \leftarrow S_t \cup \mathsf{SAMPLE}(M^* \setminus S_t, \Delta) \end{array} \right. \\ & \text{else} \\ \left\lfloor \begin{array}{c} \text{ if } \Delta < 0 \ \text{ then} \\ \left\lfloor \begin{array}{c} S_t \leftarrow \mathsf{HEAD}(S_t, T_t) \end{array} \right. \\ & M^* \leftarrow M^* \setminus S_t \end{array} \right. \end{array} \right.$$

14 Stage 3: Domain Proportion Verification foreach domain  $d \in \mathcal{D}$  do

17 Stage 4: Output Generation foreach split  $t \in$ {train, dev, test} do

 $SHUFFLE(S_t)$ 18

9 return 
$$\{S_{\text{train}}, S_{\text{dev}}, S_{\text{test}}\}$$

		Overall						
Method	fil	id	lo my		vi	(Size / BLEU-chrF)		
MERIT-3B								
+ SFT-LTP	8,000	8,000	8,000	8,000	8,000	40,000 / 15.53		
+ GRPO-SAR	$1,851_{\downarrow 76.9\%}$	$1,779_{\downarrow 77.7\%}$	$2,058_{\downarrow 74.2\%}$	$2,462_{\downarrow 69.2\%}$	$976_{\downarrow 87.8\%}$	9,126 <sub>↓77.2%</sub> / 18.23 <sub>↑17.4%</sub>		
+ LLME	$2,891_{\downarrow 63.9\%}$	$3,104_{\downarrow 61.2\%}$	3,300,58.8%	3,764,53.0%	$2,193_{\downarrow 72.6\%}$	$15,252_{\downarrow 61.9\%}$ / $19.94_{\uparrow 28.4\%}$		

Table 4: Training size comparison across five low-resource languages for MERIT-3B. Overall column shows total training data (with percentage reduction relative to initial 40,000) and BLEU-chrF score (with percentage improvement relative to the SFT-LTP stage).



Figure 3: Training loss and reward evolution across SFT and GRPO strategies.

Algorithm 2 Data Integrity Validation **Input:** Candidate splits  $\{S_t\}$ ; Anomaly threshold  $\tau_{\text{anom}} = 0.7$ ; Validity threshold  $\tau_{\text{valid}} = 0.8$ ; PCA dimension k; Regularization  $\lambda$ . **Output:** Validation ∈ {True, False} 20 Stage 1: Feature Extraction  $F_{\text{outlier}} \leftarrow \emptyset \ X \leftarrow$ 21 foreach sample  $s \in S_{train}$  do CNNENCODER(s)22 е  $\leftarrow$ BERTEMBED(s)с  $\leftarrow$  $\mathbf{v}$ L2NORMALIZE( $[\mathbf{e}|\mathbf{c}]$ )  $X \leftarrow X \cup \{\mathbf{v}\}$ 23 STANDARDIZE(X) PCA(X, k)24 Stage 2: Cluster Anomaly Discovery HDBSCAN(X) foreach *cluster* C $\leftarrow$  $C_k \in C$  do  $|C_k|/|X|$ Σ 25  $\rho_k$ COSINESIMILARITY( $C_k$ ) foreach sample  $x_i \in C_k$  do  $s_i \leftarrow 1 - \frac{1}{|C_k|} \sum_{j \in C_k} \Sigma_{ij} \text{ if } s_i > \tau_{anom}$ then  $\lfloor F_{\text{outlier}} \leftarrow F_{\text{outlier}} \cup \{x_i\}$ 26 27  $w_k \leftarrow \rho_k \cdot \left(1 - |F_{\text{outlier}} \cap C_k| / |C_k|\right)$ 28 29 Stage 3: Composite Validity Score  $V_{\text{geo}}$  $\overline{\prod_k w_k^{\rho_k} V_{\text{pen}} \leftarrow e^{-\lambda |F_{\text{outlier}}|} V \leftarrow V_{\text{geo}} \cdot V_{\text{pen}}}$ Stage 4: Validation Check if  $V < \tau_{valid}$  then 30 is Valid ← False; PURGECORRUPTED-31  $DATA(S_t)$ 32 else

is Valid  $\leftarrow$  True

33

### A.6 Feature Extraction

Let H', W' denote the height and width of the final 1046 CNN feature map after L layers. Let  $d_e$  be the 1047 number of CNN output channels, and let  $d_c$  be the 1048 BERT hidden dimension (which equals the token 1049 embedding size  $d_m$ ). 1050

1045

1051

1052

1053

1054

1060

Let the training set be  $S_{train} = \{s_i\}_{i=1}^{N}$ . Each sample  $s_i$  comprises an image  $s_i^{img}$  and text  $s_i^{text}$ , processed as follows:

### 1. CNN Encoding:

$$F_i^{(0)} = s_i^{\text{img}}$$
 (8) 10

$$F_i^{(l)} = \text{ReLU}\big(W^{(l)} * F_i^{(l-1)} + b^{(l)}\big), \quad (9)$$
 105

$$F_i^{(L)} \in \mathbb{R}^{d_e \times H' \times W'} \tag{10}$$

$$e_i^j = \frac{1}{H'W'} \sum_{h=1}^{H'} \sum_{w=1}^{W'} F_i^{(L)}(j,h,w) \quad (11)$$
 1050

$$\mathbf{e}_i = [e_i^1, \dots, e_i^{d_e}]^\top \in \mathbb{R}^{d_e} \tag{12}$$

## 2. BERT Embedding:

$$X_i^{(0)} = \left[ E_{\text{tok}}(w_t) + E_{\text{pos}}(t) \right]_{t=1}^T$$
 (13) 1061

$$X_i^{(\ell)} = \text{TransformerLayer}^{(\ell)} \left( X_i^{(\ell-1)} \right), \tag{14}$$

$$\mathbf{c}_i = X_i^{(L')}[1] \in \mathbb{R}^{d_c} \tag{15}$$

1068

1069

1070

1071

1072

1073

1074

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

# 3. Concatenation & Normalization:

065 
$$\mathbf{u}_i = \begin{bmatrix} \mathbf{e}_i \\ \mathbf{c}_i \end{bmatrix} \in \mathbb{R}^{d_e + d_c}$$
 (16)

 $\mathbf{v}_i = \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|_2} \tag{17}$ 

## 4. Matrix Assembly:

 $X = \begin{pmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_N^\top \end{pmatrix} \in \mathbb{R}^{N \times (d_e + d_c)} \qquad (18)$ 

5. PCA Reduction:

$$Z = \mathrm{PCA}_{k}(X) = \begin{pmatrix} \mathbf{z}_{1}^{\top} \\ \vdots \\ \mathbf{z}_{N}^{\top} \end{pmatrix} \in \mathbb{R}^{N \times k} \quad (19)$$

A.7 Cluster Anomaly Detection

# 1. DBSCAN Clustering:

$$\mathcal{N}_{\epsilon}(\mathbf{z}_{i}) = \{\mathbf{z}_{j} : \|\mathbf{z}_{j} - \mathbf{z}_{i}\|_{2} \le \epsilon\}$$
(20)  
$$\mathbf{z}_{i} \iff |\mathcal{N}_{\epsilon}(\mathbf{z}_{i})| \ge \text{MinPts}$$
(21)

$$\mathbf{z}_i \rightsquigarrow \mathbf{z}_j \quad \iff \exists (\mathbf{z}_{i_0}, \dots, \mathbf{z}_{i_m})$$
 (23)

$$\mathbf{z}_p, \mathbf{z}_q \Leftrightarrow \iff \exists \mathbf{z}_o : \mathbf{z}_p, \mathbf{z}_q \rightsquigarrow \mathbf{z}_o \quad (24)$$

2. Clustering Procedure:

- (a) Mark all  $\mathbf{z}_i$  unvisited, set cluster counter  $c \leftarrow 0$ .
- (b) For each unvisited  $\mathbf{z}_i$ :
  - i. Mark  $\mathbf{z}_i$  visited; let  $N \leftarrow \mathcal{N}_{\epsilon}(\mathbf{z}_i)$ .
  - ii. If |N| < MinPts, label  $\mathbf{z}_i$  as noise.
  - iii. Else:

$$c \leftarrow c+1, \quad C_c \leftarrow \{\mathbf{z}_i\} \quad (25)$$

(c) expand(C, N):

$$\mathbf{z}_{j}: \begin{cases} \text{if unvisited: mark visited} \\ N' \leftarrow \mathcal{N}_{\epsilon}(\mathbf{z}_{j}) \\ \text{if } |N'| \geq \text{MinPts}, \quad N \leftarrow N \cup N' \\ (26) \end{cases}$$

1089 if  $\mathbf{z}_j \notin C$ , then  $C \leftarrow C \cup \{\mathbf{z}_j\}$  (27) 1090 (d) Resulting clusters:  $C_1, \dots, C_K$ 

## 3. Cosine Similarity:

$$\Sigma_{ij} = \frac{\mathbf{z}_i^\top \mathbf{z}_j}{\|\mathbf{z}_i\|_2 \|\mathbf{z}_j\|_2}$$
(28) 1092

1091

1097

# 4. Anomaly Score: 1093

$$s_i = 1 - \frac{1}{|C_k|} \sum_{j \in C_k} \Sigma_{ij}$$
 (29) 1094

5. Outlier Set: 1095

$$\mathcal{F}_{\text{outlier}} = \{ i \mid s_i > \tau_{\text{anom}} \}$$
(30) 1096

## 6. Cluster Weighting:

$$\rho_k = \frac{|C_k|}{N} \tag{31}$$

$$\delta_k = 1 - \frac{|\mathcal{F}_{\text{outlier}} \cap C_k|}{|C_k|} \tag{32}$$

$$w_k = \rho_k \,\delta_k \tag{33}$$