A Single Model Ensemble Framework for Neural Machine Translation using Pivot Translation

Anonymous ACL submission

Abstract

Despite the significant advances in neural ma-001 chine translation, performance remains subpar for low-resource language pairs. Ensembling multiple systems is a widely adopted technique 005 to enhance performance, often accomplished by combining probability distributions. However, the previous approaches face the challenge of 007 high computational costs for training multiple models. Furthermore, for black-box models, averaging token-level probabilities at each decoding step is not feasible. To address the problems of multi-model ensemble methods, we present a pivot-based single model ensemble. The proposed strategy consists of two steps: pivot-based candidate generation and post-hoc aggregation. In the first step, we generate candidates through pivot translation. This can be 017 018 achieved with only a single model and facilitates knowledge transfer from high-resource pivot languages, resulting in candidates that are not only diverse but also more accurate. Next, in the aggregation step, we select k high-quality candidates from the generated candidates and merge them to generate a final translation that outperforms the existing candidates. Our experimental results show that our method produces translations of superior quality by leveraging candidates from pivot translation to capture the subtle nuances of the source sentence.

1 Introduction

033

037

041

Neural machine translation (NMT) models exhibit outstanding capabilities when a large volume of the parallel corpus is available (e.g., translate from and to English). However, their performance still falls short in cases involving low-resource languages (e.g., Basque) and translating between non-English languages from different language families (e.g., German-Russian) (Artetxe et al., 2018). Top-performing large language models (LLMs), such as GPT models (Ouyang et al., 2022), also demonstrate suboptimal translation performance in low-resource language pairs (Robinson et al., 2023; Moslem et al., 2023). The scarcity of parallel data, primarily due to limited cultural interaction, makes the low-resource translation task more challenging. 042

043

044

047

048

053

054

056

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

In many generation tasks, ensembling multiple systems has proven to be a successful strategy for performance enhancement. In NMT, traditional ensemble methods average probability distributions over output tokens from multiple models during decoding. However, the expensive cost of training multiple models is the primary shortcoming of ensemble decoding. Additionally, computing tokenlevel probabilities at each decoding step is not feasible with recent black-box models such as GPT-40 and Gemini (OpenAI, 2024; Team, 2023).

Ensemble methods that can be utilized even when token-level probabilities cannot be computed have also been proposed. Selection-based ensemble method involves generating candidates from multiple models and then selecting the best candidate among them (Wang et al., 2022; Hendy et al., 2023). However, in this ensemble fashion, the final output space is limited to the existing candidate pool. In contrast, the generation-based ensemble, such as LLM-Blender (Jiang et al., 2023), creates improved outputs using candidates obtained from multiple models. This approach aims to generate a final output superior to the existing candidates. Nonetheless, the main drawback from the notably high cost of generating candidates through multiple models remains, inducing computational overhead. As the size of the models used in the ensemble increases, the cost proportionally escalates, becoming more burdensome. In addition, due to the varying performance of MT systems, the quality of some candidates can be significantly lower than others, leading to a degradation in the overall performance.

To alleviate the problems above, we propose **Pivot**-based single model Ensemble (PIVOTE), a novel generation-based approach. Our intuition of a single model ensemble primarily stems from pivot translation, which can produce diverse and more accurate translations. Pivot translation (Wu and Wang, 2007; Utiyama and Isahara, 2007) is a method that splits the end task into two sequential steps: source→pivot and pivot→target. Pivoting has been employed to enhance low-resource translation by transferring knowledge from high-resource pairs. In many cases, English, being a resource-rich language, serves as the intermediate language. However, we employ not only English but various pivot languages for candidate generation, thereby producing diverse hypotheses using a single model.

084

100

101

102

103

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

In the next aggregation step, we select the top candidates for the ensemble and merge them to generate the final output. Since the quality of candidates directly impacts the results of the ensemble, it is important to select high-quality candidates. Given that the best pivot language for translation varies with each source sentence, we hence select the top-*k* candidates for each source sentence via quality estimation (QE). By leveraging diverse candidates from pivot translation and knowledge of the merging module, PIVOTE generates final translations that accurately convey the meaning and subtle nuances of the source sentence, superior to selecting from pre-existing candidates. Our contributions can be summarized as follows:

- We propose a simple but effective pivot-based single model ensemble method, PIVOTE, to improve low-resource MT.
- We show that a single model can effectively generate diverse and accurate hypotheses and that leveraging these candidates in an ensemble process can enhance translation quality while reducing computational overhead.
- The empirical results on various language pairs demonstrate that we consistently outperform state-of-the-art methods, validating the effectiveness of the pivot-based ensemble.

2 Related Work

Pivot-based approaches. Pivot translation is an 123 approach that decomposes the translation task into 124 two sequential steps (Wu and Wang, 2007; Utiyama 125 and Isahara, 2007). By transferring knowledge 126 127 from high-resource pivot languages, pivoting is especially effective in translation between low-128 resource languages (Zoph et al., 2016; Aji et al., 129 2020; He et al., 2022). In this study, pivot translation enables us to obtain high-quality candidates 131

for the ensemble. Kim et al. (2019) discusses a pivot-based transfer learning technique where source \rightarrow pivot and pivot \rightarrow target models are first trained separately, then use pre-trained models to initialize the source \rightarrow target model, allowing effective training of a single, direct NMT model. Zhang et al. (2022) further investigate the transfer learning approach by utilizing auxiliary monolingual data. 132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

155

156

157

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

Pivot translation typically employs English as the bridge language. Nonetheless, previous studies have explored the use of diverse pivot languages, taking into account factors such as data size and the relationships between languages (Paul et al., 2009; Dabre et al., 2015). By leveraging the ability of pivot translation to produce diverse outputs, several studies have focused on generating paraphrases (Mallinson et al., 2017; Guo et al., 2019). More recently, Mohammadshahi et al. (2024) uses pivot translation for ensemble, but it requires computing token-level probabilities and fails to improve translation. Our work shares the motivation with these studies, generating translations depending on the pivot path to obtain a variety of candidates.

Ensemble in NLG tasks. Ensemble learning is a widely adopted strategy to obtain more accurate predictions by employing multiple systems (Sagi and Rokach, 2018). In NMT, the traditional approach involves averaging the probability distributions of the next target token, which is predicted at each decoding step by multiple models (Bojar et al., 2014) or by different snapshots (Huang et al., 2017). When multiple sources are available, an ensemble can be conducted with predictions obtained by different sources (Firat et al., 2016). Also, a token-level ensemble through vocabulary alignment across LLMs has also been proposed (Xu et al., 2024). However, these methods are not applicable to recent black-box models as they cannot compute token-level probabilities at decoding time.

Selection-based ensemble has also been explored, which chooses the final output among the existing candidates. This can be achieved through majority voting by selecting the most frequent one (Wang et al., 2022) or selecting the best candidate with QE (Fernandes et al., 2022; Hendy et al., 2023). Recently, MBR decoding (Goel and Byrne, 2000; Farinhas et al., 2023), which aims to find the hypothesis with the highest expected utility, has gained attention. However, this approach limits the final output space to the existing candidate pool.

On the other hand, the generation-based ensem-



Figure 1: Overview of PIVOTE framework.

ble method involves generating a new final predic-183 184 tion. Fusion-in-Decoder (Izacard and Grave, 2021) proposes an architecture that aggregates additional 185 information with a given input. More recently, within the context of LLMs, Jiang et al. (2023) 187 and Yin et al. (2023) investigate a method of using 188 LLMs to generate multiple outputs and aggregate 189 them. Generating new output through LLMs offers 190 191 the benefit of explicitly harnessing their pre-trained knowledge within the ensemble process. 192

3 Pivot-based Single Model Ensemble

In this section, we first introduce the overview of PIVOTE framework (§3.1). Then, we describe the candidate generation process through pivot translation (§3.2) and the aggregation process (§3.3).

3.1 Overview

193

194

195

196

198

199

200

202

207

210

211

212

213

214

215

216

218

Our objective is the same as that of conventional translation tasks: converting the given source language sentence x into the target language sentence \hat{y} . PIVOTE consists of two steps: candidate generation and candidate aggregation. Figure 1 illustrates an overview of the proposed ensemble framework.

As the first step, we input x to generate candidates through a single multilingual NMT model. One translation path could be directly translating from the source to the target through the source \rightarrow target path. Alternatively, pivot translations can be achieved by employing high-resource pivot languages, enabling translation paths from source \rightarrow pivot and pivot \rightarrow target. During the pivot process, leveraging abundant parallel data enables knowledge transfer from high-resource pivot languages, thereby facilitating the generation of diverse and more accurate translations. Through these n paths, we can obtain a candidate pool $C = \{c_1, ..., c_n\}$ composed of n candidates in the target language, employing only a single model.

As the second step, a ranking process is first conducted within the candidate pool C since not all candidates contribute to the ensemble. Using the estimated quality of each candidate, we select the top-k candidates. We then generate the final output \hat{y} using the selected high-quality candidates. This generation-based approach facilitates the production of outputs superior to existing candidates.

3.2 Pivot-based Candidate Generation

In the first step, PIVOTE takes a source sentence x as input and generates n candidates. Direct translation yields only one candidate, whereas pivot translation enables the generation of multiple candidates from a single source sentence using a single model. Generating candidates through pivot translation has two major advantages: diversity and quality.

First, we can obtain diverse candidates that can act complementarily. One of the key principles for the ensemble is that the participants must be sufficiently diverse to provide various inductive biases. In PIVOTE, each source sentence is translated diversely by passing through multiple translation paths. Diverse translation paths enhance the likelihood of providing expressions that convey the accurate meaning of the source sentence. Pivotbased candidate generation shares a similar goal with a previous study that generates paraphrases through round-trip translation, aiming to generate diverse translations (Thompson and Post, 2020).

Second, by utilizing a parallel corpus of high-resource pivot languages, pivoting enables more accurate translations. For low-resource language pairs, more appropriate translations can be achieved through two-step decoding through a pivot language (He et al., 2022). Moreover, leveraging pivot languages with abundant parallel data, not limited to English, allows us to obtain better

256

345

346

347

348

349

351

352

353

354

356

357

307

308

translations (Paul et al., 2009; Dabre et al., 2015).

257

262

265

266

267

271

272

295

In addition, pivot translation with a single model offers practical benefits over employing multiple models. Firstly, it can reduce the costs of operating multiple models including LLMs. Secondly, the substantial performance disparities among models mean that using the top-performing single model for candidate generation often leads to higherquality outcomes. Lastly, it reduces inference latency by using a single model for two batched inferences, while multi-model ensembles require up to 11, causing significant overhead and limiting realtime response capability. Given that pivot translation with a single model allows for the creation of diverse and more accurate translations, we utilize an MNMT model to generate the candidates.

Selecting pivot languages. For each language pair, 273 we carefully select pivot languages based on the assumption that pivot language with abundant mutual 275 knowledge would allow us to obtain higher-quality 276 candidates. We select n top-performing paths for 277 our study based on BLEU scores on the FLORES-278 200 benchmark (Costa-jussà et al., 2022). We evaluate the outputs for each path, including direct translation and through various pivot translations. 281 NLLB (Costa-jussà et al., 2022) is used to generate candidates, and results on the FLORES-200 for selecting translation paths are in Appendix A. If pivot languages are selected based on BLEU scores, high-resource languages are predominantly chosen, rather than low-resource ones. The experiments 287 detailed in Appendix B demonstrate that overly prioritizing diversity by employing low-resource pivot 289 languages, at the expense of candidate quality, does not result in improvements in the final translation. The experiments comparing metrics for selecting translation paths are in Appendix C. As a result, we compose the candidate pool using the 4 paths. 294

3.3 Candidate Aggregation

296In the aggregation step, we take the candidate pool297C as input and output the merged final translation \hat{y} .298The post-hoc aggregation process encompasses two299stages: selecting and merging. In the first stage, we300select candidates by ranking method. There are two301approaches for selecting candidates. One approach302evaluates each translation path and selects the best303paths for all source sentences. The other approach304involves selecting the best top-k candidates for each305source sentence. After selecting k candidates, we306generate the final translation \hat{y} using the merging

module. This process enables the creation of better outputs beyond the quality of existing candidates.

Selecting the top-k candidates. The pivot language that generates the highest-quality candidate varies for each source sentence. The best output is not guaranteed from one translation path alone, as it can vary depending on factors such as the size of the parallel corpus and the relationship between languages. First, PIVOTE uses QE to rank all ncandidates from candidate pool $C = \{c_1, ..., c_n\}$. Afterward, we select top-k candidates among n candidate pool. Selecting the top-k candidates ensures the quality of the output by filtering out low-quality candidates while also efficiently reducing the cost during the merging process. We use the referencefree COMETkiwi (*wmt22-COMETkiwi-da*) (Rei et al., 2022b) for ranking candidates.

Generating the final translation. To generate the final translation \hat{y} by merging the top-k candidates, we explore methods from two categories: encoder-decoder ensemble architectures and LLM-based approach. Employing encoderdecoder architectures during the merging process offers the advantage of relatively low training costs. We conduct experiments using Fusionin-Decoder (FiD) (Izacard and Grave, 2021) and TRICE (Huang et al., 2021) architectures. The former method involves passing Translate source into <target language> referring <target language> candidate. source: <x> candidate: $\langle c_k \rangle$ through the encoder, representations are concatenated and merged in the decoder. The latter approach involves concatenating $<x></s><l_s>;<c_1></s><l_t>;...;<c_k></s><l_t>$ with language token $\langle l_{lang} \rangle$ and providing it as input. Encoder-decoder ensemble architectures are further described in detail in Appendix D.

On the other hand, the LLM-based ensemble implicitly leverages their translation capabilities during ensemble, as the source sentence is also provided. We conduct merging experiments with GEN-FUSER (Jiang et al., 2023), Llama-3 (AI@Meta, 2024), and GPT models (Ouyang et al., 2022; OpenAI, 2023, 2024). When employing GENFUSER, we construct the input by concatenating top-*k* candidates to the prompt, as presented in Jiang et al. (2023). For merging with Llama-3 and GPT, we use the prompt template in Appendix E. By leveraging a variety of candidates, each with different strengths, the aggregation process can effectively mitigate errors in a complementary manner.

360

373

376

377

384

391

.

4

-

periments.

Experiments

F			
4.1	Datasets		

We use NVIDIA RTX 3090 or 4090 GPUs for ex-

Long poir	Detect	# Sentences				
Lang-pan	Dataset	Train	Dev	Test		
$KO \leftrightarrow IT$	TED 2020 v1 (Reimers and Gurevych, 2020)	357,733	2,000	2,000		
$AR \leftrightarrow PT$	WikiMatrix v1 (Schwenk et al., 2019)	153,441	2,000	2,000		

Table 1: Datasets statistics.

We conduct experiments on the linguistically distant languages within pairs: not in the same language family and using different scripts. We select 2 language pairs, resulting in 4 translation directions in total, Korean (*Koreanic*) \leftrightarrow Italian (*Romance*) and Arabic (*Arabic*) \leftrightarrow Portuguese (*Romance*). The language family grouping is defined by the criteria presented in Fan et al. (2020).

We validate our approach across various domains. For Korean↔Italian pair, we run experiments on TED2020 (Reimers and Gurevych, 2020). For Arabic↔Portuguese pair, we use WikiMatrix (Schwenk et al., 2019). All the datasets are obtained from the OPUS¹ (Tiedemann, 2012) project. The statistics for the datasets are listed in Table 1.

4.2 Evaluation Metrics

We assess the translation quality using BLEU (Papineni et al., 2002), chrF++ (Popović, 2017), and reference-based COMET (*wmt22-COMET-da*) (Rei et al., 2022a). For reporting BLEU, *SacreBLEU* (Post, 2018) is used with ko-mecab tokenizer for Korean and 13a tokenizer for the others.

4.3 Baselines

As an encoder-decoder NMT model, we use NLLB-200-distilled-600M (Costa-jussà et al., 2022). When training NLLB, we use the Transformers library from HuggingFace (Wolf et al., 2020). AdamW optimizer (Loshchilov and Hutter, 2019) is used with a learning rate of 2e-5, batch size of 2, and dropout with a probability of 0.1. When validation BLEU was not improved for 3 checkpoints, with 30k steps between them, we stopped training.

For open-source LLMs, we use Vicuna 13B (Chiang et al., 2023), Baize 13B (Xu et al., 2023a), and Llama-3-8B-Instruct (AI@Meta, 2024) as the baseline. We fine-tuned these LLMs with

¹https://opus.nlpl.eu

QLoRA (Dettmers et al., 2024); r=16, $\alpha=64$, dropout=0.1 for all linear layers. For blackbox LLMs, we use GPT-4 (OpenAI, 2023) and GPT-40 (OpenAI, 2024). The version of gpt-4-1106-preview and gpt-40-2024-08-06 are employed for GPT-4 and GPT-40, respectively. For GPT models, *temperature* is set to 0.0 for stable responses (Peng et al., 2023) and *top_p* is set to 0.1 to ensure reproducibility. For LLMs, we use the prompt template of Hendy et al. (2023), as presented in Appendix E. 398

399

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

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

As state-of-the-art ensemble baselines, we employ LLM-Blender (Jiang et al., 2023), EVA (Xu et al., 2024), and MBR (Farinhas et al., 2023). For LLM-Blender and EVA, we fine-tuned the same open-source LLMs used in each study utilizing the parallel corpus described in Table 1. The list of the LLMs is in Appendix F. *temperature* is set to 0.1 to mitigate hallucination for low-resource pairs (Guerreiro et al., 2023). For MBR, we generate a set of 5 hypotheses using GPT-4. When generating hypotheses, *temperature* was set to 0.0 for its optimal performance, based on the results of our pilot experiments in Appendix G and prior study (Peng et al., 2023). Other configurations are the same as in the original work (Farinhas et al., 2023).

4.4 Implementation Details

In the candidate generation step of PIVOTE, we employ NLLB. For each source-target language pair, we use an NLLB fine-tuned for the language pair in Table 1 to generate the directly translated candidates. For the merging module, we use Llama-3, GPT-4, and GPT-40. For all models used in PIVOTE, including NLLB, Llama-3, GPT-4, and GPT-40, we apply the same settings in §4.3.

As detailed in §3.3, we explore two approaches in the ensemble process: one dynamically selects the top-k (k=3) candidates and another uses candidates obtained from fixed paths. To select the top-kcandidates for each source sentence, we use the reference-free COMETkiwi as described in §3.3. When selecting candidates from fixed paths, we used directly translated candidates and Englishpivoted candidates, which were the top-performing paths on the FLORES-200 benchmark.

4.5 Main Results

Table 2 reports the overall performance of PIVOTEand other methods. The results demonstrate thatPIVOTE consistently outperforms baselines acrossall language pairs. While standalone NMT systems

Madal	Korean→Italian		Italian→Korean		Arabic → Portuguese			Portuguese → Arabic				
Widder	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
Standalone NMT System												
NLLB (Costa-jussà et al., 2022)	16.27	41.14	84.60	17.40	23.39	87.33	27.25	50.35	84.21	13.50	40.90	84.24
Vicuna (Chiang et al., 2023)	10.11	31.15	70.29	10.60	16.51	72.29	17.64	38.44	76.01	8.40	27.38	79.18
Baize (Xu et al., 2023a)	10.62	31.87	73.62	10.38	16.44	76.63	16.56	36.67	76.87	8.50	27.28	79.18
Llama-3 (AI@Meta, 2024)	11.79	34.82	77.37	13.82	18.95	85.80	18.78	40.20	78.73	12.25	35.16	82.79
GPT-4 (OpenAI, 2023)	14.07	42.22	86.80	17.23	22.96	86.94	25.82	51.89	85.46	15.11	41.39	83.99
GPT-40 (OpenAI, 2024)	15.11	42.59	85.93	17.20	22.82	85.31	27.28	52.57	85.90	16.28	42.40	83.82
Prior Ensemble Method												
LLM-Blender (Jiang et al., 2023)	8.77	28.74	82.80	0.03	0.85	42.77	11.80	29.85	67.95	0.94	2.69	46.49
EVA (Xu et al., 2024)	2.53	15.26	39.00	1.51	3.57	37.17	9.77	28.40	68.75	7.99	27.00	73.15
MBR (Farinhas et al., 2023)	14.10	42.24	86.70	17.14	23.00	87.53	25.45	51.78	85.55	14.66	41.11	83.93
Proposed Method												
PIVOTE (Llama-3; top3)	15.60	39.86	84.10	14.56	19.92	87.34	23.41	45.95	81.66	14.27	38.25	81.80
PIVOTE (Llama-3; D, E)	13.85	37.36	69.96	14.97	20.21	85.42	21.35	43.75	79.71	12.37	36.51	82.09
PIVOTE (GPT-4; top3)	16.66	42.85	86.82	17.95	23.84	87.50	27.22	51.73	85.65	16.53	42.41	84.46
PIVOTE (GPT-4; D, E)	17.10	43.29	85.92	18.18	24.05	88.74	27.98	52.41	85.27	17.02	43.02	84.82
PIVOTE (GPT-40; top3)	17.77	43.38	85.46	18.08	23.98	88.15	28.62	52.53	85.87	16.92	42.93	84.52
PIVOTE (GPT-40; D, E)	18.02	43.46	86.19	18.31	24.32	88.33	29.50	53.16	86.03	17.66	43.73	84.27

Table 2: Main results. The best scores in each pair are marked **bold**. Within parentheses in the proposed method, the parts separated by semicolons denote the merging module and the candidates used. *D* and *E* represent candidates obtained from direct translation and English pivot, respectively.

Model	Korean→Italian					
Widdel	BLEU	chrF++	COMET			
Candidate						
NLLB (direct)	16.27	41.14	84.60			
NLLB (Portuguese pivot)	13.13	37.57	83.21			
NLLB (Spanish pivot)	13.87	38.47	83.71			
NLLB (English pivot)	14.77	39.39	81.48			

Table 3:	Quality	of c	candidates	used	for the	ensemble.

rely solely on their pre-trained knowledge, PIVOTE explicitly leverages candidates during the ensemble. Even when training an open-source LLM, Llama-3, we can enhance translation capability by utilizing candidates obtained via pivoting. Compared to using LLMs for translation, we can improve performance with only the minimal cost of utilizing a small 0.6B model. Table 3 presents the quality of candidates utilized in the ensemble. We will further elaborate with a case study, showing that PIVOTE achieves better translations by leveraging candidates to capture subtle nuances of the source sentence. We report experiments in a setting that does not use training data in Appendix H and experiments with other GPT models in Appendix K. The analysis of the proportion of top-k candidates and performance variation with k are in Appendix J.

Comparison with multi-model ensemble. We compare PIVOTE with LLM-Blender (Jiang et al., 2023) and EVA (Xu et al., 2024), state-of-theart ensemble methods utilizing multiple models. LLM-Blender employs N (N=11) LLMs for candidate generation, picks top-3 candidates with PAIR-RANKER, and fuses them with GENFUSER. EVA is a token-level ensemble method that leverages vocabulary alignment across multiple models. Results in Table 2 show that PIVOTE outperforms multi-model ensemble baselines by a considerable margin. LLM-Blender was unable to improve outputs compared to its candidate LLMs in non-English translation tasks. Additionally, LLMs used for generating candidates in LLM-Blender, such as Vicuna and Baize, exhibit subpar performance on given tasks. These results align with recent work (Xu et al., 2023b); open-source LLMs often struggle when not translating into English. 475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

EVA is not only ineffective on the given tasks but also has several limitations inherent to its design as a token-level ensemble. First, EVA is unable to use black-box models such as GPT-4. Second, it is memory-intensive, as it requires loading multiple models into memory simultaneously. While multimodel ensemble methods generate candidates using up to 11 LLMs (with sizes up to 13B), PIVOTE generates candidates with a significantly smaller single model (0.6B), thereby greatly reducing computational overhead.

Results on all language pairs. To validate generalizability, we report the results for all language pairs we experimented with, including those within the same language family. Distant pairs refer to languages that belong to different families and use different scripts, while similar pairs belong to the same family and share the same script. The statistics for each language pair are in Appendix I. Language pairs used in the experiments are as follows:

- Distant language pairs: Portuguese↔Russian, Dutch↔Russian, and French↔Ukrainian
- Similar language pairs: Spanish↔Portuguese 506 and Ukrainian↔Russian 507

472

473

474

#	Туре	Example
	Source Sentence	그래서 그동안 자문해왔습니다. 왜 우리는 질병들과 싸우기에 더 현명하고, 정확하며 더욱 적합한
		(English Translation: So we've been asking ourselves, why should we limit this smarter, more precise, more appropriate)
	Target Reference	Quindi mi sono chiesta: perché dovremmo limitare questo modo intelligente, preciso, migliore
1	Top-1 Candidate	Quindi ci siamo chiesti, perché dovremmo limitare questo modo più intelligente, più preciso e più appropriato
1	Top-2 Candidate	Quindi ci siamo chiesti: perché dovremmo limitare questo metodo più intelligente, più preciso e più adatto
	Top-3 Candidate	Quindi nel corso di questo tempo, abbiamo chiesto: perché dovremmo limitare questo metodo più intelligente, più preciso e più adeguato
	GPT-4	Quindi, abbiamo cercato consigli fino ad ora. Perché dobbiamo limitare questo metodo, che è più intelligente, preciso e più adatto
	ΡινότΕ	Quindi ci siamo chiesti: perché dovremmo limitare questo metodo più intelligente, più preciso e più adatto
	Source Sentence	많은 사람들이 헤드폰을 사용하는데 이것의 문제점은 3가지 큰 건강 이슈를 가져온다는 것입니다.
		(English Translation: The trouble with widespread headphone use is it brings three really big health issues.)
	Target Reference	Il problema dell'utilizzo diffuso degli auricolari è che scatenano tre grandi problemi di salute.
2	Top-1 Candidate	Il problema è che molte persone usano le cuffie, e questo Porta a tre grandi problemi di salute.
4	Top-2 Candidate	Il problema è che molte persone usano le cuffie, e questo Porta a tre grandi problemi di salute.
	Top-3 Candidate	Il problema è che molte persone usano cuffie, e questo è ciò che causa tre problemi di salute principali.
	GPT-4	Molte persone utilizzano le cuffie, ma il problema è che ciò comporta tre importanti questioni di salute.
	PIVOTE	Il problema è che molte persone usano le cuffie, e questo porta a tre grandi problemi di salute.

Table 4: Case study. Parts with the same meanings as the source and mistranslated parts are highlighted in blue and red, respectively. English translation of the source sentence is obtained from another pair within the same dataset.

Model	BLEU	chrF++	COMET	BLEU	chrF++	COMET		
	Distant Language Pairs							
	Port	uguese→F	Russian	Rus	sian→Port	uguese		
NLLB	25.17	51.77	90.12	29.69	55.81	86.01		
GPT-4	26.50	52.76	91.11	25.51	54.05	86.69		
PivotE	27.48	53.49	91.74	30.82	56.73	88.37		
	D	utch->Rus	sian	R	ussian→D	utch		
NLLB	22.95	50.21	89.92	25.56	53.60	88.18		
GPT-4	24.37	51.32	91.28	24.46	53.85	88.58		
PivotE	25.45	52.16	91.47	28.05	55.80	89.35		
	Fre	nch→Ukr	ainian	Ukrainian→French				
NLLB	14.58	37.11	82.99	20.69	44.04	80.61		
GPT-4	13.84	39.03	84.12	23.30	47.13	83.43		
ΡινότΕ	17.20	39.82	86.55	24.35	47.17	84.36		
		Simil	ar Languag	e Pair (Re	omance)			
	Spai	nish→Port	uguese	Port	uguese→S	panish		
NLLB	32.38	56.97	86.88	33.63	57.61	85.13		
GPT-4	29.94	55.26	84.84	34.70	58.63	86.75		
PivotE	34.06	58.11	87.70	36.03	59.32	86.92		
	Similar Language Pair (Slavic)							
	Ukı	ainian→R	ussian	Rus	sian→Ukr	ainian		
NLLB	22.16	45.41	89.82	19.67	43.35	89.87		
GPT-4	24.41	47.59	89.43	22.42	45.61	90.39		
PivotE	24.64	47.51	90.78	22.09	45.40	90.70		

Table 5: Results on all language pairs.

Table 5 shows the results with the topperforming baselines, NLLB (Costa-jussà et al., 2022) and GPT-4 (OpenAI, 2023). PIVOTE consistently exhibits superior performance compared to strong baselines on distant language pairs. Surprisingly, it also showed improvements in similar language pairs, such as Spanish↔Portuguese.

510

511

512

513

514

515

516

517

518

519

520

522

523

524

525

526

Case study. We conduct a qualitative analysis to verify the impact of candidates on the final translation. We compare the output of GPT-4, used as the merging module, with PIVOTE, which utilizes candidates for the ensemble process. In Table 4, we provide two examples along with the source and target sentences, as well as the top-3 candidates.

Through the first example, we can observe that PIVOTE can appropriately translate homonyms within the context. In Korean, "자문" has the meaning of both "consultation" and "asking oneself". Considering the context, the expression should be

Condidate Concretion	# Cand	Korean→Italian			
Candidate Generation	# Cand.	BLEU	chrF++	COMET	
LLMs (of LLM-Blender)	11	14.75	41.29	86.20	
LLMs + NLLB (direct)	12	16.08	42.38	86.22	
NLLB (pivot, ours)	4	16.66	42.85	86.82	

Table 6: Comparison of candidate generation methods.

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

translated to convey the meaning of "asking ourselves", as also shown in the English translation. However, GPT-4 mistranslated the source sentence, converting the phrase "자문해왔습니다" to "abbiamo cercato consigli" ("seeking consultation from others"). On the other hand, PIVOTE accurately translates with the expression "ci sono chiesti" that means "asking ourselves", aligning well with the context by leveraging information from candidates.

In the second sample, GPT-4 translates the source sentence by translating the noun "이合" into "questioni". However, given the topic of discussing potential health risks, this translation does not fit well with the overall context. By contrast, the ensemble result of PIVOTE, generated using the identical model, improves translation quality by using a more accurate expression "problemi", despite having access to the same pre-trained knowledge. Additionally, when more suitable expressions (e.g., "ne vale la pena") appear in candidates, PIVOTE utilizes them to refine the final translation.

4.6 Analysis

Candidate generation. To validate the effectiveness of PIVOTE, we conduct experiments only varying the candidate generation method, while using the same merging module, GPT-4. We compare a candidate pool of size 4 obtained through pivot translation (PIVOTE) with a candidate pool

Model	Korean→Italian			
Widder	BLEU	chrF++	COMET	
Standalone NMT System				
NLLB (Costa-jussà et al., 2022)	16.27	41.14	84.60	
Encoder-Decoder				
PIVOTE (FiD)	13.74	36.78	78.98	
PIVOTE (TRICE)	15.89	41.98	84.06	
LLM-based				
PIVOTE (GENFUSER)	14.56	39.32	80.07	
PIVOTE (GPT-4)	16.66	42.85	86.82	

Table 7: Evaluation of merging module variants.

of size 11 obtained using 11 LLMs as employed in LLM-Blender (Jiang et al., 2023).

As shown in Table 6, the proposed method of generating candidates through pivot translation achieves the highest performance, despite using the smallest candidate pool. From the perspective of direct translation in NLLB, leveraging 3 candidates obtained through pivot translation yields higher scores than incorporating candidates generated by 11 LLMs. These results demonstrate that using stable-quality candidates generated by a single model via pivot translation outperforms the use of multiple models with performance disparities.

Candidate aggregation. We first investigate whether PIVOTE shows improvement when utilizing other merging modules. As detailed in §3.3, we run experiments with three architectures: FiD (Izacard and Grave, 2021), TRICE (Huang et al., 2021), and GENFUSER (Jiang et al., 2023). When implementing FiD, we replace the backbone of FiD to mT5_{BASE} (Xue et al., 2021). TRICE is a method proposed for multi-source translation. Since TRICE was not originally intended for ensemble use, we repurposed it by training on the following two tasks: The first task is the original translation which converts source sentences into target sentences. The second task is refining candidates that are paired with target references. In the case of TRICE, only the highest quality candidates, which are the directly translated ones, are used due to its architecture. FiD and GENFUSER use top-3 candidates.

Table 7 shows that the ensemble methods using encoder-decoder architectures and GENFUSER do not yield improved results. These methods struggle to leverage additional information from the candidates and, consequently, do not enhance performance. In contrast, using GPT-4 as the merging module leads to better performance compared to the standalone NMT system. We also compare ranking methods COMETkiwi and PAIR-RANKER (Jiang et al., 2023). While the perfor-

Cotogory	Mathad	Korean→Italian			
Category	Methou	BLEU	chrF++	COMET	
	PAIRRANKER	15.61	40.62	84.46	
(top 1)	COMETkiwi	15.61	40.71	84.10	
(top-1)	COMET* (ideal)	17.77	42.81	84.83	
Generation-based	ΡινότΕ	16.66	42.85	86.82	

Table 8: Comparison with selection-based ensemble. Note that COMET* is the ideal baseline, as it requires references. Best scores including COMET* are **bolded**, while best scores excluding it are underlined.

mance is comparable, considering the efficiency factor, we opt for COMETkiwi. Detailed experiments about the ranking method are in Appendix L.

596

597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

Comparison with selection-based ensemble. With a selection-based ensemble, we can choose one of the existing candidates as the final translation, rather than generating a new one. In this experiment, we compare our approach with a selection-based ensemble by selecting the top-1 translation using PAIRRANKER (Jiang et al., 2023) and COMETkiwi (Rei et al., 2022b). Additionally, we report results with an ideal case: selecting top-1 by considering references as well, which are not available in practice. The ideal top-1 is selected by reference-based COMET (Rei et al., 2022a).

As shown in Table 8, PIVOTE exhibits superior performance compared to the selectionbased ensemble methods. Even when we leverage reference-based COMET, which is impossible in real-world scenarios due to the necessity for references, PIVOTE outperforms it in chrF++ and COMET. These results indicate that performing a generation-based ensemble with pivoting can effectively produce final translations that surpass those selected from the existing candidate pool.

5 Conclusion

In this work, we introduced PIVOTE, a pivot-based single model ensemble framework, to enhance translation in scenarios where parallel data are scarce. By transferring knowledge from diverse pivot languages, we were able to obtain not only diverse but also high-quality candidates. And the optimal path to generating the best candidate varies per sentence, our study underscores the significance of exploiting a spectrum of pivot languages. Moreover, the single model generation process offers cost savings compared to multi-model ensemble approaches. Empirical results and qualitative analyses show that the proposed method can yield contextually suitable translations for the given source sentences by leveraging pivoted candidates.

585

591

592

595

555

743

744

745

689

690

Limitations

637

665

674

675

677

678

679

638Despite PIVOTE utilizes candidates obtained via639pivoting, limitations arise from the nature of pivot640translation. Constraining the pivot language to high-641resource languages can limit the number of candi-642dates because pivoting through low-resource lan-643guages can lead to some loss of information due644to error propagation inherent in the two-step trans-645lation. This semantic shift potentially causes a de-646crease in candidate quality. If the quality of candi-647dates declines, improvements from the ensemble648might not be significant, indicating a limitation in649the number of pivot paths.

References

AI@Meta. 2024. Llama 3 model card.

- Alham Fikri Aji, Nikolay Bogoychev, Kenneth Heafield, and Rico Sennrich. 2020. In neural machine translation, what does transfer learning transfer? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7701–7710, Online. Association for Computational Linguistics.
- Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2018. Unsupervised neural machine translation. *Preprint*, arXiv:1710.11041.
- Baichuan. 2023. Baichuan 2: Open large-scale language models. *Preprint*, arXiv:2309.10305.
- Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. Findings of the 2014 workshop on statistical machine translation. In Proceedings of the Ninth Workshop on Statistical Machine Translation, pages 12–58, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Ye Chen, Wei Cai, Liangmin Wu, Xiaowei Li, Zhanxuan Xin, and Cong Fu. 2023. Tigerbot: An open multilingual multitask llm. *Preprint*, arXiv:2312.08688.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%* chatgpt quality.
- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi,

Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *ArXiv preprint*, abs/2210.11416.

- Mike Conover, Matt Hayes, Ankit Mathur, Xiangrui Meng, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instruction-tuned llm.
- Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Raj Dabre, Fabien Cromieres, Sadao Kurohashi, and Pushpak Bhattacharyya. 2015. Leveraging small multilingual corpora for SMT using many pivot languages. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1192–1202, Denver, Colorado. Association for Computational Linguistics.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 320–335, Dublin, Ireland. Association for Computational Linguistics.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation. *Preprint*, arXiv:2010.11125.
- António Farinhas, José de Souza, and Andre Martins. 2023. An empirical study of translation hypothesis ensembling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11956–11970, Singapore. Association for Computational Linguistics.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412,

- 746 747 748 752 753 754 755 756 759 760 761 765 769 770 771 773 774 775 776 781 790 791

- 796 797 798

- Seattle, United States. Association for Computational Linguistics.
 - Orhan Firat, Baskaran Sankaran, Yaser Al-onaizan, Fatos T. Yarman Vural, and Kyunghyun Cho. 2016. Zero-resource translation with multi-lingual neural machine translation. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 268-277, Austin, Texas. Association for Computational Linguistics.
 - Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. Koala: A dialogue model for academic research. Blog post.
 - Vaibhava Goel and William J Byrne. 2000. Minimum bayes-risk automatic speech recognition. Computer Speech Language, 14(2):115–135.
 - Nuno M. Guerreiro, Duarte Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André F. T. Martins. 2023. Hallucinations in large multilingual translation models. Preprint, arXiv:2303.16104.
 - Yinpeng Guo, Yi Liao, Xin Jiang, Qing Zhang, Yibo Zhang, and Qun Liu. 2019. Zero-shot paraphrase generation with multilingual language models. Preprint, arXiv:1911.03597.
 - Zhiwei He, Xing Wang, Zhaopeng Tu, Shuming Shi, and Rui Wang. 2022. Tencent AI lab - shanghai jiao tong university low-resource translation system for the WMT22 translation task. In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 260-267, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
 - Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. Preprint, arXiv:2302.09210.
 - Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. 2017. Snapshot ensembles: Train 1, get m for free. Preprint, arXiv:1704.00109.
 - Xuancheng Huang, Jingfang Xu, Maosong Sun, and Yang Liu. 2021. Transfer learning for sequence generation: from single-source to multi-source. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 5738-5750, Online. Association for Computational Linguistics.
 - Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 874-880, Online. Association for Computational Linguistics.

Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. LLM-blender: Ensembling large language models with pairwise ranking and generative fusion. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 14165–14178, Toronto, Canada. Association for Computational Linguistics.

803

804

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

- Yunsu Kim, Petre Petrov, Pavel Petrushkov, Shahram Khadivi, and Hermann Ney. 2019. Pivot-based transfer learning for neural machine translation between non-English languages. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 866-876, Hong Kong, China. Association for Computational Linguistics.
- LAION-AI. 2023. Open assistant. https://github. com/LAION-AI/Open-Assistant.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In International Conference on Learning Representations.
- Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2017. Paraphrasing revisited with neural machine translation. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 881-893, Valencia, Spain. Association for Computational Linguistics.
- Alireza Mohammadshahi, Jannis Vamvas, and Rico Sennrich. 2024. Investigating multi-pivot ensembling with massively multilingual machine translation models. In Proceedings of the Fifth Workshop on Insights from Negative Results in NLP, pages 169-180, Mexico City, Mexico. Association for Computational Linguistics.
- NLP Team MosaicML. 2023. Introducing mpt-7b: A new standard for open-source, ly usable llms. Accessed: 2023-05-23.
- Yasmin Moslem, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models. In Proceedings of the 24th Annual Conference of the European Association for Machine Translation, pages 227–237, Tampere, Finland. European Association for Machine Translation.
- OpenAI. 2023. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- OpenAI. 2024. Gpt-40 system card. Preprint, arXiv:2410.21276.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730-27744.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-

Jing Zhu. 2002. Bleu: a method for automatic evalu-

ation of machine translation. In Proceedings of the

40th Annual Meeting of the Association for Com-

putational Linguistics, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational

Michael Paul, Hirofumi Yamamoto, Eiichiro Sumita,

and Satoshi Nakamura. 2009. On the importance of

pivot language selection for statistical machine trans-

lation. In Proceedings of Human Language Tech-

nologies: The 2009 Annual Conference of the North American Chapter of the Association for Computa-

tional Linguistics, Companion Volume: Short Papers,

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen,

Maja Popović. 2017. chrF++: words helping charac-

ter n-grams. In Proceedings of the Second Confer-

ence on Machine Translation, pages 612-618, Copen-

hagen, Denmark. Association for Computational Lin-

Matt Post. 2018. A call for clarity in reporting BLEU

scores. In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 186-

191, Belgium, Brussels. Association for Computa-

Ricardo Rei, José G. C. de Souza, Duarte Alves,

Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova,

Alon Lavie, Luisa Coheur, and André F. T. Martins.

2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Proceedings of the

Seventh Conference on Machine Translation (WMT),

pages 578–585, Abu Dhabi, United Arab Emirates

(Hybrid). Association for Computational Linguistics.

Chrysoula Zerva, Ana C. Farinha, Christine Maroti,

José G. C. de Souza, Taisiya Glushkova, Duarte M.

Alves, Alon Lavie, Luisa Coheur, and André F. T.

submission for the quality estimation shared task.

Nils Reimers and Iryna Gurevych. 2020. Making

monolingual sentence embeddings multilingual us-

ing knowledge distillation. In Proceedings of the

2020 Conference on Empirical Methods in Natural

Language Processing. Association for Computational

Nathaniel Robinson, Perez Ogayo, David R. Mortensen,

and Graham Neubig. 2023. ChatGPT MT: Competi-

tive for high- (but not low-) resource languages. In

Proceedings of the Eighth Conference on Machine

Translation, pages 392-418, Singapore. Association

Cometkiwi: Ist-unbabel 2022

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,

of chatgpt for machine translation.

Xuebo Liu, Min Zhang, Yuanxin Ouyang, and

Dacheng Tao. 2023. Towards making the most

Preprint,

Linguistics.

pages 221-224.

arXiv:2303.13780.

tional Linguistics.

Martins. 2022b.

Linguistics.

Preprint, arXiv:2209.06243.

guistics.

- 875
- 877
- 882

887

- 890
- 892
- 894

896 897

900

901 902

- 903
- 904
- 905 906

907 908

- 909 910
- 911
- 912
- 913
- for Computational Linguistics. 914

Omer Sagi and Lior Rokach. 2018. Ensemble learning: A survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 8.

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

- Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2019. Wikimatrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia. Preprint, arXiv:1907.05791.
- Stability-AI. 2023. Stablelm: Stability ai language https://github.com/stability-AI/ models. stableLM.
- Tianxiang Sun and Xipeng Qiu. 2023. Moss. https: //github.com/OpenLMLab/MOSS.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2023. Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding. In Findings of the Association for Computational Linguistics: ACL 2023, pages 4265-4293, Toronto, Canada. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https: //github.com/tatsu-lab/stanford_alpaca.
- Gemini Team. 2023. Gemini: A family of highly capable multimodal models. Preprint, arXiv:2312.11805.
- Brian Thompson and Matt Post. 2020. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. In Proceedings of the Fifth Conference on Machine Translation, pages 561–570, Online. Association for Computational Linguistics.
- Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12), Istanbul, Turkey. European Language Resources Association (ELRA).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,

Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. *Preprint*, arXiv:2307.09288.

971

972

973

974

975

976

977

978

981

982

983

993

994

995

996

997 998

1001

1005

1006

1008

1009

1010 1011

1012

1013

1015

1016

1017 1018

1019

1020

1021

1022 1023

1024

1026

1027

- Masao Utiyama and Hitoshi Isahara. 2007. A comparison of pivot methods for phrase-based statistical machine translation. In *Human Language Technologies* 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference, pages 484–491, Rochester, New York. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2022. Rationaleaugmented ensembles in language models. *Preprint*, arXiv:2207.00747.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface's transformers: State-of-the-art natural language processing. *Preprint*, arXiv:1910.03771.
- Hua Wu and Haifeng Wang. 2007. Pivot language approach for phrase-based statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 856–863, Prague, Czech Republic. Association for Computational Linguistics.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023a. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *ArXiv preprint*, abs/2304.01196.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023b. A paradigm shift in machine translation: Boosting translation performance of large language models. *Preprint*, arXiv:2309.11674.
- Yangyifan Xu, Jinliang Lu, and Jiajun Zhang. 2024. Bridging the gap between different vocabularies for LLM ensemble. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7140–7152, Mexico City, Mexico. Association for Computational Linguistics.
- 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, Online. Association for Computational Linguistics.
 - Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng Guo, Junqi Dai, Xuanjing Huang, and Xipeng Qiu.

2023. Exchange-of-thought: Enhancing large language model capabilities through cross-model communication. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15135–15153, Singapore. Association for Computational Linguistics. 1028

1029

1031

1032

1034

1035

1036

1037

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1054

- Meng Zhang, Liangyou Li, and Qun Liu. 2022. Triangular transfer: Freezing the pivot for triangular machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 2: Short Papers*), pages 644–650, Dublin, Ireland. Association for Computational Linguistics.
- Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.

Appendix

A Pivot Language Selection

Based on the results from the FLORES-200 (Costajussà et al., 2022) benchmark, we select the top-4 pivot paths as presented in Table 9. We utilize the full 2009 sentences as our test set: 997 sentences from the *dev* and 1012 sentences from the *devtest*. The pivot language pool is chosen as the *bridge languages* in Fan et al. (2020).

Divot I anguago		Lang	g-pair	
I Ivot Language	KO→IT	IT→KO	AR→PT	PT→AR
direct	14.02	18.63	27.15	15.22
arb_Arab	11.03	15.82	-	-
ben_Beng	10.79	15.44	18.65	9.76
ces_Latn	11.48	16.08	21.23	11.55
deu_Latn	12.49	17.11	22.62	12.56
ell_Grek	11.96	16.54	22.53	12.54
eng_Latn	14.82	19.34	28.40	15.92
fin_Latn	9.62	14.31	17.27	9.48
fra_Latn	13.55	17.27	24.96	13.77
heb_Hebr	10.42	14.37	20.31	10.94
hin_Deva	11.54	17.12	21.79	11.72
hun_Latn	10.54	14.96	18.64	9.65
ind_Latn	12.41	17.03	22.47	11.97
ita_Latn	-	-	24.70	14.09
jpn_Jpan	10.60	14.73	14.29	7.31
kor_Hang	-	-	16.09	7.67
lit_Latn	10.46	14.96	18.14	9.47
nld_Latn	12.27	17.10	23.22	12.94
pes_Arab	11.09	15.86	20.88	11.50
pol_Latn	11.54	15.86	21.14	11.60
por_Latn	13.80	18.01	-	-
rus_Cyrl	12.25	16.57	22.77	12.39
spa_Latn	13.89	18.39	26.60	14.91
swe_Latn	11.93	16.54	22.34	12.25
swh_Latn	10.66	14.22	19.13	10.19
tam_Taml	9.90	14.92	18.09	9.48
tur_Latn	11.25	15.92	19.53	10.04
ukr_Cyrl	11.76	16.43	21.87	12.12
vie_Latn	12.00	16.32	21.39	11.49
zho_Hans	10.00	11.51	15.29	6.82

Table 9: BLEU scores on FLORES-200 benchmark. Pivot languages are sorted in alphabetical order and top-4 pivot paths are marked **bold**.



Figure 2: Illustration of the merging process using FiD (Izacard and Grave, 2021).



Figure 3: Illustration of the merging process using TRICE (Huang et al., 2021).

B Impact of Resource-level of Pivot Languages

1055

1056

1059

1062

1063

1064

1065

1067

1068

1070 1071 Under the assumption that high-quality candidates are more adept at conveying the meaning of the source sentence, we select the top-4 paths based on scores on FLORES-200. To verify this hypothesis, we conduct experiments using mid/low-resource pivot languages. According to the WMT22², we select Ukrainian and Croatian as mid- and lowresource languages, respectively. Table 10 shows that using candidates from high-resource languages outperforms those obtained from mid/low-resource languages. The quality of candidates is presented in Table 11. In conclusion, since high-resource languages can also provide sufficient diversity, we select top-performing paths based on the results on FLORES-200.

Mathad	Korean→Italian					
Wittillu	BLEU	chrF++	COMET			
РіvотЕ (GPT-4; <i>U</i> , <i>C</i>)	15.28	41.78	85.75			
PIVOTE (GPT-4; E, S)	16.27	42.55	86.50			

Table 10: Comparison with mid/low-resource languages. *U*, *C*, *E*, and *S* represent candidates from Ukrainian, Croatian, English, and Spanish pivot, respectively.

Model	Korean→Italian						
	BLEU	chrF++	COMET				
Candidate							
NLLB (Ukrainian pivot)	11.95	35.03	82.32				
NLLB (Croatian pivot)	12.25	35.93	79.91				
NLLB (Spanish pivot)	13.87	38.47	83.71				
NLLB (English pivot)	14.77	39.39	81.48				

Table 11: Quality of candidates from each pivot path.

²https://www.statmt.org/wmt22/

translation-task.html

C Metric for Selecting Translation Paths

1072

1073

1074

1075

1077

1079

1081

1082

1083

1084

1085

1086

1087

1089

1090

1091

1092

1093

1094

1097

We conduct a comparative analysis between BLEU and COMET, used for selecting the n translation paths. The results in Table 12 indicate that the difference between metrics is marginal. We believe that this stems from the minimal difference in selected paths, as presented in Table 13. We observe some changes in order and minor differences, but the pivot languages selected by BLEU and COMET show similar compositions.

ton k	Path Salaction	Korean→Italian					
	I alli Selection	BLEU	chrF++	COMET			
	COMET	15.98	42.59	86.22			
top-1	BLEU	16.20	42.84	85.36			
	COMET	16.46	42.58	86.66			
top-2	BLEU	16.57	43.04	86.37			
ta a 2	COMET	16.39	42.41	86.04			
top-5	BLEU	16.66	42.85	86.82			

Table 12: Impact of the pivot path selection metric.

D Implementation Details of the Merging Modules

In this section, we provide detailed explanations of the merging modules. Figure 2 shows the FiD (Izacard and Grave, 2021). First, the instruction and the source sentence are concatenated with each candidate, and processed independently by the encoder. Then the decoder takes the concatenation of each representation and generates the final translation.

As shown in Figure 3, TRICE (Huang et al., 2021) is trained with a two-stage fine-tuning method. In the first fine-tuning stage, the model is trained on two different inputs and single targets: Source \rightarrow Target and Candidate \rightarrow Target. In the second fine-tuning stage, the source and the candidate are concatenated and provided as a single input.

Lang-pair	BLEU	COMET
KO→IT	English (14.82), direct (14.02), Spanish (13.89), Portuguese (13.80)	English (82.89), Spanish (82.70), Indonesian (81.62), Portuguese (81.50)
$IT \rightarrow KO$	English (19.34), direct (18.63), Spanish (18.39), Portuguese (18.01)	Spanish (87.32), English (87.07), Portuguese (87.02), French (86.14)
$AR \rightarrow PT$	English (28.40), direct (27.15), Spanish (26.60), French (24.96)	direct (85.71), English (85.57), Spanish (85.54), Indonesian (84.94)
$PT \rightarrow AR$	English (15.92), direct (15.22), Spanish (14.91), Italian (14.09)	French (82.65), direct (81.36), English (81.04), German (80.44)

Table 13: Selected top-4 pivot paths from each metric. Scores are from experiments on FLORES-200.

1098 E Prompt Templates

1099

1100

1101 1102

1103

1104

1105

1106

1119

1120

1121

1122

We use the zero-shot prompt template from Hendy et al. (2023) to instruct the LLMs for translation,

> Translate this sentence from [source language] to [target language], Source: [source sentence] Target:

when ensembling with candidates, we use the prompt template as follows,

1107 1108 1109 1110	Ensemble the [source language] sentence with the provided [target language] candidates to create the best possible [target language] translation.
1111 1112	[source language] sentence: [source sentence]
1113 1114	<pre>[target language] candidate k: [target candidate]</pre>
1115 1116 1117	Please provide only the [target language] translation and no additional text.
1118	[target language] translation:

F Open-source LLMs

In experiments with LLM-Blender and EVA, we employ the same models as used in each paper. These open-source LLMs are listed in Table 14.

Model	Model Size
LLM-Blender (Jiang et al., 2023)	
Vicuna (Chiang et al., 2023)	13B
Baize (Xu et al., 2023a)	13B
Alpaca (Taori et al., 2023)	13B
Koala (Geng et al., 2023)	13B
Open Assistant (LAION-AI, 2023)	12B
Dolly V2 (Conover et al., 2023)	12B
Flan-T5 (Chung et al., 2022)	11B
MOSS (Sun and Qiu, 2023)	7B
Mosaic MPT (MosaicML, 2023)	7B
StableLM (Stability-AI, 2023)	7B
ChatGLM (Du et al., 2022)	6B
EVA (Xu et al., 2024)	
Baichuan2-Chat (Baichuan, 2023)	7 <u>B</u>
TigerBot-Chat-V3 (Chen et al., 2023)	7B
Vicuna-V1.5 (Chiang et al., 2023)	7B
Llama-2-Chat (Touvron et al., 2023)	7B

Table 14: Open-source LLMs along with their respective model sizes.

G Impact of *temperature* in MBR

To investigate the best performance of MBR, we compared it across three different *temperature* configurations: 1.0, 0.5, and 0.0, which were used in prior works by Farinhas et al. (2023), Suzgun et al. (2023), and Peng et al. (2023), respectively.

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

Table 15 and 16 show the quality of MBR outputs and hypotheses under different *temperature* settings, respectively. Aligning with the findings of the previous study (Peng et al., 2023), we observed that a lower *temperature* setting achieved better performance. Thus, we set the *temperature* of 0.0 for MBR in our experiments.

Mathad	Korean → Italian						
	BLEU	chrF++	COMET				
MBR (temp=1.0; Farinhas et al. (2023))	13.53	42.13	86.57				
MBR (temp=0.5; Suzgun et al. (2023))	13.90	42.19	86.69				
MBR (temp=0.0; Peng et al. (2023))	14.10	42.24	86.70				

Tat	ole	15:	Impact	of	tempe	rature	in	MBR	decoding	3.
-----	-----	-----	--------	----	-------	--------	----	-----	----------	----

Mathad	Korean→Italian					
	BLEU	chrF++	COMET			
MPP hypotheses (temp=1.0)	13.47	41.71	84.94			
MBR hypotheses (<i>lemp</i> =1.0)	(±0.21)	(±0.14)	(±2.62)			
MBP hypotheses $(tamp-0.5)$	13.86	42.03	86.55			
WBR hypotheses (temp=0.5)	(±0.13)	(±0.11)	(±0.15)			
MBR hypotheses $(t_{amn} - 0.0)$	14.09	42.21	86.62			
WBR hypotheses (<i>lemp=</i> 0.0)	(±0.07)	(±0.06)	(±0.10)			

Table 16: Average quality of MBR hypotheses.

H Experiments without Training Data

In Table 17, we report the results of experiments conducted in a setting where no training data was used. When compared to the results in Table 2, this observation affirms that the quality of candidates during ensemble plays an important role in enhancing the final translation.

I Datasets Statistics

Table 18 shows the dataset statistics for each lan-1144guage pair used in the experiments in Table 5.1145





Figure 4: Proportion of pivot languages (Korean \rightarrow Italian) comprising the top-*k* candidates.

Figure 5: Impact of top-*k* values on performance.

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

Model	Korean→Italian		Italian→Korean			Arabic→Portuguese			Portuguese→Arabic			
	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
Standalone NMT System												
NLLB (Costa-jussà et al., 2022)	13.88	38.56	81.54	14.98	20.99	85.75	23.55	47.50	82.08	13.23	38.55	82.79
Vicuna (Chiang et al., 2023)	3.06	15.73	55.31	3.85	7.57	64.24	6.43	17.62	68.10	3.15	19.56	63.79
Baize (Xu et al., 2023a)	4.28	22.85	65.66	3.71	7.36	58.06	7.83	21.78	76.12	2.83	17.91	62.03
Llama-3 (AI@Meta, 2024)	8.46	34.83	82.59	0.61	2.07	29.03	18.37	43.32	81.75	7.31	25.67	75.45
GPT-4 (OpenAI, 2023)	14.07	42.22	86.80	17.23	22.96	86.94	25.82	51.89	85.46	15.11	41.39	83.99
GPT-40 (OpenAI, 2024)	15.11	42.59	85.93	17.20	22.82	85.31	27.28	52.57	85.90	16.28	42.40	83.82
Prior Ensemble Method												
LLM-Blender (Jiang et al., 2023)	3.75	19.58	61.26	0.17	1.58	39.04	6.92	22.14	64.01	1.97	3.42	44.62
EVA (Xu et al., 2024)	3.65	23.54	63.94	1.14	3.23	46.99	8.23	26.87	56.68	3.46	21.88	60.30
MBR (Farinhas et al., 2023)	14.10	42.24	86.70	17.14	23.00	87.53	25.45	51.78	85.55	14.66	41.11	83.93
Proposed Method												
PIVOTE (Llama-3; top3)	13.38	39.99	84.21	12.19	19.40	41.52	24.92	49.57	85.17	15.33	40.77	84.26
PIVOTE (Llama-3; D, E)	12.74	38.47	80.04	12.15	18.22	42.51	24.29	48.17	82.75	13.15	38.67	83.08
PIVOTE (GPT-4; top-3)	16.35	42.45	86.32	17.27	23.17	86.06	26.28	51.06	85.22	15.08	41.19	84.37
PIVOTE (GPT-4; D, E)	16.29	42.61	86.75	17.55	23.50	88.02	26.60	51.45	85.68	15.12	41.34	83.87
PIVOTE (GPT-40; top3)	17.39	42.89	85.30	17.15	23.27	87.68	27.88	51.98	85.58	15.91	42.10	84.58
PIVOTE (GPT-40; D, E)	17.22	42.78	85.29	16.99	23.04	87.96	27.38	51.68	85.64	15.60	41.64	84.80

Table 17: Results on Korean⇔Italian and Arabic⇔Portuguese, in a setting where no training data was used.

Long noin	Detect	# Sentences					
Lang-pair	Dataset	Train	Dev	Test			
	Distant Language Pa	irs					
	news-commentary v18.1	66 742	2 000	2 000			
$P1 \leftrightarrow RU$	(Tiedemann, 2012)	00,743	2,000	2,000			
	news-commentary v18.1	80 724	2 000	2 000			
$NL \leftrightarrow KU$	(Tiedemann, 2012)	80,724	2,000	2,000			
	WikiMatrix v1	166 063	2 000	2 000			
	(Schwenk et al., 2019)	100,005	2,000	2,000			
	Similar Language Pa	irs					
ES () PT	TED 2020 v1	215 462	2 000	2 000			
$E3 \leftrightarrow P1$	(Reimers and Gurevych, 2020)	515,402	2,000	2,000			
	TED 2020 v1	107 078	2 000	2 000			
$\mathbf{UK} \leftrightarrow \mathbf{KU}$	(Reimers and Gurevych, 2020)	197,978	2,000	2,000			

Table 18: Number of sentences in the corpus and data split for each language pair.

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

J Additional Analysis on Candidates

We conduct experiments to investigate the impact of the value of k in the top-k candidates and its composition. Figure 4 illustrates the proportion of pivot languages composing the top-k candidates. Top-k candidates, selected by the QE metric, are composed of diverse candidates obtained through various pivot languages. We also observe the same tendency in other datasets. This suggests that generating diverse candidates through multiple paths helps acquire higher-quality candidates. Figure 5 presents BLEU scores for different values of the k. The highest BLEU is achieved when k is set to 3. These results demonstrate that more candidates in the aggregation process enhance diversity, thereby increasing the likelihood of providing contextually appropriate information. However, it shows convergence around top-3. We attribute this to the inclusion of candidates with lower estimated scores, such as degenerated sentences. Hence, as k increases, the improvement reaches a plateau.

K Results with Additional Models

We report results with diverse GPT models, GPT-3.5 and GPT-40-mini, in Table 19. The version of gpt-3.5-turbo-1106 and gpt-40-mini-2024-07-18 are employed for GPT-3.5 and GPT-40-mini, respectively.

L Impact of Ranking Strategies for Candidate Selection

In this experiment, we compare the case of using1175the PAIRRANKER (Jiang et al., 2023) and the case1176of using COMETkiwi (Rei et al., 2022b) for a1177ranking stage. Table 20 compares the results af-1178ter selecting the top-3 using PAIRRANKER and1179COMETkiwi. As shown in the results, the differ-1180

Model	K	orean→It	alian	Ita	Italian→Korean		Arabic → Portuguese			Portuguese→Arabic		
Model	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
Standalone NMT System												
GPT-3.5 (Ouyang et al., 2022)	12.77	40.13	82.58	15.28	21.13	85.91	25.40	50.23	85.06	14.73	40.81	84.37
GPT-4 (OpenAI, 2023)	14.07	42.22	86.80	17.23	22.96	86.94	25.82	51.89	85.46	15.11	41.39	83.99
GPT-4o-mini (OpenAI, 2024)	13.39	41.35	85.45	17.15	22.80	85.22	23.98	50.25	84.47	15.00	40.72	84.44
GPT-40 (OpenAI, 2024)	15.11	42.59	85.93	17.20	22.82	85.31	27.28	52.57	85.90	16.28	42.40	83.82
Proposed Method												
PIVOTE (GPT-3.5; top3)	15.07	39.87	83.13	15.08	21.12	86.03	26.83	50.23	85.46	15.65	41.10	84.13
PIVOTE (GPT-3.5; D, E)	16.44	41.53	85.23	16.60	22.64	87.84	27.49	50.63	84.35	16.12	41.59	83.82
PIVOTE (GPT-4; top3)	16.66	42.85	86.82	17.95	23.84	87.50	27.22	51.73	85.65	16.53	42.41	84.46
PIVOTE (GPT-4; D, E)	17.10	43.29	85.92	18.18	24.05	88.74	27.98	52.41	85.27	17.02	43.02	84.82
PIVOTE (GPT-4o-mini; top3)	16.25	41.43	83.87	17.06	23.10	88.09	27.10	50.80	85.26	16.16	41.79	84.27
PIVOTE (GPT-40-mini; D, E)	16.58	41.84	85.16	17.44	23.51	87.12	28.19	51.44	84.65	17.21	42.63	83.11
PIVOTE (GPT-40; top3)	17.77	43.38	85.46	18.08	23.98	88.15	28.62	52.53	85.87	16.92	42.93	84.52
PIVOTE (GPT-40; D, E)	18.02	43.46	86.19	18.31	24.32	88.33	29.50	53.16	86.03	17.66	43.73	84.27

Table 19: Results on Korean↔Italian and Arabic↔Portuguese, with diverse models.

Mathad	Korean→Italian						
Wiethou	BLEU	chrF++	COMET				
PAIRRANKER (Jiang et al., 2023)	16.74	42.82	85.92				
COMETkiwi (Rei et al., 2022b)	16.66	42.85	86.82				

Table 20: Impact of candidate ranking strategies.

ence in the final ensemble scores using the two 1181 ranking methods is not significant. We believe this 1182 is because the candidates selected by both ranking 1183 methods are similar. There are 979 out of 2000 test 1184 sets (48.95%) where the top-3 candidates selected 1185 by both ranking methods are the same. In cases 1186 with 2 out of 3 matches, there were 1533 instances 1187 (76.65%). Given the similarity in predictions by 1188 both ranking methods, the final scores exhibit com-1189 parable performance, except in the case of COMET. 1190 From the cost perspective, PAIRRANKER requires 1191 comparisons for $O(N^2)$ unique pair combinations 1192 depending on the number of candidates N. How-1193 ever, COMETkiwi only needs to sort the scores 1194 of N candidates, resulting in a time complexity of 1195 $O(N \log N)$. Therefore, due to its computational 1196 efficiency, we use COMETkiwi to rank candidates. 1197