

Beyond English: The Impact of Prompt Translation Strategies across Languages and Tasks in Multilingual LLMs

Itai Mondshine^a, Tzuf Paz-Argaman^a, and Reut Tsarfaty^a

^aBar-Ilan University, Israel,

{mondshil, tzuf.paz-argaman, reut.tsarfaty}@biu.ac.il

Abstract

Despite advances in the multilingual capabilities of Large Language Models (LLMs) across diverse tasks, English remains the dominant language for LLM research and development. This has led to the widespread practice of *pre-translation*, i.e., translating non-English task prompts into English before inference. *Selective pre-translation*, a more surgical approach, focuses on translating specific prompt components. However, its current use is sporadic and lacks a systematic research foundation. Consequently, the optimal *selective pre-translation* strategy for various multilingual settings and tasks remains unclear. In this work, we aim to uncover the optimal setup for *selective pre-translation* by systematically assessing its use. Specifically, we view the prompt as a modular entity, composed of four functional parts: instruction, context, examples, and output, either of which could be translated or not. We evaluate pre-translation strategies across 35 languages covering both low and high-resource languages, on various tasks including Question Answering (QA), Natural Language Inference (NLI), Named Entity Recognition (NER), and Abstractive Summarization. Our experiments show the impact of factors as similarity to English, translation quality, and the size of pre-trained data, on the model performance. We suggest practical guidelines for choosing optimal strategies in various multilingual settings.¹

1 Introduction

Large language models (LLMs) demonstrate impressive capabilities across various natural language processing tasks, including machine translation (Kocmi et al., 2023), natural language understanding (Saba, 2024) and complex reasoning

¹We launched a user-friendly HuggingFace Space for generation and use of selective pre-translation prompts https://huggingface.co/spaces/naacl-anonymous/selective_pre_translation. Appendix D provides further details and illustrations.

tasks (Huang and Chang, 2022). These exceptional capabilities of LLMs stem, to a large extent, from the vast amounts of data they were trained on (Kaplan et al., 2020). Current LLMs are primarily trained on English data but also include data from other languages, i.e., GPT-3 was trained on 119 languages, but only 7% of the tokens are from non-English languages.² With over 7,000 languages spoken worldwide (Anderson, 2010), the increasing pace of globalization has amplified the need for LLMs that understand and respond in diverse languages.

One common strategy to respond to a task presented in a language different than English is *pre-translation*, which involves translating the complete prompt into English before querying the model (Ahuja et al., 2023; Shi et al., 2022), allowing to leverage the robust capabilities acquired in English across different languages. At the same time, this approach introduces complexities and risks of information loss (Nicholas and Bhatia, 2023). Also, it is unclear whether this approach is uniformly effective across languages and tasks, especially tasks requiring region-specific or culturally-appropriate knowledge.

In contrast to pre-translation, recent studies show that *direct inference*, i.e., prompting the model directly in the (non-English) *source* language spoken by the user, outperforms pre-translation for tasks like QA (Intrator et al., 2024). However, it is unclear whether this approach is optimal, considering that the model was trained on limited data in the source language. It is also unclear how much information is shared across languages during pre-training. Be that as it may, as we show in Sec. 3.2.1, direct inference results still remain suboptimal.

In view of these shortcomings, various studies propose to use *selective pre-translation*, a more nuanced method compared to the de-facto standard

²https://github.com/openai/gpt3/blob/master/dataset_statistics

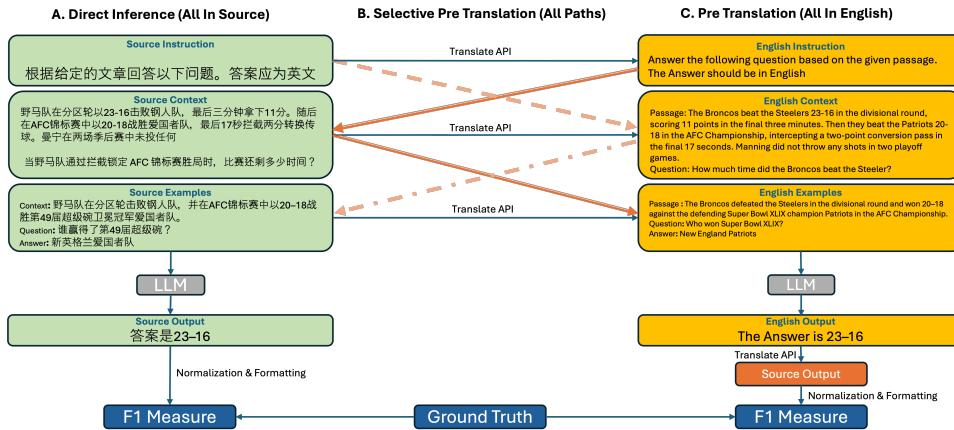


Figure 1: Prompting Strategies: *Direct Inference*, *Selective Pre-Translation*, and *Pre-Translation*

pre-translation approach, which calls for translating only specific parts of the prompt (Ahuja et al., 2023; Kim et al., 2023; Kim et al.). For example, Liu et al. (2024) show that translating only the context to English outperforms direct-inference in summarization and NLI. Ahuja et al. (2023) translated few-shot examples to English while keeping the context in the source language. Kim et al. (2023) used the selective approach when prompting different cross-lingual compositions of in-context examples. However, the selective approach lacks a systematic evaluation of more complex setups, e.g., instruction in English and output in the source language. Consequently, the efficacy of *selective pre-translation* and the optimal prompt configurations for various multilingual settings and tasks remain unclear. To fill this gap, in this paper we set out to examine the impact of selective pre-translation, a commonly used method, across diverse tasks, in order to devise effective prompting strategies for multilingual LLMs.

Concretely, we define a formal *configuration* for a prompt — consisting of four functional parts: instruction, context, examples, and output — either of which could be *selectively* pre-translated or not (see also Winata et al. (2021); Ahuja et al. (2023)). We exhaustively assess all configurations of cross-lingual prompt translation into English from different source languages. Figure 1 presents an overview of our approach, demonstrating the various *selective pre-translation* strategies compared against *pre-translation* and *direct inference* in the source language.

Through a comprehensive evaluation involving 35 languages, four tasks, six dataset collections, and three models, our results demonstrate that *se-*

lective pre-translation consistently outperforms both *pre-translation* and *direct inference* in the source language, establishing the efficacy of *selective pre-translation* strategies (Section 3). Additionally, we analyze the considerations in determining which component to translate, and illustrate the optimal strategies across tasks and languages with varying resource levels. Moreover, we examine how factors such as language similarity to English, training size, and language script affect task performance, and show the effectiveness of selective pre-translation method in mitigating various translation issues, by choosing which prompt components to translate (Section 4).

More specifically, our findings demonstrate that in extractive tasks such as QA or NER, where the output overlaps with the provided context and no generation is needed, the model is either agnostic to the context language in the case of high-resource languages or prefers context in the source language in the case of low-resource languages. Surprisingly, we have discovered that low-resource languages yield better results even when the model’s output is required in English, e.g., in NER (Section 3). Moreover, we show that translation quality significantly affects model performance and that the *selective pre-translation* approach essentially mitigates the negative effect of suboptimal translations, which are in turn specifically problematic in lower-resourced languages (Section 4).

All in all, our extensive and systematic evaluation of pre-translation strategies facilitates generalization across a broader range of languages and tasks, beyond the specific ones herein, towards more robust LLM-use in multilingual settings.

2 The Proposal: Formalizing Prompts Selective Pre-Translation Strategies

Current practices of *prompting* generative LLMs, such as the GPT models family (Ouyang et al., 2022) and Gemini (Team et al., 2023), uncover two remarkable capabilities in performing language processing tasks: (i) *chain of thought* (Wei et al., 2022), where LLMs solve complex tasks through a series of intermediate reasoning steps, and (ii) *in-context learning* (Brown et al., 2020), allowing the model to adapt to new tasks based on limited examples, without weights updates. These capabilities are built on top of the notion of a *prompt*, which serves as a prefix for the LLM’s response. These capabilities are powered by the complex nuanced structure of nowadays prompts, consisting of four components: *instruction*, *context*, *examples*, and *output*.

Let us first define these four components, as follows. The **Instruction** (I) provides a natural language guidance to the model, explaining the task to be performed. The **Context** (X) represents the task data that the model operates on in performing the task. **Examples** (E) are optional illustrations of context:output pairs, that can be used for in-context learning. Overall, we define $\langle I X \rangle$ as a zero-shot prompt and a few-shot prompt as $\langle I X E \rangle$. The prompt is processed by a model M to yield an **Output** (O), where the instruction can include a request for the model to generate the output in a specific language or format.

Each component, i.e., the instruction, the context, the examples, or the output, may be pre-translated or not. We denote a pre-translation decision $l \in \{e, s\}$, where e stands for English pre-translation and s for the source language. Standardly the prompt is composed as $\langle I X E \rangle$ and is delivered to the model M , which in turn emits an output O . We define a specific *pre-translation configuration*³ as $c = \langle I^l, X^{l_x}, E^{l_e}, O^{l_o} \rangle$ where the subscript $l \in \{e, s\}$ indicates the language of the component. Having defined the pre-translation configurations, we evaluate them in different settings.

3 Selective Pre-Translation Evaluation

3.1 Experimental Setup

Goal We set out to compare *selective pre-translation* to both *pre-translation* and *direct infer-*

³See Appendix A.2 for specific configuration examples.

Affinity	Class	Range (% of tokens)	Avg. #tokens (M)	STD
High Resource	A	$p \geq 0.1\%$	1,240	1,156
Medium Resource	B	$0.01\% < p < 0.1\%$	72	49
Low Resource	C	$0\% < p \leq 0.01\%$	5.07	5.41
Unrepresented	D	$p = 0\%$	0	0

Table 1: Language categorization based on the percentage (p) of tokens per language in GPT-3’s training data. Avg. token count (millions), STD: standard deviation.

ence, and to assess the impact of the selected configuration on task performance across languages.

Prompt Configuration We assess *selective pre-translation* in both zero-shot and few-shot settings. In the zero-shot settings, with no examples, we considered 2^3 configurations. For the few-shot scenario, with four components, each is either translated to English or retains the source language, we get 2^4 configurations. All in all we experiments with 24 configurations per language and task.⁴

Prompt Creation And Output Normalization

Based on the prompt configuration, we used the Google Translate API⁵ to translate the components that required translation. After querying the model, we normalized and formatted its output, then translated it to match the language of the gold standard. See Appendix A.1 for further implementation.

Models We conducted experiments on several LLMs: (1) Standard generative models—GPT-3.5-turbo (Ouyang et al., 2022), Mistral-8x7B (Jiang et al., 2023), and Gemini-1.0-pro (Team et al., 2023) — with context sizes of 16k, 32k, and 8k, respectively; (2) multilingual - bloomz-7b1-mt (Muennighoff et al., 2022), with a 2k context.⁶

Language Selection and Categorization We selected ~11 languages per task, ensuring a balanced representation across resource levels (high, medium, low). Due to the lack of precise pre-training distribution for the LLMs we use, we employed the GPT-3 distribution as a proxy, as it is the only distribution publicly shared, to our knowledge.⁷ The GPT-3’s multilingual coverage enables us to categorize languages into classes based on their data ratios. Following Lai et al. (2023), we categorized the tested languages into four classes based on data ratio: High-Resource (A), Medium-

⁴NLI has 12 configurations, with output always in the instruction language, due to its particular, fixed, output format.

⁵pypi.org/project/easygoogletranslate/

⁶See Appendix A.1 for details on the models we used.

⁷https://github.com/openai/gpt-3/blob/master/dataset_statistics

Task	Dataset	Languages
NLI	XNLI	<i>High</i> : Spanish, German, Chinese <i>Medium</i> : Greek, Turkish, Arabic <i>Low</i> : Bulgarian, Hindi, Thai, Swahili, Urdu
QA	XQuAD	<i>High</i> : German, Russian, Romanian <i>Medium</i> : Arabic, Greek, Vietnamese
	IndicQA	<i>Low</i> : Hindi, Malayalam, Bengali, Telugu <i>Unrepresented</i> : Assamese
NER	MasakhaNER	<i>Unrepresented</i> : Bambara, Ewe, Hausa, Yoruba
	WikiANN	<i>High</i> : French, Chinese, Italian, Portuguese, Swedish <i>Medium</i> : Serbian, Slovak
Summarization	XL-Sum	<i>High</i> : French, Japanese, Spanish, Portuguese <i>Medium</i> : Korean, Turkish, <i>Low</i> : Azerbaijani, Nepali, Persian, Uzbek

Table 2: Experiment Setup: Tasks, datasets, languages. Languages are separated by their resource-type affinity.

Resource (B), Low-Resource (C), and Unrepresented (D).⁸ Class D, which includes languages unseen during training. Table 1 summarizes this classification with basic properties.⁹

Tasks and Datasets We assess model performance on 4 tasks, NLI, QA, NER, and Summarization, which we detail in turn. (i) *Natural Language Inference (NLI)* determines if a hypothesis entails, contradicts, or is neutral to a premise. We use the XNLI dataset (Conneau et al., 2018), with sentence pairs in 11 languages, and measured prediction accuracy. (ii) *Question Answering (QA)*: We focus on extractive QA, where the model identifies the answer span in a given context. We evaluated performance on XQuAD (Artetxe et al., 2019) and IndicQA (Doddapaneni et al., 2022) for Indic languages, using the F1 score to assess performance. (iii) *Named Entity Recognition (NER)*: We sampled languages from two datasets: WikiANN (Pan et al., 2017), which includes Wikipedia sentences annotated with Location, Person, and Organization tags in 176 languages; and MasakhaNER (Adelani et al., 2021), for African languages. While both datasets use the BIOES scheme to delineate entity boundaries, we recast the task as generative, prompting the model to generate the list labeled named entities for a given input context. Model performance has been evaluated using F1 scores. (iv) *Abstractive Text Summarization* involves generating short summaries of long contexts, rather than extracting existing sentences. We used the XL-Sum dataset (Hasan et al., 2021), which offers news article summaries in diverse languages. We sampled 10 languages from the dataset and evaluated with ROUGE. We conducted experiments

⁸See Table 12 for list of languages, codes and data ratios.

⁹Alternative criteria such as speakers ratio, as proposed by Joshi et al. (2020), do not reflect language diversity in LLMs, which is affected by availability of data rather than speakers.

on a sample of 250 examples¹⁰ from the test sets for each language.¹¹In total, the datasets we use encompass 35 languages across 4 tasks. Table 2 lists the datasets used, covering ~11 languages per task, ensuring a balanced representation of {Low, Mid, High} resource categories for each task.

Analysis Methods To analyze the empirical results and detect the most influential components, we use three methods: (i) *Correlation analysis* – Assessing the relationship between the model’s prediction scores and the language selection per component. (ii) *Association Rule Learning (ARL)* and *the Apriori algorithm* – While correlation analysis provides a preliminary understanding of the relationship between individual components and model performance, it does not capture non-linear relationships, i.e., the combined effect of multiple translation decisions on performance. To address this limitation, we utilize ARL with the Apriori algorithm (Piatetsky-Shapiro, 1991).¹² (iii) *Performance Gap* – We computed the average difference of k configuration pairs c_i and c_j such that they differ only in the language of one component, e.g., $\langle I^e, X^e, E^e, O^e \rangle$ and $\langle I^s, X^e, E^e, O^e \rangle$. We then calculated this average to determine the performance gap for each specific task: $\frac{1}{k} \sum_{(i,j)=1}^k (\text{Eval}(c_i) - \text{Eval}(c_j))$, where $\text{Eval}()$ denotes the task evaluation score, and k is the number of distinct pairs.

3.2 Results

In Section 3.2.1, we present the results of *selective pre-translation* demonstrating their advantage over both *direct inference* (source language only) and *pre-translation* (English only). Subsequently, in Section 3.2.2, we identify the optimal configurations for each task and analyze the impact of each component on the overall performance, emphasizing key considerations for effective prompting. We start off with GPT-3.5-Turbo and proceed to verify that our results generalize to other models.

3.2.1 Selective Pre-Translation Advantage

Table 3 shows each language’s highest-performing configuration score among all 24 distinct configurations.¹³ Additionally, we display the improvement

¹⁰We selected 250 that followed the can fit into the context of the model, i.e., $< 16K$.

¹¹For tasks without public test sets (XQuAD, IndicQA), we used the validation data.

¹²Appendix A.3.1 details the algorithm and implementation.

¹³Appendix E displays the full-fledged table of results.

Question Answering (QA, F1)					Summarization (ROUGE)				Named Entity Recognition (NER, F1)				Natural Language Inference (NLI, Acc.)						
Lng	Cls.	↑Top	Src. (%)	Eng. (%)	Lng	Cls.	↑Top	Src. (%)	Eng. (%)	Lng	Cls.	↑Top	Src. (%)	Eng. (%)	Lng	Cls.	↑Top	Src. (%)	Eng. (%)
en	A	0.77	N.A	N.A	en	A	30.23	N.A	N.A	en	A	0.65	N.A	N.A	en	A	0.69	N.A	N.A
de	A	0.85	18%	9%	fr	A	35.12	16%	10%	sr	B	0.77	52%	265%	sw	C	0.73	58%	28%
hi	C	0.82	32%	182%	ja	A	32.47	17%	14%	it	A	0.75	9%	41%	bg	C	0.72	57%	8%
ar	B	0.74	84%	138%	fa	C	29.34	21%	0%	sk	B	0.72	15%	36%	el	B	0.71	24%	30%
vi	B	0.73	0%	58%	es	A	28.28	10%	3%	po	A	0.72	18%	20%	es	A	0.69	20%	18%
ro	A	0.69	0%	9%	po	A	27.40	8%	0%	fr	A	0.72	23%	24%	ar	B	0.67	28%	23%
ru	A	0.69	6%	305%	tr	B	20.87	18%	0%	hau	C	0.70	62%	51%	hi	C	0.64	59%	8%
el	B	0.69	0%	2%	ne	C	19.58	31%	28%	ee	D	0.68	46%	81%	de	A	0.64	19%	8%
bn	D	0.68	44%	423%	as	D	15.79	17%	7%	sv	A	0.68	12%	9%	zh	B	0.63	16%	4%
as	D	0.56	138%	450%	uz	C	15.72	58%	24%	zh	B	0.63	90%	121%	th	B	0.57	49%	10%
te	C	0.53	210.10%	253.30%	ko	B	11.84	36.99%	11.78%	bam	D	0.33	33.25%	80.24%	ur	C	0.57	29.96%	9.08%
ml	C	0.49	104.30%	600.00%						yor	D	0.32	66.02%	49.19%	tr	B	0.57	0.00%	8.14%

Table 3: For each Language, we present the top-performing selective configuration score over all other configurations (*Top*) along with its relative improvement (%) over *direct inference* (*Src*) and *pre-translation* (*Eng*).

QA					Summarization					NER					NLI							
lng	cls.	instruction	context	examples	output	lng	cls.	instruction	context	examples	output	lng	cls.	instruction	context	examples	output	lng	cls.	instruction	context	examples
ru	A	-0.08**	0.35**	0.12**	0.09**	ja	A	-0.33**	-0.08	-0.02*	0.00	fr	A	-0.11*	0.10*	-0.01	0.01	de	A	-0.03	-0.02	-0.01
de	A	-0.03**	0.30**	0.08	0.03*	fr	A	0.01	0.020	-0.04	0.06	it	A	0.02	-0.04	-0.04	0.01	es	A	-0.03	0.02	-0.03*
ro	A	-0.03	0.12**	0.04	0.02	po	A	-0.08*	0.05*	-0.03*	0.10*	po	A	-0.15	0.09*	0.1	0.01	el	B	-0.04	0.01	0.07
vi	B	0.04	0.40**	0.10**	0.10	es	A	-0.09*	0.03*	-0.03	0.05	sv	A	-0.11*	0.06*	-0.03**	0.01	zh	B	0.01	-0.06	-0.06
ar	B	-0.07**	0.20**	0.13**	0.04*	tr	B	-0.14**	0.10	-0.1	-0.03*	zh	B	-0.26**	0.44**	0.00	0.07	ar	B	0.00	-0.02	-0.04
el	B	-0.06	0.48**	0.03	0.07*	ko	B	-0.10**	0.13	0.01	0.05	sr	B	-0.26**	0.44**	0.09**	0.05	th	B	-0.03	0.03	-0.14*
bn	C	-0.10**	0.38**	0.03	0.03	uz	C	-0.42**	0.14	0.03	-0.12*	sk	B	-0.11**	0.30**	-0.1*	0.01	tr	B	-0.02	0.00	0.02
ma	C	-0.14**	0.30**	0.01	0.03	fa	C	-0.37**	0.05	-0.07**	-0.04	bam	D	0.03	0.44**	-0.11*	0.02	ur	C	0.01	0.01	-0.08*
te	C	-0.10**	0.38**	0.03	0.03	ne	C	-0.35**	-0.09	0.07**	-0.14	ewe	D	-0.01	0.38**	-0.12**	0.01	bg	C	0.01	0.05	-0.13*
hi	C	-0.07**	0.30**	0.05	0.01	az	C	-0.30**	0.04	-0.00	-0.05	yo	D	-0.01	0.36**	0.01*	0.03	sw	C	0.12	-0.06	-0.09
as	D	-0.04**	0.30**	0.06	0.06							hau	D	-0.04	0.30**	0.08*	0.02	hi	C	-0.03	-0.09	-0.09**

Table 4: For each language, we present the Point-biserial correlation (τ) between the individual component’s language selection (English/Source), and the model performance score across all the configurations samples that use it. Positive $|\tau|$ values correlate with the source language, and negative $|\tau|$ values correlate with English. Significant correlations are indicated by * $p < 0.05$ and ** $p < 0.01$. Bold values denote correlations ($|\tau| > 0.3, p < 0.01$).

Resource	Model	QA				NER				Summarization				NLI			
		I	X	E	O	I	X	E	O	I	X	E	O	I	X	E	O
High	GPT	N	S	S	S	N	S	S	S	S	S	N	N	N	S	S	E
	Gemini	S	S	S	S	N	S	S	S	E	E	Z	N	N	N	N	E
	Mixtral	N	S	S	S	N	S	S	S	S	S	Z	S	N	S	S	E
Low	GPT	N	S	S	S	N	S	S	E	E	S	E	N	E	N	E	S
	Gemini	S	S	S	S	E	S	S	E	S	S	Z	N	N	N	E	E
	Mixtral	N	S	S	S	E	S	S	E	E	E	E	N	S	S	E	E
High	Bloomz	S	S	S	S	S	S	S	S	E	E	E	E	E	E	E	E
Low	Bloomz	S	S	S	S	S	S	S	S	E	E	E	E	E	E	E	E

Table 5: The Top-performing configurations based on Apriori-based Association Rules for High/Low resource level. Comparing Standard LLMs (GPT/Gemini/Mixtral) and Multilingual LLM. *confidence* > 0.8, *support* > 0.15. **S** / **E** - source/English language, **Z** - zero-shot (no examples), **N** - neutral (same performance for English/Source).

(%) over *direct inference* and *pre-translation* for each language. The results indicate that 92% of the tested languages show an improvement over the basic *pre-translation* configuration. Particularly for low-resource languages like Malayalam and Telugu, the gains with *selective pre-translation* are substantial, exceeding 200% in relative improvement. Overall, when comparing selective pre-translation to *complete pre-translation*, the average improvement in low-resource languages is 65% greater than the average improvement in high-resource languages.

The results further reveal that 90% of the languages show improvement over basic *direct inference*. Similar to the pre-translation approach,

low-resource languages like Telugu and Assamese demonstrate a relative improvement of over 100%. High-resource languages also show impressive improvement, albeit smaller, e.g., French and Portuguese show an improvement of over 20% in NER.

Overall, the table shows that *selective pre-translation* can outperform both *pre-translation* and *direct inference*, particularly for languages considered low-resource during pre-training.

3.2.2 The Holy Grail of Optimal Configuration

Having established the advantage of *selective pre-translation* in general, we now study the effects of component language selection on model performance and provide general guidelines for multilingual scenarios.¹⁴ Table 4 shows the Point-biserial correlation between individual component selection and model performance for all the 24 configurations per language/task.¹⁵ Table 5 presents the top-performing configurations, based on the highest-scoring apriori rules for multiple component selections, henceforth *optimal configurations*.

¹⁴For the instruction component language, except for a slight preference for English as demonstrated in Table 4, we did not observe a strong affinity for any language selection.

¹⁵We calculate the correlation between performance and a binary vector indicating whether the component is in English.

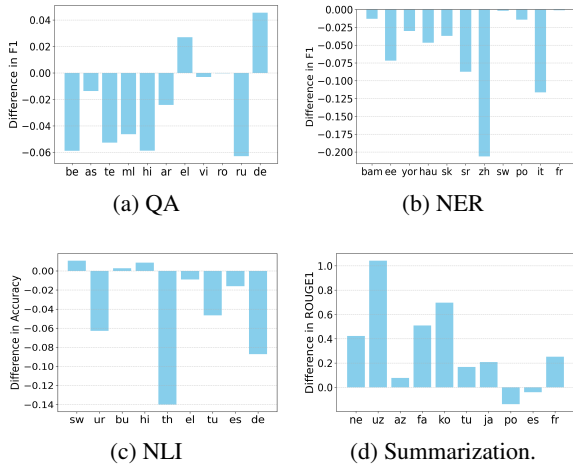


Figure 2: Performance Gap Analysis for the Examples language (English minus Source). Left to right X-axes order indicates Low to High Resource Level. Y-axes indicate language preference: positive values for the source language, and negative for the English language.

Context Language Impact Table 4 shows that source language selection correlates most with model performance score in extractive tasks, such as QA (average of 0.33) and NER (average of 0.32), particularly for low-resource languages (Class C/D), which demonstrate a 70% higher correlation coefficient compared to high-resource languages. In contrast, in tasks like abstractive summarization and NLI, we found no correlation (average of 0.05) to context language selection. Our rule-association analysis in Table 5 further underscores the importance of source-language context in extractive tasks, especially with low-resource languages.

Examples Impact In general, the top-performing configurations of the GPT model in Table 5 show that the optimal configurations are those that include examples, i.e., a few-shot rather than zero-shot setup (Appendix A.4.1 further supports incorporating examples in prompts, especially for high-resource languages). Concretely concerning the language selected for the examples, the optimal configurations in Table 5 show that extractive tasks as NER perform better with source-language examples, possibly due to NER’s dependence on region-specific or culturally-relevant knowledge. Also, the performance gap analysis in Figure 2, shows that, for extractive tasks, prompts with examples in the source language perform better than those with English examples, especially for low-resource languages (See 2(a)/2(b)).

Output Language Impact Unlike context and examples, the output depends on the model generations’s grammaticality and fluency. The best-performing configurations for the GPT model in Table 5 indicate that for extractive tasks, source-language output is beneficial across all languages. Interestingly, despite context mismatches, NER in low-resource languages also benefits from English output. For generative tasks such as summarization, model output in English performs better due to the model’s stronger capabilities in English, *even though* we back-translate the output to source prior to evaluation. Thus, while it is fine in such generative tasks to instruct the model to generate outputs in the source language for high-resource languages, it appears better to generate in English in the low-resource case.

3.3 Beyond Configuration: Key Factors

Having analyzed the impact of the components’ language selection, we discuss key additional factors influencing the efficacy of our approach.¹⁶

Pre-Training Data Size Impact Table 3 presents the optimal prompt configuration scores per language and task. For QA, summarization, and NER, the general trend indicates that even for the optimal pre-translation configuration, classes A and B (High-Medium resource) achieve better results than classes C and D (Low resource), However, a few exceptions exist, i.e., in Hausa and Ewe (Class C/D) we see better results on the NER task compared to Swedish and Chinese (Class A/B). For the NLI task we found no trend where a class C/D languages outperform A/B languages. So, while pre-trained data distribution matters, selective pre-translation can help low-resource languages match the results of higher-resource ones in specific tasks.

Linguistic Similarity to English Impact We used pre-computed syntactic similarities to English from the URIEL dataset (Littell et al., 2017) and calculated the Pearson correlation for each task between the best-performing configuration scores (top score per language) and the syntactic similarity of these languages to English. A moderate correlation (0.42) for summarization shows that syntactic similarity to English positively correlates with performance. NER also shows moderate correlations, suggesting models better identify entities

¹⁶See also Appendix A.4.3 for script impact.

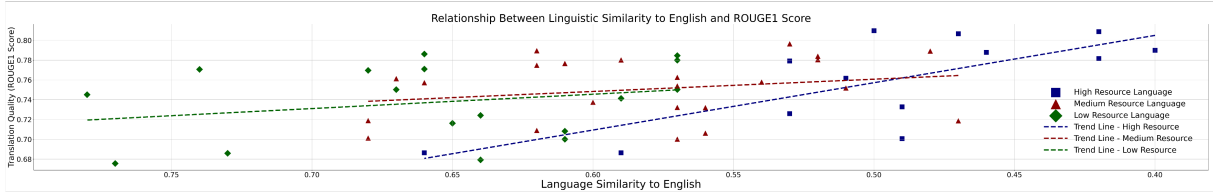


Figure 3: Scatter plot showing the relationship between syntactic similarity to English (further right is more similar) and translation quality (ROUGE) for four language resource subsets (represented as distinct four colored shapes). Each dot represents a different language. Positive linear regression shows an upward trend.

when texts share syntactic features with English. Appendix A.4.3 further details these results.

Standard Model Impact In the previous section, we assessed the selective pre-translation strategies using the GPT model. In this section we check whether these strategies generalize to other LLMs. Table 5 displays the optimal configurations per task and language for Gemini and Mixtral. We see that the *preference for source language in extractive tasks* (for context, examples, and output) holds across all three models. Additionally, outputting in English while keeping the *context in the source language for NER* in low-resource languages is consistent. In NLI, models are agnostic to instruction language. However, surprisingly, in abstractive summarization, we found no clear pattern.¹⁷

Multilingual Model Impact In addition to the standard LLMs, we evaluated BLOOMZ-7b1-mt, known for its multilingual capabilities (Muenighoff et al., 2022). Table 5 displays the optimal configurations for BLOOM across all tasks. We found no distinction between resource types for this model. As shown, the preference for source language in QA is relevant here as well. Interestingly, NER can be answered in the source language, highlighting its multilingual strength. However, for generative tasks and NLI, this multilingualism diminishes, as the model tends to favor English prompts. Overall, the top-performing configurations for BLOOMZ indicate that it performs better with single-language prompts rather than with selective pre-translation prompts.

4 Translation: Key to Pre-Translation

In the previous Section we examined various factors affecting model performance for selective pre-translation strategies. Since translation forms the foundation of selective pre-translation, this Section

¹⁷See Appendix E for additional results.

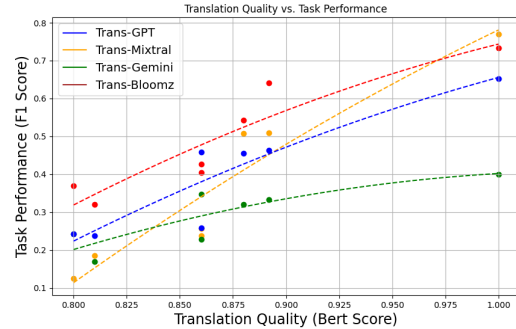


Figure 4: Correlation between translation quality (BERTScore) and accuracy (F1) for Pre-Translation-Zero-shot prompting, each dot is a different language.

focuses on a key question: Are these factors primarily due to the limitations of LLMs, or are due to the quality of the pre-translations themselves? To address this, we first isolate the impact of these factors on translation quality through a controlled experiment. Subsequently, we investigate how translation quality, independently of other factors, influences downstream tasks in our setup.

4.1 Experimental Setup

Factors Affecting Translation Quality We evaluated Google Translate performance on the FLORES-200 validation set (Guzmán et al., 2019), analyzing 91 languages across all resource levels, each with 997 sentences paired with their English translations. We compared machine translations to human-generated references using: (1) n-gram matching metrics – Meteor (Banerjee and Lavie, 2005), ROUGE (Lin, 2004), and BLEU (Papineni et al., 2002), and (2) neural network-based evaluation metrics – BertScore (Zhang et al., 2019) and Comet (Rei et al., 2020). Additional experiments using other translation models are in Appendix C.1.

Impact on Downstream tasks We used the QA dataset XQuAD, with 300 parallel sentences in English and other languages. We compared the

zero-shot *pre-translation* using output in the source language with the *selective pre-translation* optimal approach using GPT-3.5-Turbo (0125).

4.2 Results

Factors Affecting Translation Quality Our analysis focuses on the factors influencing model performance discussed in Sec. 3 and we explore their impact on translation quality. We focus on two factors: *language resource levels (high/low)* and *linguistic similarity to English*. First, we found that the average quality for high-resource languages was higher compared to low-resource languages (0.75 vs 0.73), although correlation between resource level and quality wasn't significant. As for *linguistic similarity to English* we used the pre-computed linguistic similarities from the URIEL dataset (Littell et al., 2017)¹⁸ and calculated the correlation between syntactic similarity to English and the translation quality for each language. Figure 3 shows a positive correlation (coefficient = 0.33, p-value = 0.01) between syntactic similarity to English and translation quality (ROUGE-1). This correlation is particularly strong for high-resource languages (coefficient=0.73, p-value=0.004). Additional correlation results are detailed in Appendix C.2.

Impact On Downstream Tasks To isolate translation quality as the sole direct factor influencing model performance, we translated the entire prompt into English. Since the original input differed only in language, not content, any variation in the processed input can be attributed solely to the quality of the translation. This approach allows us to directly measure the correlation between translation quality and model performance across various tasks. Figure 4 shows the correlation between translation quality (BERTScore) and model performance (accuracy) for each language. Our results show that *higher translation quality goes hand in hand with improved task performance*. The overall Pearson correlation is 0.233 ($p < 0.001$). However, when assessing the same tasks with *selective pre-translation* instead of a *completely pre-translated* prompt, we found a low correlation of 0.05 ($p < 0.05$) between the translation quality and task performance, while selective pre-translation outperforms the fully translated prompt. This disparity shows that the *selective pre-translation* method effectively neutralizes translation issues. By strategi-

cally choosing which prompt components to translate, we can make pre-translation useful for languages with lower translation quality.

In sum, our findings demonstrate that factors influencing downstream tasks, such as high resource level and similarity to English, are positively correlated with translation quality. We further show that *selective pre-translation* can mitigate the negative effects of poor translation quality. These two findings underscore the importance of investing in high-quality translation, and on the other hand, prioritizing the *selective pre-translation* approach in languages where machine translation is sub-optimal.

5 Related Work

With over 7,000 languages spoken globally (Anderson, 2010), the growing use of diverse languages have fueled the demand for multilingual LLMs. Progress in this field stems from two primary efforts: (1) developing dedicated monolingual models for low-to-medium-resource languages (Seker et al., 2022; Cui et al., 2023; Andersland, 2024), and (2) creating multilingual LLMs with pre-trained data encompassing multiple languages (Qin et al., 2024; Jiang et al., 2024).

The ability of Multilingual LLMs to operate in different languages (Raffel et al., 2020; Conneau et al., 2019; Chowdhery et al., 2023) comes from two sources: (1) training or fine-tuning on multilingual data in order to achieve multilingual proficiency (Xue et al., 2020; Chen et al., 2021; Le Scao et al., 2023; Shaham et al., 2024; Muenighoff et al., 2022), and (2) utilizing prompting techniques to harness the model's inherent multilingual capabilities without modifying parameters during inference (Brown et al., 2020). This latter approach has gained popularity due to its efficiency and applicability to a wider range of use cases.

For the latter, to improve the multilingual capabilities of LLMs researchers developed various prompting methods. Huang et al. (2023) introduced XLT, a cross-lingual prompt that directs LLMs to function as experts in a specific language through a process involving problem-solving and cross-lingual thinking. Zhao and Schütze (2021) employed discrete and soft prompting techniques and showed that few-shot non-English prompts outperform finetuning in cross-lingual transfer. Shi et al. (2022) found that chain-of-thought (CoT) prompts lead to multilingual reasoning abilities in LLMs, even in under-represented languages. An-

¹⁸<https://github.com/antonisa/lang2vec>

other strategy is *pre-translation* which translates the entire prompt to English (Chowdhery et al., 2023; Qin et al., 2023; Ahuja et al., 2023). A more nuanced approach, *selective pre-translation*, translates part of the prompt into English, for instance, Liu et al. (2024) translated only the instruction, and Ahuja et al. (2023) translated the few shot examples. While these use cases lack a systematic research foundation, in this study, we systematically study pre-translation configurations to provide evidence-based recommendations for optimal use.

6 Conclusion

In this work we formalize and comprehensively assess selective pre-translation prompting strategies for LLMs in multilingual settings. With four tasks, six dataset collections, three models, and 35 languages, we deliver the first systematic evaluation, to our knowledge, of all existing prompt configurations of pre-translation. We demonstrate that *selective pre-translation* consistently outperforms both *pre-translation* of the entire prompt and *direct-inference* in the source language, establishing the efficacy of *selective pre-translation* in both the high- and low-resource cases. Additionally, we show that translation quality significantly affects performance and that selective pre-translation can mitigate the negative effects of suboptimal translations.

Limitations

Subset of LLMs This study aims to systematically assess the effectiveness of various prompting strategies across different tasks and LLMs. Due to resource limitations, it was not possible to evaluate more advanced models such as GPT-4 or GPT-4o. However, we endeavored to cover several LLMs representing different architectures. Additionally, the choice of Bloom as our multilingual model is based on previous works (Bawden and Yvon, 2023; Nezhad and Agrawal, 2024). We make our evaluation framework, code, configurations, and execution pipeline, for open public use, allowing to extend the investigation to more and newer models.

LLM Adherence and Impact on the Output In our evaluation, we attempted to influence the output by instructing the model to generate a response in a specific language. However, the model occasionally did not follow these instructions, producing output in a different language, which could impact

the results. Appendix B provides error analysis of the various issues we encountered.

Evaluation Metrics based on n-gram matching, such as ROUGE (Lin, 2004), are commonly used for evaluating summarization quality in English. However, these metrics can be problematic when applied to morphologically rich languages (MRL) such as Persian, which have more flexible word order compared to English. Additionally, their morphological richness means that the same concept can be expressed in multiple ways due to variations in prefixes, suffixes, and root conjugations.

Translation Quality’s Impact on Downstream tasks Our analysis of the impact of translation quality impact on downstream tasks in section 4.2 was constrained by the scarcity of datasets with parallel splits for English and other languages, limiting our evaluation to the QA task. Future research should incorporate a wider array of datasets and tasks to validate and expand upon our findings.

Pretrained Data Distribution Details In our experiments, we evaluated four models and grouped the languages based on GPT-3’s pre-training data distribution information. Ideally, we would split the languages according to each model’s data distribution. However, to our knowledge, only GPT-3’s pre-training data distribution is publicly shared. Explicitly testing different language distributions is desired but resource intensive, and is left for future research.

Acknowledgements

This research has been funded by a grant from the Israel Science Foundation (ISF) grant number 670/23 as well as a grant from the Israeli Ministry of Science and Technology (MOST), and a KAMIN grant from the Israeli Innovation Authority, for which we are grateful. We are further grateful for a generous VATAT grant to the BIU NLP team which contributed resources for computation, annotation and human evaluation in this project.

References

David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Constantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, et al. 2021. Masakhaner: Named entity recognition for african languages. *Transactions of the Association for Computational Linguistics*, 9:1116–1131.

- Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. 1993. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM SIGMOD international conference on Management of data*, pages 207–216.
- Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, et al. 2023. Mega: Multilingual evaluation of generative ai. *arXiv preprint arXiv:2303.12528*.
- Michael Andersland. 2024. Amharic llama and llava: Multimodal llms for low resource languages. *arXiv preprint arXiv:2403.06354*.
- Stephen R Anderson. 2010. How many languages are there in the world. *Linguistic Society of America*, pages 1–12.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2019. On the cross-lingual transferability of monolingual representations. *CoRR*, abs/1910.11856.
- 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 machine translation and/or summarization*, pages 65–72.
- Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of bloom. *arXiv preprint arXiv:2303.01911*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Guanhua Chen, Shuming Ma, Yun Chen, Li Dong, Dongdong Zhang, Jia Pan, Wenping Wang, and Furu Wei. 2021. Zero-shot cross-lingual transfer of neural machine translation with multilingual pretrained encoders. *arXiv preprint arXiv:2104.08757*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. [XNLI: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.
- Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.
- Sumanth Doddapaneni, Rahul Aralikatte, Gowtham Ramesh, Shreya Goyal, Mitesh M Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2022. Indicxtreme: A multi-task benchmark for evaluating indic languages. *arXiv preprint arXiv:2212.05409*.
- Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc’Aurelio Ranzato. 2019. [The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English](#). 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 6098–6111, Hong Kong, China. Association for Computational Linguistics.
- Tahmid Hasan, Abhik Bhattacharjee, Md Saiful Islam, Kazi Samin, Yuan-Fang Li, Yong-Bin Kang, M Sohel Rahman, and Rifat Shahriyar. 2021. Xl-sum: Large-scale multilingual abstractive summarization for 44 languages. *arXiv preprint arXiv:2106.13822*.
- Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu Wei. 2023. Not all languages are created equal in llms: Improving multilingual capability by cross-lingual-thought prompting. *arXiv preprint arXiv:2305.07004*.
- Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*.
- Yotam Intrator, Matan Halfon, Roman Goldenberg, Reut Tsarfaty, Matan Eyal, Ehud Rivlin, Yossi Matias, and Natalia Aizenberg. 2024. Breaking the language barrier: Can direct inference outperform pre-translation in multilingual llm applications? *arXiv preprint arXiv:2403.04792*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the nlp world. *arXiv preprint arXiv:2004.09095*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Sunkeyoung Kim, Dayeon Ki, Yireun Kim, and Jinsik Lee. Translating qa is enough: A key to unlocking in-context cross-lingual performance. In *ICML 2024 Workshop on In-Context Learning*.
- Sunkeyoung Kim, Dayeon Ki, Yireun Kim, and Jinsik Lee. 2023. Boosting cross-lingual transferability in multilingual models via in-context learning. *arXiv preprint arXiv:2305.15233*.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023. Findings of the 2023 conference on machine translation (wmt23): Lms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42.
- Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023. Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning. *arXiv preprint arXiv:2304.05613*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2023. Bloom: A 176b-parameter open-access multilingual language model.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Patrick Littell, David R Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin. 2017. Uriel and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 8–14.
- Chaoqun Liu, Wenxuan Zhang, Yiran Zhao, Anh Tuan Luu, and Lidong Bing. 2024. Is translation all you need? a study on solving multilingual tasks with large language models. *arXiv preprint arXiv:2403.10258*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Sina Bagheri Nezhad and Ameeta Agrawal. 2024. What drives performance in multilingual language models? *arXiv preprint arXiv:2404.19159*.
- Gabriel Nicholas and Aliya Bhatia. 2023. Lost in translation: Large language models in non-english content analysis. *arXiv preprint arXiv:2306.07377*.
- 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.
- Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages. In *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 1946–1958.
- 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.
- Gregory Piatetsky-Shapiro. 1991. Discovery, analysis, and presentation of strong rules. *Knowledge Discovery in Data-bases*, pages 229–248.
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*.
- Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S Yu. 2024. Multilingual large language model: A survey of resources, taxonomy and frontiers. *arXiv preprint arXiv:2404.04925*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*.
- Walid S Saba. 2024. Lms’ understanding of natural language revealed. *arXiv preprint arXiv:2407.19630*.
- Amit Seker, Elron Bandel, Dan Bareket, Idan Brusilovsky, Refael Greenfeld, and Reut Tsarfaty. 2022. Alephbert: Language model pre-training and evaluation from sub-word to sentence level. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 46–56.

Uri Shaham, Jonathan Herzig, Roei Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Multilingual instruction tuning with just a pinch of multilinguality. *arXiv preprint arXiv:2401.01854*.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Genta Indra Winata, Andrea Madotto, Zhaojiang Lin, Rosanne Liu, Jason Yosinski, and Pascale Fung. 2021. Language models are few-shot multilingual learners. *arXiv preprint arXiv:2109.07684*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Mengjie Zhao and Hinrich Schütze. 2021. Discrete and soft prompting for multilingual models. *arXiv preprint arXiv:2109.03630*.

A Selective Pre-Translation Evaluation

A.1 Experimental Setup

Models To query GPT-3.5-turbo (0125), we used the Azure platform via the API¹⁹. For Mixtral-8x7B-287 Instruct-v0.1, we utilized the API platform provided by Together.ai²⁰. For Gemini-1.0-pro, we accessed the API through Google AI Studio²¹. Lastly, for bloomz-7b1-mt, we used deployed the model on Hugging Face²²

¹⁹<https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/models>

²⁰<https://www.together.ai/>

²¹<https://aistudio.google.com>

²²<https://huggingface.co/bigscience/bloomz-7b1-mt>

All Configurations

<i>Instruction</i>	<i>Context</i>	<i>Examples</i>	<i>Output</i>
Source	Source	-	English
English	Source	-	English
Source	Source	Source	English
English	English	English	Source
Source	English	Source	English
English	Source	Source	Source
English	Source	-	Source
Source	Source	Source	Source
English	Source	English	Source
Source	Source	-	Source
English	English	Source	Source
English	Source	English	English
English	English	-	English
Source	English	-	English
Source	English	-	Source
Source	Source	English	English
English	English	-	Source
English	Source	Source	English
Source	Source	English	Source
Source	English	Source	Source
English	English	English	English
Source	English	English	English
English	English	Source	English
Source	English	English	Source

Table 6: All Valid Configurations (24)



Figure 5: Examples of 3 configurations of German. Each configuration is in the following format <Instruction, Context, Examples, Output>

Prompt Creation For constructing the prompts we used the LangChain library²³ which enables us to build and validate prompts dynamically for both zero-shot and few-shot templates. For creating the instructions, we initially used ChatGPT to generate them and then fine-tuned them based on quality analysis from our experiments.

Normalization And Formatting Before evaluation, we normalized the model’s output, with each task following a unique normalization process. For the QA task, for example, we converted the text to lowercase and removed punctuation, articles, and extra whitespace. In the Summarization task, we removed prefixes like ‘The Summary:’. For the NER task, we converted the model’s output into a list of tuples, each in the format (Tag, Entity). After normalization, additional formatting was applied if necessary. For instance, in the NER task, we transformed the normalized output into a list in the BIOES format, identifying the entities in the original sentence and converting each entity prediction to its correct format based on its position (e.g., B-ORG for the first entity tagged as ‘ORG’).

A.2 Configuration Format

We define a specific *selective pre-translation configuration* as $C_i = \langle I^l, X^l, E_n^l, O^l \rangle$, $n \geq 0$, $l \in \{e, s\}$. Each configuration contains 4 components: instruction, context, examples, and output. Figure 5 displays examples for 3 configurations in the

²³<https://pypi.org/project/langchain/>

German language. See Table 6 for a list of all the configurations.

Python Libraries In Use For evaluation of the different models, we used the most common ROUGE package for non-English papers²⁴. For loading and processing the data, we used NumPy²⁵. For help with writing the code, we used assistance from ChatGPT.

A.3 Analysis Methods

A.3.1 Rule Association And Apriori Algorithm

Association rule mining, one of the most important and well-researched techniques of data mining, was first introduced by Agrawal et al. (1993). It aims to extract interesting correlations, frequent patterns, associations, or casual structures among sets of items in the transaction databases or other data repositories.

Apriori algorithm The Apriori algorithm is a popular approach for mining association rules. It works by identifying frequent itemsets, which are groups of items that appear together in a dataset with a frequency above a specified threshold. The algorithm then generates association rules from these frequent itemsets, highlighting the likelihood of one item being present given the presence of another item. Apriori uses a bottom-up approach,

²⁴https://github.com/csebuettlp/xl-sum/tree/master/multilingual_rouge_scoring

²⁵<https://pypi.org/project/numpy/>

gradually building larger itemsets from smaller ones while pruning those that do not meet the minimum support threshold.

In our analysis, we reported the following measures: (i) **Support**: $s(X) = \frac{\sigma(X)}{N}$, where $\sigma(X)$ is the number of transactions in which X appears and N is the total number of transactions.

(ii) **Confidence**: $c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$, measures the probability of occurrence of itemset Y with itemset X .

Implementation Details To implement the Rule Association algorithm, we created a DataFrame for each task’s results using pandas DataFrames²⁶. Each DataFrame contains the results for all the configurations for every language. Subsequently, we binned each score column into three bins - high, medium, and low, based on the 30th and 60th percentiles. Later, we merged all the data frames based on the configuration name. Then we used the apriori algorithm from the efficient-apriori²⁷ library, which produces two outputs - itemsets and rules. Later, we filtered weak rules (support > 0.05 & confidence > 0.75).

A.4 Prompting

Question Answering Answer the following <Question> based only on the given <Context>. Follow these instructions:

- Include only words from the given context in your answer.
- Keep the answer as short as possible.
- Provide the answer in *expected output language*.

Named Entity Recognition You are an NLP assistant whose purpose is to perform Named Entity Recognition (NER). You need to assign each entity a tag from the following:

1. PER means a person.
2. ORG means an organization.
3. LOC means a location entity.

The output should be a list of tuples in the format:

[(Tag, Entity), (Tag, Entity)]

for each entity in the sentence. The entities should be in the *expected output language*.

²⁶<https://pypi.org/project/pandas/>

²⁷<https://pypi.org/project/efficient-apriori/>

Summarization Write a summary of the given <Text> The output should be in *expected output language*. The output must be up to 2 sentences maximum.

Natural Language Inference You are an NLP assistant whose purpose is to solve Natural Language Inference (NLI) problems. NLI is the task of determining the inference relation between two texts: entailment, contradiction, or neutral. Your answer should be one word from the following: entailment, contradiction, or neutral.

A.4.1 The Holy Grail of Optimal Configuration

Few-Shot Examples Impact Figure 6 demonstrates that for all tasks, using a few-shot setting over a zero-shot setting yields better results. Interestingly, For all tasks, except for NLI, high-resource languages achieved better improvement when considering a few-shot setting over low-resource languages.

Output Selection Effects Figure 7 demonstrates that while in extractive QA the output should be in the source language, and in the summarization task, the output should be in English; in NER, the output is ambiguous.

A.4.2 Factors Explaining Performance

A.4.3 Script Impact

Figure 8 presents the performance improvement achieved by the highest-performing prompt configuration among all configurations compared to the pre-translation prompt, for each language. Notably, the language family (as categorized by scripts) reveals a relatively even distribution of performance gains within the same language family. For example, languages using the Cyrillic script show greater improvement than those using the Latin script. Interestingly, languages in the same script family sometimes show varying results; for example, Spanish and Ewe belong have Latin script, but Ewe shows greater improvement over Spanish.

Linguistic Similarity To English We used the lang2vec²⁸ library to obtain syntactic similarity scores for each language. The Pearson correlation was calculated based on two vectors: one representing language similarities (ranging from 0 to 1) and the other representing model performance scores for each language across tasks. This correlation

²⁸<https://github.com/antonisa/lang2vec>

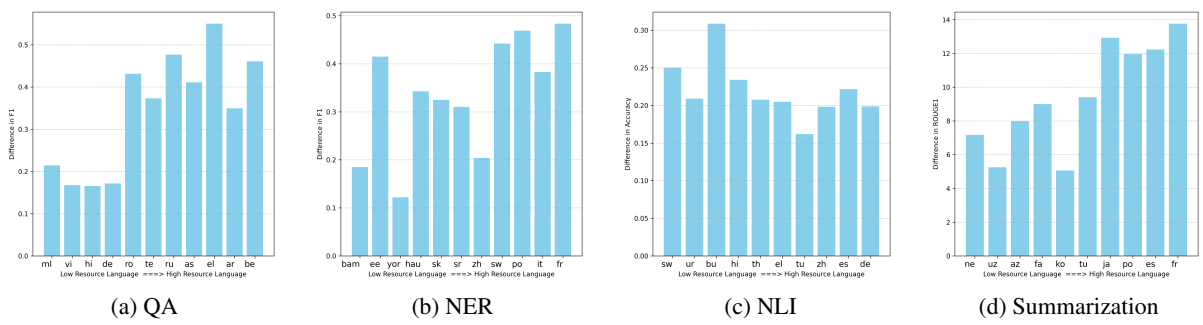


Figure 6: Few-Shot and Zero-Shot Performance Gap (Few-Shot - Zero-Shot) for each task/language.

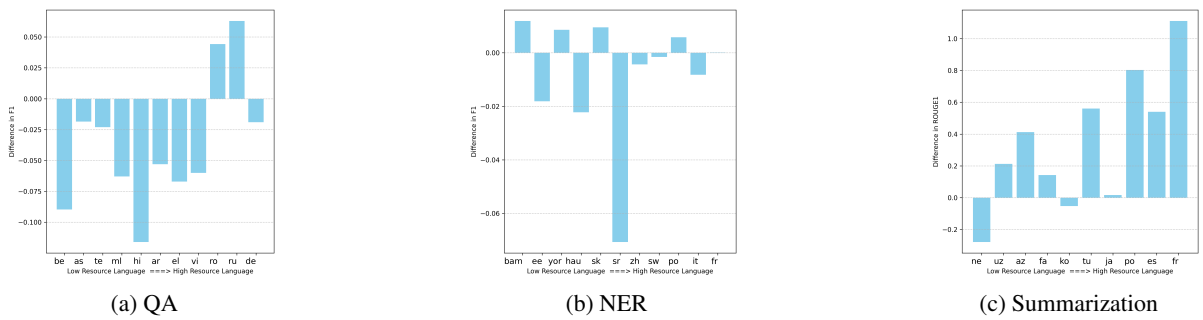


Figure 7: Output Performance Gap (English - Source) for each task/language

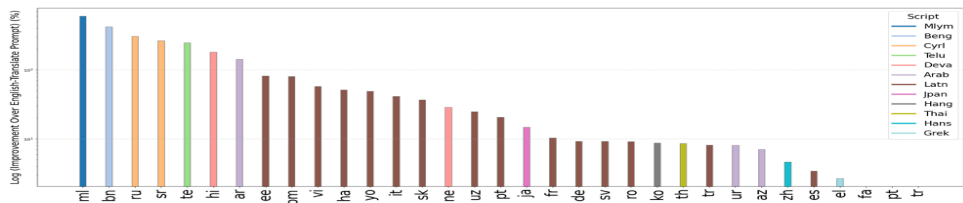


Figure 8: Percentage improvement over *pre-translation* approach, when using the highest configuration for each task For GPT-3.5-Turbo. The bars are color-coded based on the language family script.

Task	Model Output	Expected Output	Explanation	Phenomenon
NER	['LOC: 新北市', 'LOC: 平溪']	[(LOC, '新北市'), (LOC, '平溪')]	List of strings, instead of list of tuples.	Format Inconsistency
	[PER: Hiei]n- [PER: Hinata]n	['PER', 'Hiei'], ('PER', 'Hinata')	New line between each entity.	Format Inconsistency
	Ner Tags: ['PER: LL Cool J']	[(PER: LL Cool J)]	Redundant words in the beginning.	Extraneous information
	[] (No entities found in the sentence)	[]	Redundant words in the end.	Extraneous information
QA	Since the last sentence is in English, I will provide the NER tags in English as well	[(PER: ПАВИЛИВ ГРУЗИНСКИ)]	Refusing to output in the desired language.	Unwarranted Refusal
	[The united states]	The united states	List of string instead of a string.	Format Inconsistency
NLI	The question cannot be answered as the answer is not provided in the given context	[Luke Kuechly]	Insufficient information	Unwarranted Refusal
	The second statement neutral because it does not provide any information that contradicts vinculación	neutral entailment	Unnecessary justification for the choice. Spanish word for entailment instead of English.	Extraneous information Wrong Language
Summarization	Resumo: O ministro de Emergências da Rússia, Sergei Shoigu ...	O ministro de Emergências da Rússia, Sergei Shoigu ...	Redundant words ('Resumo' - Summary in Portuguese) in the beginning.	Extraneous information

Table 7: Error analysis of unexpected model outputs and observed in various tasks/languages.

Model	QA	NER	Summarization	Model	QA	NER	Summarization	NLI
GPT	56	60	96	GPT	0.14*	0.28**	0.42**	0.01
Mixtral	60	61	78	Mixtral	0.13*	0.2**	0.31**	0.08
Gemini	63	61	96	Gemini	0.1*	0.19**	0.25**	0.01

Table 8: Percentage of success of expected output languages for each model/task

Table 9: Pearson correlation between linguistic syntactic similarity to English and task performance for GPT, Mixtral, and Gemini. * $p < 0.05$, ** $p < 0.01$

was calculated at the instance level. Table 9 shows a positive correlation between model performance and syntactic similarity to English, especially for the summarization task, indicating that syntactic similarity to English significantly improves performance in this task. Additionally, NER also exhibits positive correlations, suggesting that models can better identify and classify entities in languages that share syntactic features with English.

B Error Analysis

B.1 Format Issues

Automatic evaluation requires consistent output formatting, especially in tasks like Named Entity Recognition (NER), which must adhere to a pre-defined format rather than free text. A common practice involves prompting the model to generate results in a specific format, such as a list of tuples representing entities and their types (e.g., *(Loc, NewYorkCity)*). However, achieving perfect consistency can be challenging. Models may not always adhere to the requested format, leading to difficulties in evaluation.

Qualitative Analysis We analyzed unexpected model outputs in various tasks and languages. For each task, we noted common phenomena observed and the expected model output. The results in Table 7 reveal that for the NER task, due to its rigid format, the model exhibited many error types. The models showed phenomena such as format inconsistency and extraneous introduction, which require a more generative normalization method to handle. An interesting phenomenon that made our modular selective pre-translating approach difficult to im-

plement is unwarranted refusal, where the model refuses to output in the required language.

B.2 Incorrect Output Language

Table 8 summarizes the percentage of accurately outputted language for all tasks (except NLI, due to its index-based format) across all models. The results reveal that in extractive tasks such as extractive QA and NER, where the output overlaps with the context, the model struggles the most to output in the desired language. However, in abstractive summarization, a generation task, the model had better success.

C Translation: Key to Pre-Translation

C.1 Machine Translation Engines Comparison

To evaluate machine translation tools, we compared Google Translate API and Bing Translator. We excluded multilingual LLMs from consideration, as *zhu2023multilingual* found that these models still lag behind commercial systems like Google Translate, especially for low-resource languages. As shown in Table 11 and Figure 12, Google Translate outperformed Bing Translator across all metrics, demonstrating superior performance. Notably, Welsh and Maltese, both low-resource languages, achieved the highest scores. .

C.2 Linguistic Similarity To English

The results in Table 10 demonstrate the correlation between the syntactic similarity to English of the language and the ROUGE translation score of the language. The results show that the most significant

Group	Pearson Correlation	P-value
High Resource	0.73	0.05
Medium Resource	0.48	0.07
Low Resource	0.06	0.78
Extremely Low Resource	-0.34	0.30

Table 10: Correlation between syntactic similarity to English and the ROUGE score (by language subset).

correlation was observed in languages belonging to the high-resource category, and this correlation decreases as the class of the language becomes low-resource.

D Selective Pre-Translation Prompt Generator

We have launched a space on Hugging Face. The space makes it easy for the community to receive recommended configurations based on the type of task and language. In Figure 9, we can see an overview of the application and an example of a recommended configuration. Figures 10 and 11 provide examples of generating prompts for zero-shot and few-shot settings.

E Detailed Results

The results across all tasks, languages and models are included in our benchmarking exercise are provided in Table 15 (for XQuAD), 13 (for indciQA), 16 (for WikiANN), 17 (for MasakhNER), 18 (for XL-Sum), 14 (for XNLI). The result of the correlation for Gemini are included in Table 19, for Mixtral in Table 20, and for bloomz in Table 21

Language	Google Translate API			Bing Translator		
	ROUGE	Meteor	BLEU	ROUGE	Meteor	BLEU
Welsh	0.86	0.86	0.63	0.85	0.85	0.61
Maltese	0.84	0.83	0.59	0.83	0.82	0.56
Danish	0.81	0.79	0.51	0.81	0.79	0.49
Swedish	0.81	0.80	0.51	0.80	0.79	0.51
Portuguese	0.81	0.79	0.52	0.80	0.78	0.50
Catalan	0.80	0.79	0.49	0.78	0.77	0.45
Spanish	0.79	0.64	0.30	0.69	0.64	0.30
Serbian	0.79	0.77	0.48	0.03	0.07	0.01
Bulgarian	0.79	0.77	0.45	0.75	0.72	0.37
French	0.79	0.77	0.48	0.78	0.76	0.48
Nepali (macrolanguage)	0.79	0.78	0.46	0.74	0.72	0.38
Macedonian	0.78	0.77	0.46	0.72	0.70	0.35
Swahili (macrolanguage)	0.78	0.79	0.51	0.74	0.74	0.43
Hebrew	0.78	0.77	0.47	0.76	0.75	0.44
German	0.78	0.76	0.46	0.79	0.76	0.46
Indonesian	0.78	0.77	0.46	0.78	0.77	0.44
Romanian	0.78	0.76	0.45	0.77	0.75	0.43
Punjabi	0.78	0.77	0.46	0.74	0.72	0.4
Bosnian	0.78	0.76	0.45	0.75	0.73	0.39
Hindi	0.78	0.76	0.45	0.76	0.74	0.41
Turkish	0.77	0.75	0.43	0.76	0.73	0.41
Armenian	0.77	0.75	0.43	0.68	0.65	0.28
Irish	0.77	0.76	0.47	0.76	0.74	0.42
Gujarati	0.77	0.76	0.44	0.73	0.69	0.35
Telugu	0.77	0.76	0.44	0.73	0.71	0.38
Slovak	0.76	0.74	0.42	0.75	0.72	0.40
Italian	0.76	0.68	0.34	0.72	0.68	0.34
Galician	0.76	0.74	0.43	0.74	0.71	0.39
Estonian	0.76	0.74	0.41	0.74	0.71	0.37
Czech	0.76	0.74	0.42	0.76	0.73	0.4
Marathi	0.76	0.74	0.41	0.72	0.69	0.35
Uzbek	0.75	0.74	0.39	0.67	0.63	0.28
Urdu	0.75	0.72	0.39	0.71	0.68	0.33
Ukrainian	0.75	0.73	0.41	0.74	0.72	0.4
Malayalam	0.75	0.73	0.4	0.71	0.68	0.34
Sinhala	0.75	0.73	0.39	0.69	0.66	0.32
Bengali	0.74	0.73	0.39	0.74	0.70	0.36
Croatian	0.74	0.72	0.39	0.73	0.70	0.36
Lao	0.74	0.73	0.39	0.69	0.66	0.30
Haitian	0.74	0.73	0.41	0.67	0.65	0.30
Hungarian	0.74	0.72	0.38	0.74	0.71	0.37
Kazakh	0.73	0.72	0.38	0.67	0.62	0.27
Russian	0.73	0.70	0.38	0.72	0.69	0.36
Vietnamese	0.73	0.72	0.39	0.72	0.71	0.36
Slovenian	0.73	0.71	0.38	0.69	0.67	0.32
Zulu	0.73	0.74	0.43	0.65	0.65	0.32
Tamil	0.73	0.71	0.37	0.70	0.68	0.33
Finnish	0.73	0.70	0.36	0.72	0.68	0.33
Kannada	0.72	0.71	0.37	0.71	0.68	0.33
Lithuanian	0.72	0.70	0.36	0.68	0.63	0.28
Icelandic	0.72	0.70	0.37	0.72	0.7	0.36
Southern Sotho	0.72	0.71	0.40	0.62	0.6	0.27
Korean	0.71	0.68	0.32	0.69	0.66	0.31
Basque	0.71	0.68	0.34	0.67	0.63	0.27
Thai	0.71	0.66	0.3	0.69	0.65	0.28
Chinese	0.70	0.67	0.31	0.68	0.64	0.29
Georgian	0.70	0.66	0.31	0.64	0.58	0.21
Xhosa	0.70	0.70	0.38	0.63	0.63	0.28
Dutch	0.70	0.66	0.31	0.72	0.67	0.34
Japanese	0.69	0.66	0.30	0.69	0.65	0.29
Polish	0.69	0.65	0.30	0.68	0.64	0.29
Burmese	0.69	0.65	0.30	0.63	0.58	0.22
Khmer	0.68	0.65	0.30	0.64	0.59	0.23
Kinyarwanda	0.68	0.67	0.34	0.61	0.6	0.23
Samoaan	0.67	0.65	0.33	0.62	0.59	0.26
Somali	0.66	0.66	0.32	0.59	0.57	0.22
Faroese	0.62	0.60	0.28	0.65	0.62	0.28
Lingala	0.60	0.59	0.24	0.60	0.58	0.23
Azerbaijani	0.31	0.29	0.05	0.11	0.10	0.00
Fijian	0.16	0.16	0.02	0.51	0.47	0.12

Table 11: Comparison between Google Translate API and Bing Translator.

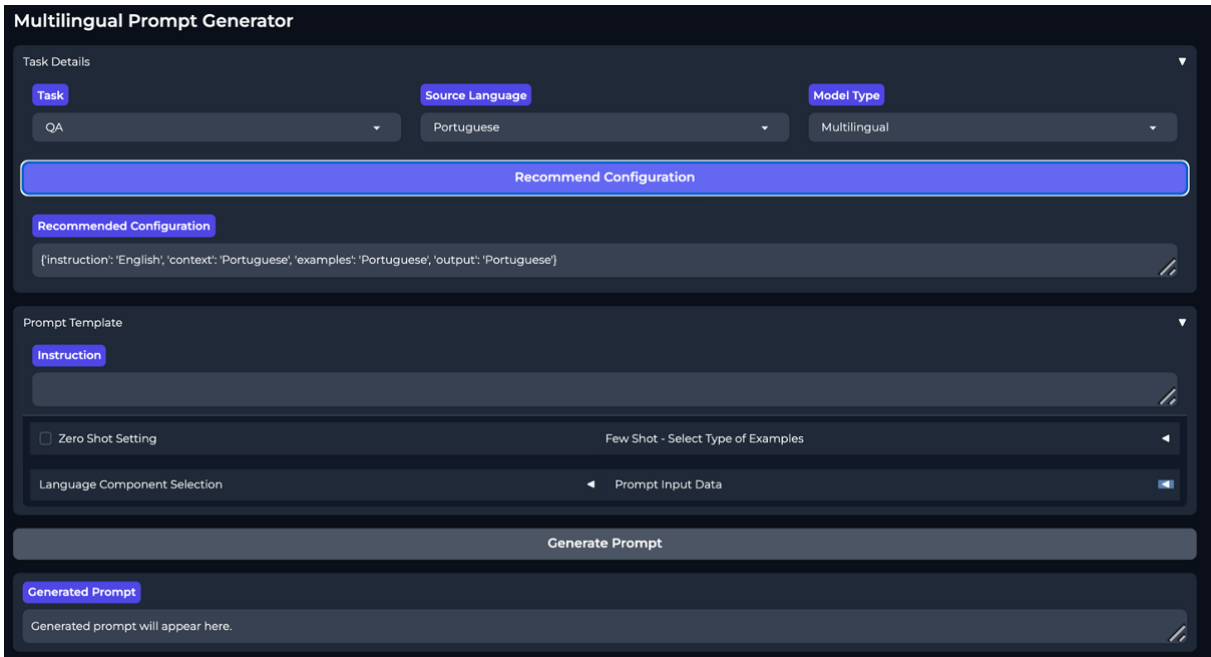


Figure 9: *Recommended Configuration*: Overview of our application. From top to bottom: Task Details—general configuration about the task (Task, Language, Model). Clicking on "Recommended Configuration" provides a suggested selective pre-translation configuration.

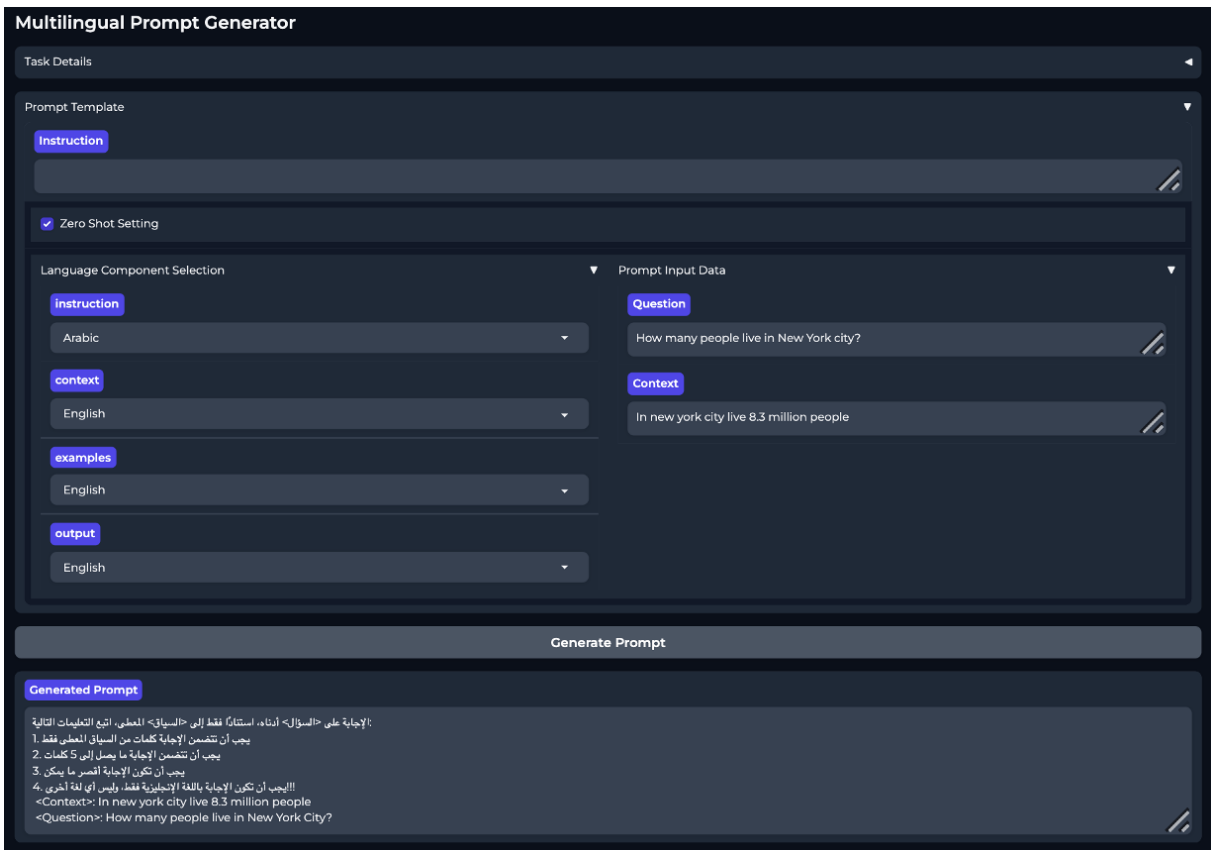


Figure 10: *Generating selective pre-translation prompt for zero-shot*: The user needs to configure the instruction (optional) and the languages for the components under "Language Component Selection": instruction, context, examples, and output. Additionally, under "Prompt Input Data," the user must configure the relevant input data or task, such as the question and context for QA in this example. Clicking on "Generate Prompt" provides a zero-shot pre-translation prompt

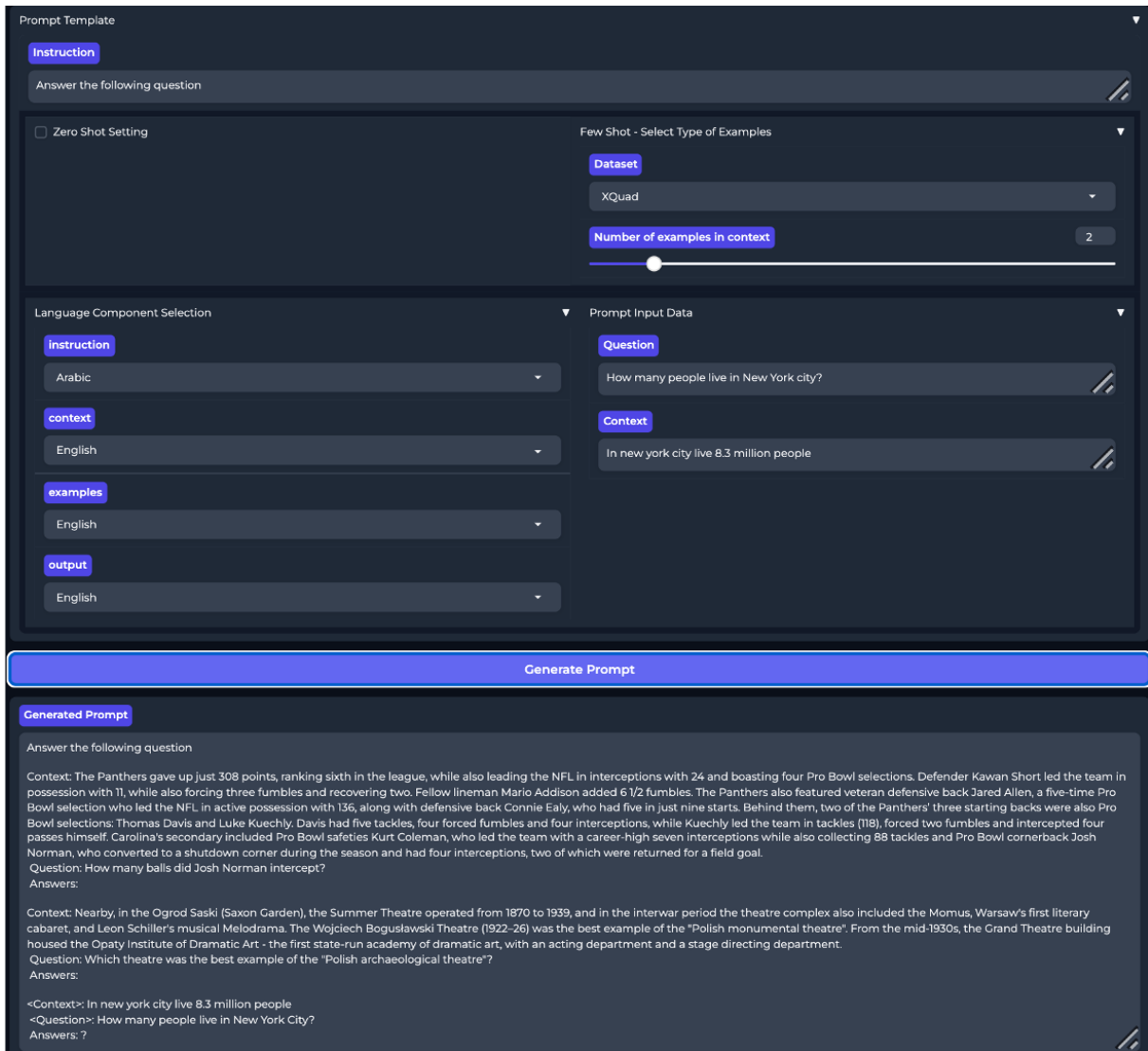
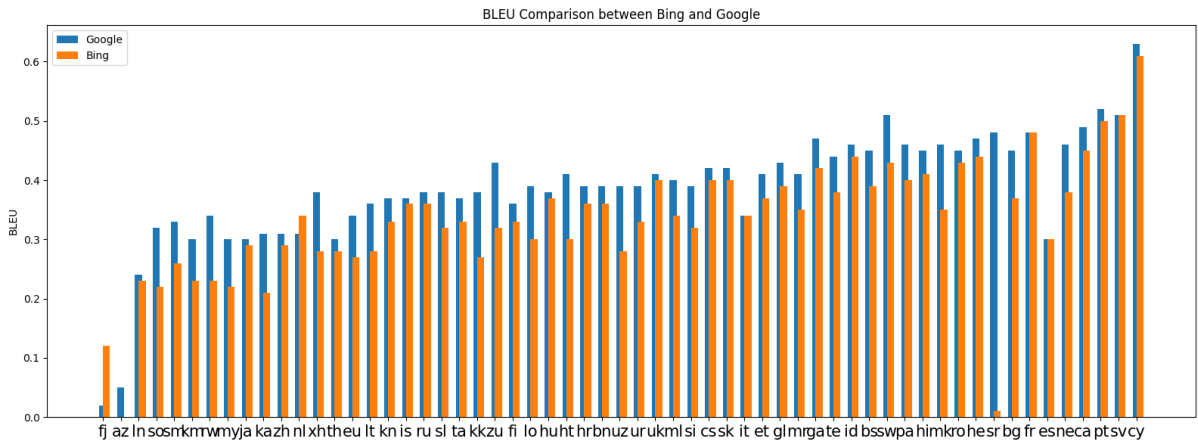
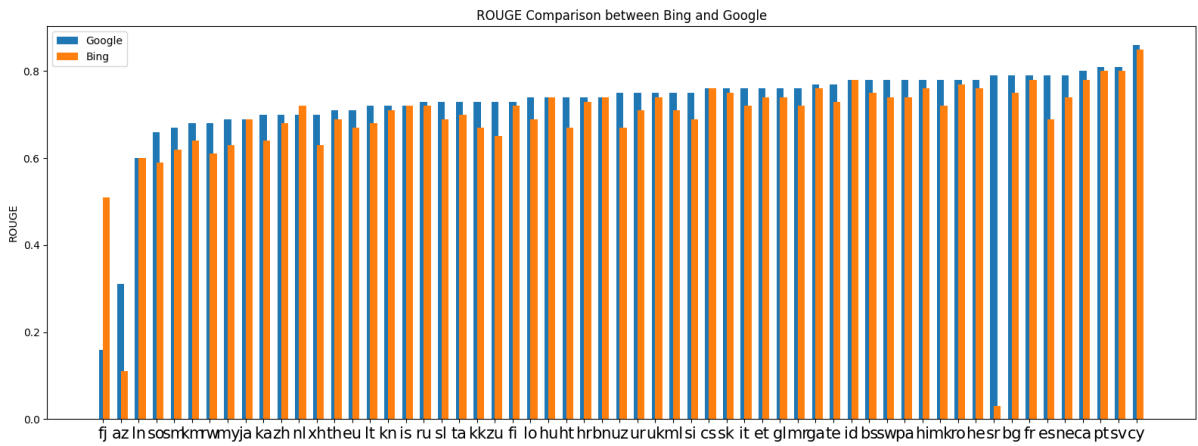


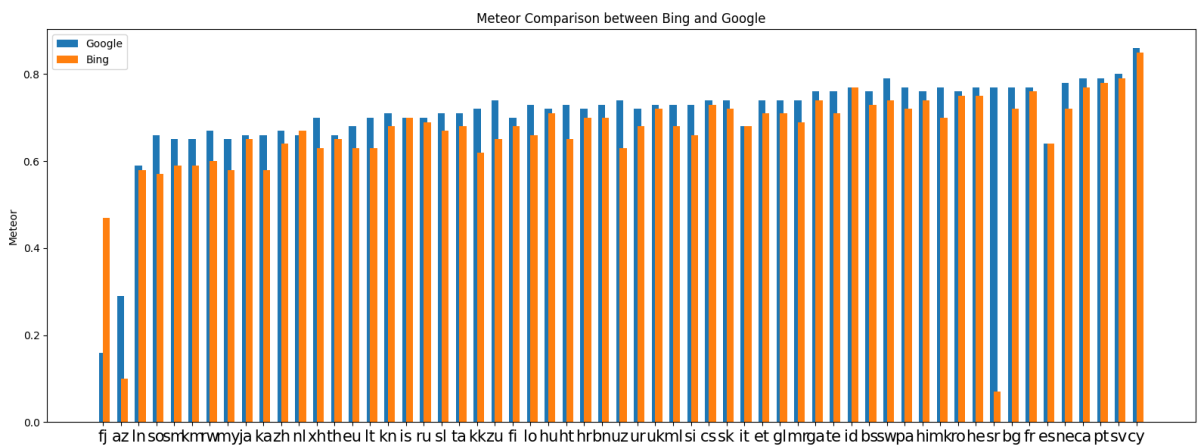
Figure 11: *Generating selective pre-translation prompt for few-shot*: Here, the user must also configure the few-shot settings: the dataset to use (from which the few-shot examples are taken) and the number of examples to use (default = 1).



(a) BLEU



(b) ROUGE



(c) Meteor

Figure 12: Google Translate API vs Bing Translator Comparison

Language	Lang Code	Number of Tokens (M)	Percentage of Tokens	Class
English	en	181,015	92.64%	A
French	fr	3,553	1.81853%	A
German	de	2,871	1.46937%	A
Spanish	es	1,510	0.77289%	A
Italian	it	1,188	0.60793%	A
Portuguese	po	1,025	0.52483%	A
Russian	ru	368	0.18843%	A
Romanian	ro	308	0.15773%	A
Swedish	sv	221	0.11307%	A
Japanese	ja	217	0.11109%	A
Chinese	zh	194	0.09905%	B
Indonesian	id	117	0.05985%	B
Turkish	tr	116	0.05944%	B
Vietnamese	vi	83	0.04252%	B
Greek	el	62	0.03153%	B
Arabic	ar	61	0.03114%	B
Serbian	sr	53	0.02706%	B
Korean	ko	33	0.01697%	B
Slovak	sk	28	0.01431%	B
Thai	th	27	0.01372%	B
Slovenian	sl	26	0.01333%	B
Persian	fa	17	0.00856%	C
Hebrew	he	15	0.00769%	C
Hindi	hi	9	0.00483%	C
Bulgarian	bg	6	0.00303%	C
Bengali	bn	3	0.00154%	C
Malayalam	ml	3	0.00165%	C
Azerbaijani	az	2	0.00128%	C
Telugu	te	2	0.00084%	C
Uzbek	uz	1.5	0.00075%	C
Nepali	ne	1.1	0.00057%	C
Urdu	ur	0.7	0.00035%	C
Swahili	sw	0.6	0.00030%	C
Assamese	as	0	0.00000%	D
Bambara	bam	0	0.00000%	D
Ewe	ee	0	0.00000%	D
Hausa	hau	0	0.00000%	D
Yoruba	yor	0	0.00000%	D

Table 12: List of languages, language codes, number of tokens in pre-trained GPT-3 data, data ratios. The languages are grouped into four classes based on their data ratios in the GPT-3 pre-trained data: High Resource ($H > 0.1\%$), Medium Resource ($M > 0.01\%$), and Low Resource ($L < 0.01\%$), and extremely low resource for unrepresented languages.

Configuration				Chinese			French			Italian			Portuguese			Serbian			Slovak			Swedish		
<i>P</i>	<i>I</i>	<i>C</i>	<i>O</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>
E	E	E	E	0.00	0.01	0.23	0.53	0.58	0.15	0.38	0.53	0.23	0.56	0.60	0.41	0.20	0.21	0.00	0.48	0.53	0.19	0.59	0.62	0.41
E	E	E	S	0.00	0.00	0.18	0.37	0.63	0.17	0.31	0.52	0.23	0.55	0.60	0.42	0.22	0.26	0.01	0.48	0.50	0.21	0.59	0.58	0.38
E	E	S	E	0.00	0.01	0.21	0.54	0.6	0.25	0.66	0.66	0.52	0.56	0.58	0.29	0.17	0.18	0.01	0.45	0.51	0.31	0.59	0.59	0.42
E	E	S	S	0.00	0.06	0.22	0.55	0.61	0.00	0.59	0.68	0.48	0.55	0.61	0.38	0.13	0.13	0.02	0.46	0.48	0.24	0.57	0.56	0.21
E	E	Z	E	0.00	0.00	0.22	0.45	0.54	N.A	0.51	0.63	0.49	0.45	0.57	0.34	0.13	0.12	0.02	0.47	0.47	0.28	0.51	0.56	0.32
E	E	S	S	0.00	0.00	0.20	0.44	0.53	0.0	0.44	0.63	0.50	0.43	0.54	0.35	0.09	0.11	0.03	0.38	0.49	0.29	0.44	0.59	0.36
E	S	E	E	0.00	0.00	0.23	0.55	0.60	0.25	0.34	0.48	0.26	0.55	0.60	0.37	0.21	0.22	0.01	0.39	0.50	0.18	0.60	0.60	0.27
E	S	E	S	0.00	0.00	0.28	0.53	0.59	0.21	0.30	0.48	0.26	0.55	0.61	0.40	0.17	0.21	0.01	0.47	0.52	0.30	0.58	0.64	0.26
E	S	S	E	0.01	0.02	0.22	0.42	0.59	0.26	0.65	0.63	0.52	0.58	0.61	0.40	0.14	0.17	0.01	0.46	0.49	0.25	0.57	0.61	0.22
E	S	S	S	0.01	0.03	0.24	0.45	0.59	0.28	0.61	0.67	0.49	0.56	0.60	0.37	0.16	0.16	0.01	0.46	0.36	N.A	0.55	0.55	0.22
E	S	Z	E	0.00	0.00	0.24	0.48	0.55	0.36	0.53	0.64	0.50	0.47	0.55	0.35	0.12	0.12	0.02	0.46	0.48	0.29	0.52	0.56	0.33
E	S	Z	S	0.01	0.00	0.20	0.45	0.55	0.38	0.52	0.62	0.48	0.46	0.53	0.33	0.11	0.12	0.03	0.41	0.50	0.29	0.50	0.57	0.36
S	E	E	E	0.06	0.02	0.09	0.61	0.69	0.32	0.40	0.53	0.33	0.60	0.68	0.49	0.57	0.20	0.05	0.66	0.62	0.26	0.63	0.64	0.14
S	E	E	S	0.11	0.05	0.07	0.68	0.71	0.31	0.36	0.53	0.29	0.64	0.66	0.54	0.68	0.64	0.3	0.61	0.61	0.31	0.60	0.64	0.24
S	E	S	E	0.61	0.61	0.00	0.64	0.72	0.32	0.69	0.74	0.64	0.71	0.69	0.53	0.72	0.75	0.47	0.73	0.71	0.39	0.69	0.67	0.34
S	E	S	S	0.59	0.63	0.00	0.63	0.66	0.35	0.69	0.75	0.62	0.76	0.70	0.48	0.77	0.68	0.46	0.69	0.72	0.54	0.66	0.68	0.35
S	E	Z	E	0.07	0.07	0.02	0.48	0.59	0.48	0.54	0.66	0.56	0.50	0.63	0.48	0.22	0.42	0.25	0.57	0.65	0.49	0.50	0.57	0.37
S	E	Z	S	0.16	0.05	0.00	0.55	0.60	0.48	0.48	0.67	0.57	0.48	0.62	0.48	0.47	0.5	0.32	0.51	0.62	0.45	0.50	0.62	0.39
S	S	E	E	0.17	0.02	0.01	0.64	0.68	0.29	0.38	0.54	0.32	0.66	0.66	0.48	0.53	0.17	0.05	0.67	0.29	0.27	0.61	0.63	0.18
S	S	E	S	0.14	0.02	0.07	0.59	0.68	0.31	0.35	0.59	0.32	0.69	0.65	0.47	0.63	0.61	0.29	0.65	0.38	0.52	0.62	0.63	0.25
S	S	S	E	0.60	0.61	0.00	0.63	0.70	0.45	0.70	0.74	0.59	0.70	0.72	0.51	0.72	0.77	0.45	0.72	0.58	0.40	0.67	0.67	0.53
S	S	S	S	0.58	0.62	0.01	0.64	0.69	0.31	0.69	0.71	0.61	0.68	0.72	0.48	0.77	0.72	0.44	0.69	0.56	0.40	0.65	0.66	0.53
S	S	Z	E	0.12	0.08	0.02	0.57	0.57	0.45	0.57	0.67	0.56	0.51	0.63	0.48	0.28	0.46	0.29	0.55	0.65	0.47	0.54	0.58	0.37
S	S	Z	S	0.13	0.04	0.00	0.49	0.58	0.48	0.55	0.68	0.58	0.48	0.61	0.48	0.42	0.50	0.33	0.52	0.63	0.45	0.53	0.60	0.39

Table 16: Comparing performance of different models on all languages in WikiANN. Metric: F1 Score.

Configuration				Bambara			Ewe			Hausa			Yoruba		
<i>P</i>	<i>I</i>	<i>C</i>	<i>O</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>	<i>gemi</i>	<i>gpt</i>	<i>mixtral</i>
E	E	E	E	0.16	0.18	0.10	0.42	0.47	0.22	0.45	0.24	0.26	0.06	0.09	0.03
E	E	E	S	0.17	0.17	0.07	0.43	0.48	0.23	0.46	0.61	0.31	0.07	0.09	0.05
E	E	S	E	0.15	0.17	0.12	0.44	0.52	0.22	0.45	0.52	0.20	0.07	0.08	0.04
E	E	S	S	0.14	0.15	0.07	0.39	0.5	0.17	0.43	0.39	0.10	0.06	0.08	0.03
E	E	Z	E	0.13	0.09	N.A	0.41	0.44	0.23	0.52	0.59	0.26	0.05	0.08	0.04
E	E	S	S	0.13	0.16	0.06	0.41	0.47	0.21	0.45	0.60	0.29	0.04	0.08	0.03
E	S	E	E	0.16	0.18	0.10	0.44	0.47	0.22	0.48	0.59	0.32	0.07	0.08	0.03
E	S	E	S	0.16	0.18	0.08	0.40	0.47	0.15	0.48	0.58	0.34	0.06	0.08	0.05
E	S	S	E	0.15	0.15	0.06	0.45	0.50	0.28	0.50	0.54	0.22	0.08	0.10	0.05
E	S	S	S	0.11	0.15	0.06	0.39	0.49	0.15	0.45	0.42	0.11	0.06	0.07	0.06
E	S	Z	E	0.15	0.28	0.19	0.40	0.46	0.23	0.54	0.62	0.27	0.08	0.08	0.05
E	S	Z	S	0.17	0.13	0.05	0.38	0.44	0.25	0.47	0.28	0.29	0.06	0.09	0.06
S	E	E	E	0.09	0.24	0.00	0.31	0.47	0.01	0.46	0.43	0.00	0.08	0.20	0.05
S	E	E	S	0.27	0.27	0.06	0.40	0.52	0.04	0.61	0.58	0.10	0.09	0.23	0.10
S	E	S	E	0.28	0.32	0.14	0.56	0.68	0.47	0.56	0.70	0.30	0.26	0.32	0.15
S	E	S	S	0.26	0.32	0.15	0.54	0.66	0.51	0.55	0.70	0.32	0.23	0.29	0.20
S	E	Z	E	0.07	0.09	0.06	0.31	0.39	0.00	0.48	0.70	0.00	0.20	0.20	0.01
S	E	Z	S	0.21	0.20	0.17	0.43	0.42	0.42	0.57	0.69	0.07	0.17	0.26	0.07
S	S	E	E	0.10	0.25	0.00	0.27	0.39	0.00	0.45	0.43	0.00	0.08	0.21	0.04
S	S	E	S	0.23	0.27	0.05	0.40	0.51	0.08	0.54	0.60	0.09	0.13	0.22	0.05
S	S	S	E	0.29	0.33	0.21	0.61	0.67	0.43	0.58	0.67	0.34	0.27	0.30	0.19
S	S	S	S	0.27	0.32	0.21	0.61	0.63	0.45	0.59	0.69	0.41	0.26	0.31	N.A
S	S	Z	E	0.08	0.18	0.00	0.31	0.37	0.00	0.46	0.46	0.01	0.03	0.22	0.03
S	S	Z	S	0.23	0.28	0.04	0.39	0.46	0.06	0.06	0.64	0.05	0.07	0.26	0.07

Table 17: Comparing performance of different models on all languages in MasakhaNER. Metric: F1 Score.

