

Shortcomings of LLMs for Low-Resource Translation: Retrieval and Understanding are Both the Problem

Anonymous ACL submission

Abstract

This work investigates the in-context learning abilities of pretrained large language models (LLMs) when instructed to translate text from a low-resource language into a high-resource language as part of an automated machine translation pipeline. We conduct a set of experiments translating Southern Quechua to Spanish and examine the informativity of various types of information retrieved from a constrained database of digitized pedagogical materials (dictionaries and grammar lessons) and parallel corpora. Using both automatic and human evaluation of model output, we conduct ablation studies that manipulate (1) context type (morpheme translations, grammar descriptions, and corpus examples), (2) retrieval methods (automated vs. manual), and (3) model type. Our results suggest that even relatively small LLMs are capable of utilizing prompt context for zero-shot low-resource translation when provided a minimally sufficient amount of relevant linguistic information. However, the variable effects of prompt type, retrieval method, model type, and language community-specific factors highlight the limitations of using even the best LLMs as translation systems for the majority of the world’s 7,000+ languages and their speakers.

1 Introduction

The field has made great progress improving the quality of machine translation (MT) systems, but constraints on the amount and kinds of data available in the majority of the world’s 7,000+ languages have led to yet another disparity in access and support for speakers of these languages: low-resource MT continues to be a major challenge (Hendy et al., 2023; Stap and Araabi, 2023; Robinson et al., 2023; Nicholas and Bhatia, 2023). While many of these languages lack the kinds of large, standardized corpora necessary for traditional methods, recent work shows it may be possible to leverage a smaller amount of existing resources, for example pedagogical materials

used for language instruction, with Large Language Models (LLMs), albeit with varying results (Tanzer et al., 2024; Zhang et al., 2024; Elsner and Needle, 2023). These materials are often the result of community-driven or government-led initiatives to support language revitalization, reclamation, and mother-tongue education (Schreiner et al.; Riestenberg et al., 2024; Liu et al., 2022). Such discrepancies in the needs and priorities of academic, commercial, and community-led efforts to develop digital resources and language technologies is what Gessler (2022) terms the “NLP Gap”.

In this study, we investigate one way to lessen the NLP Gap, comparing LLMs’ in-context learning abilities when translating from a low-resource language (a Peruvian variety of Southern Quechua) to a high-resource language (Spanish) using information retrieved from a database of pedagogical materials. We replicate results of earlier studies on a new language pair by comparing the effects of morpheme translations, sentences from a parallel corpus, and passages from a grammar instruction document on translation quality. We then conduct a more focused analysis by annotating translation outputs by hand using a modified MQM error typology (Burchardt, 2013). Finally, we conduct an ablation study on the effects of automated retrieval by manually constructing prompts using the same set of materials.

Our results suggest that while, unsurprisingly, translation quality improves with model size, such improvements seem to primarily be the result of previous exposure to the low-resource language during model pretraining, rather an improved ability for the model to utilize prompt context, as evidenced by high scores in response to baseline (zero-shot) translation prompts. However, we also find evidence that in-context learning abilities may be inconsistent across different models of similar size. As found in previous studies, prompts containing morpheme and word-level translations reliably

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

084 improve model outputs, but information from the
085 grammar and corpus have a null or even negative
086 effect on results. Human evaluation on a selec-
087 tion of outputs from two models – GPT-3.5 Turbo
088 and GPT-4o – align with the quantitative measures
089 we obtain using BLEURT (Sellam et al., 2020)
090 as an automatic metric. Quantitative results also
091 show an effect of automated retrieval on translation
092 quality that is most evident in prompts containing
093 morpheme translations and for models with lower
094 baseline scores. Finally, we highlight a number
095 of ethical concerns and limitations that arise from
096 the proposed methods that are supported by our
097 findings, and discuss the potential risks and chal-
098 lenges LLM-based methods for low-resource MT
099 face moving forward.

100 2 LLMs for Machine Translation

101 Modern LLMs are now capable of translating many
102 high-resource languages, but lack sufficient cov-
103 erage of even modestly resourced languages to
104 achieve comparable results without additional sup-
105 port (Kocmi et al., 2023). Retrieval-augmented
106 generation (Rubin et al., 2022) may provide such
107 support in the form of parallel sentences (Agrawal
108 et al., 2022), dictionary definitions (Ghazvinine-
109 jad et al., 2023; Lu et al., 2023) or other linguistic
110 meta-knowledge such as a grammatical description.
111 Retrieval-augmented methods offer exciting possi-
112 bilities for low-resource translation, since the LLM
113 might (in principle) be able to “teach itself” the
114 language from learner-oriented resources produced
115 by community members or language specialists.

116 Studies to date (Reid et al., 2024; Zhang et al.,
117 2024; Elsner and Needle, 2023) experiment with
118 four dimensions of variability: source language,
119 LLM, type(s) of information retrieved, and retrieval
120 method. Since the source languages in these stud-
121 ies have relatively little presence in public corpora
122 or on the web, differing results across LLMs can
123 tentatively be attributed to differences in their in-
124 context learning and instruction following abilities.

125 All studies find that word-level translations are
126 helpful additions to prompts. Zhang et al. (2024)
127 and Tanzer et al. (2024) also add sentence pairs
128 from a parallel corpus, while Elsner and Needle
129 (2023) add usage examples from a dictionary. Each
130 improve results, although to a lesser degree. Elsner
131 and Needle (2023) and Zhang et al. (2024) experi-
132 ment with small fixed “grammar lesson” passages
133 to provide explicit syntactic instruction, but find

134 these ineffective. Tanzer et al. (2024) uses passages
135 retrieved from a grammar book, also with relatively
136 disappointing results. Reid et al. (2024) use the en-
137 tire grammar book and a very long-context model
138 to obtain better translations, but without exploring
139 the role explicit grammar instruction actually plays
140 in doing so.

141 Zhang et al. (2024) find that sentences from the
142 corpus retrieved using BM25 embeddings work
143 better than random ones. Tanzer et al. (2024), how-
144 ever, report that retrieval with longest common
145 substring (LCS) matching outperforms embedding-
146 based retrieval. Overall, the question of how to
147 best retrieve relevant passages containing grammar
148 material or sentences in a low-resource language is
149 still open. This also complicates the interpretation
150 of the mostly-negative results found for grammar
151 passages. It is not clear whether these stem from
152 poor retrieval, from the LLMs’ inability to process
153 the retrieved content, or both. Moreover, although
154 Reid et al. (2024) conducts human evaluation of
155 the results for quality, to the best of our knowledge
156 no study to date systematically investigates specific
157 grammatical errors in the output.

158 3 Quechuan Languages

159 Quechua is a family of languages indigenous to the
160 Andes in South America. This study focuses on
161 varieties of Southern Quechua (S. Quechua, also
162 known as *urin quechua* or *quechua sureño*) spoken
163 in parts of Peru.¹ While previous studies investi-
164 gated language/LLM pairs for which the baseline
165 LLM lacked any pretrained knowledge, we find
166 that newer LLMs can translate some S. Quechua
167 sentences in a zero-shot setting. We expect this to
168 be typical of many low-resource languages which,
169 while often endangered, still may have some pres-
170 ence on the web.

171 Quechuan languages have by far the largest rep-
172 resentation of all indigenous Latin American lan-
173 guages in NLP research (Tonja et al., 2024) and
174 are often included in ACL-affiliated workshops,
175 datasets, and shared tasks (Ebrahimi et al., 2022,
176 2023; Cotterell et al., 2020). S. Quechua has a ro-
177 bust language toolkit (Rios, 2015), including the
178 morphological parser we use in our pipeline. It has
179 also been the subject of numerous studies on MT
180 for both text and speech, developed in conjunction
181 with monolingual and parallel corpora (Rios, 2015;

¹Unless noted otherwise, we use *Quechua* in this study to refer Southern Quechua and related varieties, following the practices of native speakers with whom we have relationships.

[TAREA] Traduce la siguiente frase del quechua al español.
Responde sólo con la traducción:
quechua: kay wasiqa turyipam
español:

Figure 1: Example BASELINE prompt. English: [TASK] Translate the following sentence from Quechua to Spanish. Respond only with the translation: Quechua: kay wasiqa turyipam; Spanish:

Cardenas et al., 2018; Ortega et al., 2020; Zevallos et al., 2022). Nonetheless, such tools continue to face challenges, and Quechuan languages continue to lack the resources necessary to develop most of today’s state of the art models.

4 Methods

4.1 Data

We conduct experiments on a collection of 50 pairs of S. Quechua - Spanish sentences sourced from one of the author’s personal notes. These were selected to highlight a range of specific grammatical phenomena at multiple levels of difficulty—they include simple clauses and tenses (e.g., *qam allinta tusunki* (tu bailas bien) ‘you dance well’) as well as more advanced constructions such as those involving past participles (e.g., *awasqay waliqa sumaqmi* (la falda que tejí es linda) ‘the skirt I knit is lovely’ and simultaneous events (e.g., *qamqa takita uyarispa wasiykita pichachkanki* (tú estás limpiando tu casa escuchando música) ‘you’re cleaning your house listening to music’). The first author, a foreign-language student of S. Quechua, received permission from her instructor to use notes from their lessons for the study. All sentence pairs were inspected by the instructor, a native bilingual speaker of both S. Quechua and Peruvian Spanish, to eliminate any errors and confirm the accuracy of all reference translations.

4.2 Prompt Construction

As a baseline, each sentence is inserted into a prompt template that instructs the model in Spanish to translate the S. Quechua sentence into Spanish and respond only with the translation (Figure 1). We automate a process for building on this template and compare the effects of adding information from three different sources to the prompt context.

4.2.1 Morpheme Translations (MORPH)

We use a morphological parser (Rios, 2015) to segment each word of the source segment into canonical morphemes, each with gloss symbols and a

Spanish translation.² Some morphemes have multiple candidate meanings, all of which are retrieved. As an example, the word *rantikuq* is segmented as *ranti-ku-q* and glossed as “comprar.DB.VRoot-DB.VDeriv.+RflxInt-+Ag.NS.” While numerous orthographic standards have been developed and promoted across Quechuan-speaking communities in South America, considerable variation in orthographic conventions may be found even within a particular community or variety (Rios and Castro Mamani, 2014). We discuss the implications of this for our results in Section 4.2.5.

We supplement the output from the parser using a Quechua-Spanish bilingual dictionary (Qheswa Simi Hamut’ana Kurak Suntur, 2005). We retrieve any dictionary entry whose headword exactly matches a canonical morpheme in our segmentation. By default, we include all senses and any usage examples or contextual information in the dictionary entry as part of the prompt. We then concatenate the output of the parser with the retrieved dictionary entries and include this MORPH information as prompt context preceding the source sentence and baseline translation prompt.

4.2.2 Grammar Descriptions (GRAMMAR)

We also experiment with the inclusion of grammar lessons found in student-facing pedagogical materials, retrieving grammatical explanations relevant to each source sentence from a PDF document developed for students and teachers of S. Quechua. (Pinto Tapia et al., 2005). The document is organized into short sections (1-3 sentences, plus paradigm tables or usage examples) that describe the particular grammatical concept associated with an affix in Quechua. For each source sentence, we retrieve sections associated with any affix listed in the document that is an exact match of a canonical morpheme and include this in prompts using contextual information from the grammar. This improves on the methods described in Tanzer et al. (2024), who use LCS-based retrieval over an entire textbook, and Elsner and Needle (2023); Zhang et al. (2024), whose grammatical description remains consistent across prompts regardless of the source text being translated.

²We set aside valid concerns regarding the theoretical status of the *morpheme* for this study and define a morph(eme) loosely as a recognizable form-meaning pair that recurs in a language.

4.2.3 Parallel Usage Examples (CORPUS)

Finally, we experiment with sentence-level examples from a S. Quechua-Spanish parallel corpus designed for traditional NLP tasks. We combine data made available via the AmericasNLP 2021 Shared Task on Open Machine Translation and the 2023 IWSLT shared task on low-resource SLT (Mager et al., 2021; Agarwal et al., 2023; Agić and Vulić, 2019; Ortega et al., 2020; Tiedemann, 2012). For each source sentence, we retrieve the three best matches from the corpus using a LCS search against the full source sentence.

4.2.4 Combined prompt types

Combinations of information from all three sources yields 8 total conditions, including the baseline. An example prompt from each information source is given in Appendix E.

4.2.5 Manually Revised Prompts

To compute a soft upper bound on the improvements possible with better retrieval, we conduct an additional set of experiments using manually revised prompts. We first examine the content retrieved from the morphological parser, dictionary, and grammar document and remove all instances of ambiguity and irrelevant or misleading information from the prompt context.³

For example, many S. Quechua speakers use the term *runasimi* (lit: ‘people mouth’, ‘the people’s language’), as an endonym for the language. The parser, however, returns only the literal decomposition (*runa* ‘ser humanos’/‘people’ and *simi* ‘boca’/‘mouth’), and the dictionary does not list *runasimi* as a headword but rather as one of eight different senses of *simi*. We thus remove all such irrelevant examples and translations from the prompt and retain only the content indicating a translation of *runasimi* in the linguistic sense.

We also manually retrieve content from the dictionary and grammar documents that were overlooked by the automated retriever. For example, the verb *yanuy* ‘to cook’ does not appear as a headword in the dictionary, but rather as a regional variant of *wayk’uy* ‘to cook’. We also eliminate content from the grammar that was retrieved because of syncretism, or mistakes that cascaded from the morphological parser to result in irrelevant retrievals.

³We do not experiment with retrieval methods for corpus examples, which were retrieved using LCS match in both conditions. Improving on LCS-based retrieval remains an open question in low-resource LLM-MT, and we leave this for future work.

We manually parse each source sentence to only retrieve and include relevant information in the prompt context. All content in the revised prompts is sourced from the same material available to the automated retriever systems, and we do not add any additional information or use supplemental materials of any sort to create the revised prompts.

4.3 Models

We experiment with three proprietary models, GPT-3.5 Turbo (gpt-3.5-turbo-0125, Brown et al., 2020), GPT-4o (gpt-4o, Achiam et al., 2023), and Gemini 1.5 Pro (gemini-1.5-pro, Reid et al., 2024), and one open-source model, Llama 3 (llama-3-8b-instruct, AI@Meta, 2024). We use the pretrained models with their default settings, and do not adjust hyperparameters or conduct any finetuning as part of our experiments.

4.4 Evaluation

We conduct both automatic and human evaluations to identify trends in model errors and outputs in the various experimental conditions. We use BLEURT as an automatic metric, and report mean BLEURT scores across items as the primary quantitative measure of translation quality for each of the conditions and models. We also use an adapted MQM schema to conduct qualitative human evaluation of the outputs of GPT-3.5 and GPT-4o for all prompts with automatic retrieval.

Each item selected for human evaluation is annotated by at least one of the authors by comparing the model’s output to the source text and reference translation.⁴ We refer to the complete MQM typology to design our own four-dimensional framework of commonly attested errors in LLM-MT, each with a defined set of specific subtypes. Precise definitions and examples for all error categories and subtypes may be found in Appendix D.

Many of the categories in our schema are defined as in the core MQM framework. However, to capture some of the key behaviors reported in previous studies on LLM-MT and to evaluate the effects of prompt type on model outputs, we make the following adjustments. First, we utilize the Addition and Omission errors defined as Accuracy subtypes in the original MQM typology, but distinguish these from three additional subtypes: Substitution - Incorrect Subject, Substitution - Incorrect

⁴We discuss limitations on this process given the authors’ respective proficiency levels in S. Quechua and Spanish, as well as the steps we take to address them, in Section 8.

Tense/Aspect/Modality (TAM), and Substitution - Other. This is intended to capture LLM translations that differ from the source in terms of discrete lexical material or case, person, number, and/or TAM markings while otherwise maintaining the lexical and structural content needed to appropriately translate the source text.

Rather than including Mistranslation and MT Hallucination as Accuracy Error subtypes as in the original MQM typology, we define a separate Non-Translation category with three possible subtypes: Complete Mistranslation, Mistranslation with Lexical Correspondences, and Refusal. The third dimension of our typology, Model Error, was ultimately not used to classify any output in this study, but characterizes more generic model “misbehavior” such as failing to follow instructions, producing garbled text, or inappropriately generating content in the source language. Finally, Target Errors identify outputs that are ungrammatical, stylistically inappropriate, or semantically incoherent in the target language, regardless of their accuracy.

Detailed annotation guidelines were drafted and agreed upon to encourage consistency across annotators and experimental items. Annotators are instructed to identify and tag up to three specific errors for each translation output, with the exception of Target Errors, which do not count towards the three-error maximum. Each model output is also tagged for quality along a four-point scale as defined in Table 7.

Before proceeding with annotation over the larger dataset, both annotators also completed a test evaluation of the same 12 experimental items (96 sentences total) to assess inter-annotator agreement. Statistical measures ($\kappa = 0.72$ for quality judgments, $\alpha = 0.55$ for error categories) indicated some discrepancies in annotator judgments, especially for categories, since determining the three most important errors is especially subjective. These were identified and discussed, and agreement was ultimately deemed sufficient to proceed.

5 Results

5.1 Quality metrics

We present BLEURT scores for prompts generated using automated retrieval in Table 1 and summarize human quality judgments for GPT-3.5 and GPT-4o with automated retrieval in Table 2. The complete distribution of quality ratings across all prompt types for these two models is provided in Appendix

| | GPT3.5 | GPT4o | Gem. | Lla3 |
|--------|--------|-------|------|------|
| BASE | 0.19 | 0.66 | 0.56 | 0.15 |
| CORPUS | 0.27 | 0.59 | 0.49 | 0.19 |
| GRAM | 0.23 | 0.56 | 0.55 | 0.17 |
| MORPH | 0.44 | 0.54 | 0.61 | 0.39 |
| C+G | 0.26 | 0.59 | 0.54 | 0.21 |
| C+M | 0.44 | 0.59 | 0.59 | 0.36 |
| G+M | 0.41 | 0.53 | 0.61 | 0.39 |
| C+G+M | 0.43 | 0.57 | 0.61 | 0.15 |

Table 1: Mean BLEURT scores by LLM and prompt type. Shaded rows include morpheme contexts.

| LLM | GPT-3.5 | GPT-4o |
|----------------|---------|--------|
| BASE | 21 | 108 |
| CORPUS | 43 | 101 |
| GRAMMAR | 33 | 99 |
| MORPH | 79 | 102 |
| CORPUS-GRAMMAR | 41 | 101 |
| C+M | 75 | 110 |
| G+M | 68 | 100 |
| C+G+M | 77 | 109 |

Table 2: Human-annotated quality ratings summarized as $3 \times high + 2 \times med + low$. Shaded rows include morpheme contexts.

| | GPT3.5 | GPT4 | Gem. | Lla3 |
|----------|--------|------|------|------|
| G-AUTO | 0.23 | 0.56 | 0.55 | 0.17 |
| G-MAN | 0.24 | 0.58 | 0.54 | 0.15 |
| M-AUTO | 0.44 | 0.54 | 0.61 | 0.39 |
| M-MAN | 0.56 | 0.63 | 0.66 | 0.49 |
| CGM-AUTO | 0.43 | 0.57 | 0.61 | 0.15 |
| CGM-MAN | 0.54 | 0.63 | 0.63 | 0.26 |

Table 3: Comparison of mean BLEURT scores for automatic versus manual retrieval of material in GRAMMAR, MORPH, and CORPUS-GRAMMAR-MORPH prompts.

F. We find clear effects of LLM, prompt type, and retrieval method, as well as interactions among all three factors.

Gemini and GPT-4o outperform Llama 3 and GPT-3.5 for every prompt type. This gap is highest for the least informative prompts, indicating that the Llama 3 and GPT-3.5 base models have relatively poor coverage of S. Quechua, while GPT-4o and Gemini have much better coverage. The effect is evident in both BLEURT scores and human quality evaluations.

Effects of prompt type are mediated by the quality of the pretrained model. Llama 3 and GPT-3.5 show a clear improvement in quality when MORPH information is included in the prompt. Gemini also improves when this information is added, but to a lesser extent. GPT-4o, on the other hand, performs best in response to the BASELINE (zero-shot) prompts, which attain the highest BLEURT scores across all models, prompt types, and retrieval methods evaluated in this study. In other words, providing additional information in the prompt’s context actually *degrades* GPT-4o’s ability to translate from S. Quechua to Spanish in all experimental conditions.

5.2 Effects of Automated Retrieval

To highlight the effects of automated retrieval on model output, we present BLEURT scores for a selection of prompt types and all four models in Table 3 (full scores may be found in Appendix F). The effect of manual retrieval for MORPH information is positive for all models, although this gap is smallest for Gemini (probably because its performance for these prompts is already highest). The effect for GRAMMAR prompts is either minor or negative.

5.3 Human Analysis of Translation Errors

The most common error type identified by the annotators is Substitution - Other, which includes a diverse assortment of lexical and phrasal incongruencies of varying degrees of severity. These are largely item-specific and therefore hard to characterize as a group. Using the error categories described in Section 4.4, we instead identify three more clearly interpretable phenomena and provide a detailed discussion of each in the following sections. We present counts for selected prompt types in Table 4, with examples in Appendix A and counts for all errors in Appendix G.

5.4 Mistranslations

Outright mistranslations are most common for GPT-3.5, making up 30 of the 50 responses in the BASELINE condition. We also consider outputs that retain only minimal traces of the source content, which we label as Mistranslations with Lexical Correspondence. Approximately 1/3 of the 637 total errors tagged across all prompt types for GPT-3.5 are mistranslations of either type, roughly split between complete mistranslations and those with lexical correspondence (15.07% and 18.37%, respectively, of all errors tagged for GPT-3.5).

As reported in previous work, adding morpheme- and word-level translations to the prompt greatly reduces the rate of this kind of response. GPT4o also produces drastically fewer mistranslations compared to its predecessor. However, it is notable that both models produce at least one mistranslation for each prompt type. In general, complete mistranslations are in fluent Spanish and contain no overt indications that something has been misrepresented. We return to the ethical implications of these errors in the Discussion.

We also note that many of the items tagged as Mistranslation with Lexical Correspondence show correspondence only for words that were already in Spanish in the source text. For example, some sentences contain Spanish loan names for the days of the week. While some of these errors are produced using deceptively fluent Spanish as described above, we find many to be accompanied by semantic incoherence or ungrammaticality in the output. We discuss such target language fluency errors in the following section.

5.5 Target Fluency

Target Fluency errors occur when the output is not grammatical, coherent, or stylistically appropriate – for instance, if an output contains a nonsensical repetition or a verb with missing arguments. Outputs of this type bear a strong similarity to human “translationese” in that structural features of the source language may surface in the translation at the expense of naturalness (Freitag et al., 2019; Koppel and Ordan, 2011). Both GPT-3.5 and GPT-4o tend to produce more such outputs when the prompt is more informative – 10 to 20% of the time (5-10 instances per 50) in prompts with morpheme translations.

| | | BASE | MORPH | C+G+M |
|--|---------|------|-------|-------|
| Mistranslation: complete + lexical correspondence | GPT-3.5 | 45 | 11 | 12 |
| | GPT-4o | 4 | 6 | 4 |
| Target Fluency: grammar + coherence + style | GPT-3.5 | 0 | 14 | 10 |
| | GPT-4o | 3 | 13 | 9 |
| Grammatical Divergence: subject + TAM | GPT-3.5 | 0 | 24 | 31 |
| | GPT-4o | 17 | 13 | 11 |

Table 4: Counts of human-annotated error types (per 50 sentences) by LLM and prompt type.

5.6 Grammatical Divergence

We group misrendered verbal subjects and tense/aspect/morphology (TAM) markers together as Grammatical Divergence errors. Such errors are distinct from the Target Fluency errors described in the previous section—the Spanish output is grammatical, but fails to accurately reflect the syntax of the source. Although they are not the only grammatical phenomena that may be similarly misrendered, we select subject and TAM markers for analysis as they are straightforward to identify and give a good indication of how well the LLMs cope with more abstract information about the meanings of functional morphemes. TAM divergences are much more prevalent than divergences in subject; for instance, only one of GPT-4o’s 13 Grammatical Divergence errors in the MORPH condition misrender the subject marker.

Grammatical Divergence errors are annotated only for sentences that are not mistranslated outright, so GPT-3.5 produces none of these in the BASELINE condition. For more informative prompts, it is clear that GPT-4o is better than GPT-3.5 at translating both functional and lexical meanings. However, a relatively large number of sentences (over 20%) still contain such an error even with the highest performing model and prompt type. The relatively small drop in error between different prompt types for GPT-4o suggests that neither the corpus-based usage examples nor example paradigms and descriptions from the grammar document can fully prevent this type of error.

6 Discussion

We observe large differences between LLMs, both in terms of the overall quality of their generated translations as well as the effects of prompt type on their outputs. GPT-4o and Gemini, which have the highest baseline scores, benefit least from additional information—their performance with CORPUS and GRAMMAR information actually decreases.

This occurs even with manually curated prompts, suggesting it is not an effect of including irrelevant material. On the other hand, it does not represent a ceiling on quality, since both models continue to make errors (GPT-4o produces 10 LOW-quality translations in our set of 50). These results suggest that even relevant grammar explanations, when written in prose with examples, do little to help the newest generation of LLMs to translate a low-resource language such as Southern Quechua.

Although GPT-4o and Gemini results are similar in many ways, we do find evidence for differences in their in-context learning abilities. Baseline prompts and the GPT-4o model produce the highest BLEURT scores across the dataset, but these outputs still show a number of errors characteristic of LLMs, particularly lexical substitution errors that are not necessarily corrected with the inclusion of more context. In contrast, Gemini, which has near-comparable performance across prompt types, shows an increase in scores when prompts include MORPH information, regardless of retrieval type, suggesting a greater ability to identify and utilize relevant word- and morph-level translations in the prompt’s context. Previous work suggests that newer builds of GPT-4 are less capable of following instructions (Chen et al., 2023); such differences may be masked by the effects of pretraining when automatically evaluating translations. This suggests that researchers should continue to carefully select and compare among different LLMs when experimenting with retrieval-based translation.

Finally, we identify a number of translation errors of varying types that appear to be due to language-specific characteristics, for example ambiguity from syncretism in grammatical markers, polysemous lexical items, or the orthographic and lexical variation discussed in Section 4.2.5. It may be possible to moderate such effects with additional refinement of the retrieval database and methods, which we leave for future work.

588 **6.1 Ethical concerns**

589 Both our work and much of the previous work in
590 this paradigm is motivated by the desire to close the
591 “NLP Gap” among researchers, community mem-
592 bers, and software developers interested in low-
593 resource language technologies. Machine trans-
594 lation is listed as a welcome topic of research by
595 some (though not all) members of American in-
596 digenous communities (Mager et al., 2023), and is
597 potentially an important tool for language learners
598 (Jolley and Maimone, 2022). Even an imperfect
599 translation system might be a useful tool for users
600 with a clear understanding of its limitations. How-
601 ever, the systems evaluated in this work have two
602 problematic tendencies that limit their potential for
603 deployment in real community settings.

604 First, unfaithful translations often tend to be
605 highly fluent (Section 5.4). While fluency ratings
606 for older MT systems correlate well with accuracy
607 scores, and have even been used as a proxy for
608 overall translation quality (Gamon et al., 2005; Es-
609 trella et al., 2007), this correlation is reversed for
610 our systems. LLMs are well-known for making
611 false statements that seem plausible and authori-
612 tative (Bickmore et al., 2018; Dinan et al., 2021);
613 this could be particularly problematic when they
614 project illusions of expertise at the expense of an
615 already marginalized group.

616 Second, some mistranslations identified in our
617 study appear to draw on stereotypes of indigenous
618 groups (Appendix 6). These are most apparent for
619 the BASELINE system and GPT-3.5, but also (less
620 frequently) occur with more informative prompts
621 and better LLMs. Stereotypical sentences can in-
622 volve flowery language with an emphasis on tradi-
623 tion or connectedness to nature (Erhart and Hall,
624 2019), as well as the unprompted addition of in-
625 digenous Andean cultural customs and products
626 (traditional medicine, chicha) to translations that
627 are otherwise faithful to the source text. The over-
628 all effect is to exoticize Southern Quechua speakers
629 and writers in ways that the original sentences do
630 not. Similar stereotypes have also been noted in
631 LLM-generated responses to open-ended prompts
632 (Cheng et al., 2023; Delgado Solorzano and Toxtli,
633 2023; Shieh et al., 2024).

634 While we prompt models to output only the trans-
635 lation for evaluation purposes, models may have
636 some capacity to explain or qualify their transla-
637 tions and give reminders for responsible use of the
638 technology. Should a retrieval-based translation

system ever be deployed in a real-world setting for
language learning, its developers should maximize
transparency by presenting the content of any re-
trieved information and its source to the user along
with the translation, reminding users directly of po-
tential inaccuracies, and offering vetted resources
for additional fact-checking when available.

7 **Conclusion**

Our results suggest a number of key limitations
and concerns regarding the use of LLMs in a low-
resource MT context, and have greater implications
for our understanding of the seemingly “humanlike”
conceptual, analytical, and in-context learning abil-
ities of LLMs.

For the majority of the world’s language com-
munities and their speakers, powering and supply-
ing LLMs with enough pretraining data to over-
come their limitations is not feasible. We therefore
offer the following suggestions to those looking
to develop low-resource LLM-MT: (1), improve
data structures and methods for interacting with a
language-specific database for retrieval-aided gen-
eration. (2), continue analysis of the mechanisms
driving in-context learning in LLMs, for example
by comparing ICL to the effects of finetuning (Dai
et al., 2023), (3) experiment with prompt structures
and techniques, for example by altering the order
information (Liu et al., 2024) or by iteratively or
prompting the model to guide its reasoning towards
a suitable translation (Wang et al., 2022).

Finally, we wish to emphasize the continued
risks of prematurely deploying this or similar meth-
ods in any low-resource language community, par-
ticularly given the vulnerability and disproportion-
ate lack of resources many such communities face
in domains where these technologies would likely
be used. As AI research continues to rapidly de-
velop, we urge those conducting it to increase com-
munity engagement, amplify the voices of those
traditionally at a disadvantage, and collaboratively
develop research infrastructures that may lessen
the NLP Gap. While there’s still much to be done
before low-resource LLM-MT may be safely im-
plemented, we believe such a tool has the poten-
tial to empower speakers of any variety, including
nonstandard varieties of traditional “high-resource”
languages such as English, to develop technologies
that reflect their preferences and serve their unique
needs.

8 Limitations

Limitations on the scope and replicability of this work may be attributed to one or more characteristics of the data and models used in this study, in addition to limitations inherent to the respective identities of its authors. First, the BLEURT scores we report are limited in their statistical validity. We have conducted some constrained tests to explore potential variance in scores, but expenses associated with text generation using proprietary models such as those developed by OpenAI and Google on a larger dataset may be prohibitive. This is compounded by the widely-acknowledged “black box” nature of the models powering both LLMs and BLEURT, as well as an increasing opacity with respect to the exact content and methods used to pretrain modern state of the art LLMs. For this reason, we focus our discussion on those results that show clear trends in both the quantitative and human evaluations we conduct.

There are also some constraints on our study and its methodology that are largely tied to linguistic factors, such as variation in orthography (and the need for digitized text-based resources as a prerequisite) as well as the lexical and grammatical variation that may be found in all languages, particularly the low-resource varieties we wish to support. Our results suggest it may be possible to guide the outputs of LLMs towards the specific usage conventions of a given community, but this is itself limited by the content of the materials used to develop the database from which prompt contexts are retrieved. Neither of the authors is a native speaker of any Quechua or Spanish varieties, and only one is a student of these languages with relationships to Quechua speakers and communities. While we have strived to be consistent in the Quechua and Spanish varieties used in our study (both the dictionary and grammar materials were provided by the same instructor who shared and proofread the 50 sentence pairs we use, and we select a morphological parser and corpora intended for use with Southern Quechua), variation is widespread among and within Quechua-speaking communities, and we do not have access to a dictionary, grammar, morphological parser, and corpus developed by a unified and consistent set of authors. Such variation is language- and community-dependent and bound to constrain potential applications of our methods. Future work should continue to explore ways to faithfully represent the diversity of linguistic con-

ventions employed by communities interested in developing such technologies.

We acknowledge, as well, limitations that arise from the size of our dataset and database and the methods used to curate them. The 50 sentence pairs we use were selected to highlight a range of specific grammatical phenomena, not all of which were well represented in our database, and differ in their structural complexity. We are grateful for the guidance provided by the Quechua instructor whose lessons were a source for such examples and proofread the sentences before their inclusion in our experiments, but limited as well by our status as non-native speakers. Human evaluation of model outputs was for this reason primarily constructed using machine-translated English texts as references, but was inspected by the Spanish- and Quechua-speaking author to remove a small number of evaluations that reflect linguistic discrepancies between Quechua, Spanish, and English or inaccuracies in the machine translated English.

9 Ethics Statement

We consulted the first author’s Quechua instructor, who gave us permission to use the sentences from the notes in this project and verified their accuracy. The instructor will be acknowledged by name if the paper is accepted. We cite the Quechua dictionary and grammar materials used to provide prompt information, and believe that our use of these materials is consonant with their original purpose. However, we do not distribute machine-readable versions of them as a contribution of this project, since this would violate the rights of the publisher.

The authors annotated the translation output themselves, so no human subjects approval/consent/compensation was required.

There are numerous ethical issues related to the training and use of LLMs, such as labor issues and energy costs. While these issues are inextricable from the methods used in this project, we believe the potential impact of making low-resource translation viable and accessible to minority language communities who want them (our primary goal in this line of research) outweighs the problems inherent in using LLMs at all.

We discuss the potential risks of deploying systems like the ones described here further in Section 6.1 of the main text.

788
789
790
791
792
793
794

795
796
797
798
799
800
801
802
803

804
805
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

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. [Gpt-4 technical report](#). *arXiv preprint arXiv:2303.08774*. See also: <https://openai.com/index/hello-gpt-4/>.

Milind Agarwal, Sweta Agrawal, Antonios Anastasopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, William Chen, et al. 2023. [FINDINGS OF THE IWSLT 2023 EVALUATION CAMPAIGN](#). In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 1–61, Toronto, Canada (in-person and online). Association for Computational Linguistics.

Željko Agić and Ivan Vulić. 2019. [JW300: A wide-coverage parallel corpus for low-resource languages](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3204–3210, Florence, Italy. Association for Computational Linguistics.

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. [In-context examples selection for machine translation](#). *Preprint*, arXiv:2212.02437.

AI@Meta. 2024. [Llama 3 model card](#).

Timothy W Bickmore, Ha Trinh, Stefan Olafsson, Teresa K O’Leary, Reza Asadi, Nathaniel M Rickles, and Ricardo Cruz. 2018. [Patient and consumer safety risks when using conversational assistants for medical information: An observational study of siri, alexa, and google assistant](#). *J Med Internet Res*, 20(9):e11510.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165. See also: <https://openai.com/index/new-embedding-models-and-api-updates/>.

Aljoscha Burchardt. 2013. [Multidimensional quality metrics: a flexible system for assessing translation quality](#). In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Ronald Cardenas, Rodolfo Zevallos, Reynaldo Baquerizo, and Luis Camacho. 2018. [Siminchik: A speech corpus for preservation of southern quechua](#). *ISLNLP*, 2:21.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023. [Analyzing chatgpt’s behavior shifts over time](#). In *RO-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*.

Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. [Marked personas: Using natural language prompts to](#)

[measure stereotypes in language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada. Association for Computational Linguistics. 843
844
845
846
847

Ryan Cotterell, Christo Kirov, John Sylak-Glassman, Géraldine Walther, Ekaterina Vylomova, Arya D. McCarthy, Katharina Kann, Sabrina J. Mielke, Garrett Nicolai, Miikka Silfverberg, David Yarowsky, Jason Eisner, and Mans Hulden. 2020. [The conll-sigmorphon 2018 shared task: Universal morphological inflection](#). *Preprint*, arXiv:1810.07125. 848
849
850
851
852
853
854

Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2023. [Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers](#). *Preprint*, arXiv:2212.10559. 855
856
857
858
859

Cecilia Delgado Solorzano and Carlos Toxtli. 2023. [Evaluating machine perception of indigeneity: An analysis of chatgpt’s perceptions of indigenous roles in diverse scenarios](#). 860
861
862
863

Emily Dinan, Gavin Abercrombie, A. Stevie Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2021. [Anticipating safety issues in e2e conversational ai: Framework and tooling](#). *Preprint*, arXiv:2107.03451. 864
865
866
867
868

Abteen Ebrahimi, Manuel Mager, Arturo Oncevay, Vishrav Chaudhary, Luis Chiruzzo, Angela Fan, John Ortega, Ricardo Ramos, Annette Rios, Ivan Vladimir Meza Ruiz, et al. 2022. [AmericasNLI: Evaluating zero-shot natural language understanding of pre-trained multilingual models in truly low-resource languages](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6279–6299, Dublin, Ireland. Association for Computational Linguistics. 869
870
871
872
873
874
875
876
877
878

Abteen Ebrahimi, Manuel Mager, Shruti Rijhwani, Enora Rice, Arturo Oncevay, Claudia Baltazar, María Cortés, Cynthia Montañó, John E. Ortega, Rolando Coto-solano, Hilaria Cruz, Alexis Palmer, and Katharina Kann. 2023. [Findings of the AmericasNLP 2023 shared task on machine translation into indigenous languages](#). In *Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)*, pages 206–219, Toronto, Canada. Association for Computational Linguistics. 879
880
881
882
883
884
885
886
887
888
889

Micha Elsner and Jordan Needle. 2023. [Translating a low-resource language using GPT-3 and a human-readable dictionary](#). In *Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 1–13, Toronto, Canada. Association for Computational Linguistics. 890
891
892
893
894
895
896

Ryan S Erhart and Deborah L Hall. 2019. [A descriptive and comparative analysis of the content of stereotypes about native americans](#). *Race and Social Problems*, 11:225–242. 897
898
899
900

| | | | |
|-----|---|--|------|
| 901 | Paula Estrella, Andrei Popescu-Belis, and Maghi King. | Zoey Liu, Crystal Richardson, Richard Hatcher, and | 958 |
| 902 | 2007. A new method for the study of correlations be- | Emily Prud'hommeaux. 2022. Not always about | 959 |
| 903 | tween MT evaluation metrics . In <i>Proceedings of the</i> | you: Prioritizing community needs when developing | 960 |
| 904 | <i>11th Conference on Theoretical and Methodological</i> | endangered language technology . In <i>Proceedings</i> | 961 |
| 905 | <i>Issues in Machine Translation of Natural Languages:</i> | of the 60th Annual Meeting of the Association for | 962 |
| 906 | <i>Papers</i> , Skövde, Sweden. | Computational Linguistics (Volume 1: Long Papers) , | 963 |
| 907 | Markus Freitag, Isaac Caswell, and Scott Roy. 2019. | pages 3933–3944 , Dublin, Ireland. Association for | 964 |
| 908 | APE at scale and its implications on MT evaluation | <i>Computational Linguistics</i> . | 965 |
| 909 | biases . In <i>Proceedings of the Fourth Conference on</i> | Hongyuan Lu, Haoyang Huang, Dongdong Zhang, Hao- | 966 |
| 910 | <i>Machine Translation (Volume 1: Research Papers)</i> , | ran Yang, Wai Lam, and Furu Wei. 2023. Chain- | 967 |
| 911 | pages 34–44, Florence, Italy. Association for Com- | of-dictionary prompting elicits translation in large | 968 |
| 912 | putational Linguistics. | language models . <i>Preprint</i> , arXiv:2305.06575. | 969 |
| 913 | Michael Gamon, Anthony Aue, and Martine Smets. | Manuel Mager, Elisabeth Mager, Katharina Kann, and | 970 |
| 914 | 2005. Sentence-level mt evaluation without refer- | Ngoc Thang Vu. 2023. Ethical considerations for | 971 |
| 915 | ence translations: Beyond language modeling. In | machine translation of indigenous languages: Giv- | 972 |
| 916 | <i>Proceedings of the 10th EAMT Conference: Practi-</i> | ing a voice to the speakers . In <i>Proceedings of the</i> | 973 |
| 917 | <i>cal applications of machine translation</i> . | 61st Annual Meeting of the Association for Compu- | 974 |
| 918 | Luke Gessler. 2022. Closing the NLP gap: Document- | tational Linguistics (Volume 1: Long Papers) , pages | 975 |
| 919 | ary linguistics and NLP need a shared software in- | 4871–4897, Toronto, Canada. Association for Com- | 976 |
| 920 | frastructure . In <i>Proceedings of the Fifth Workshop</i> | <i>putational Linguistics</i> . | 977 |
| 921 | on the Use of Computational Methods in the Study | Manuel Mager, Arturo Oncevay, Abteen Ebrahimi, John | 978 |
| 922 | of Endangered Languages , pages 119–126, Dublin, | Ortega, Annette Rios, Angela Fan, Ximena Gutierrez- | 979 |
| 923 | Ireland. Association for Computational Linguistics. | Vasques, Luis Chiruzzo, Gustavo Giménez-Lugo, Ri- | 980 |
| 924 | Marjan Ghazvininejad, Hila Gonen, and Luke Zettle- | cardo Ramos, et al. 2021. Findings of the Ameri- | 981 |
| 925 | moyer. 2023. Dictionary-based phrase-level prompt- | casNLP 2021 shared task on open machine transla- | 982 |
| 926 | ing of large language models for machine translation . | tion for indigenous languages of the Americas . In | 983 |
| 927 | <i>Preprint</i> , arXiv:2302.07856. | Proceedings of the First Workshop on Natural Lan- | 984 |
| 928 | Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, | guage Processing for Indigenous Languages of the | 985 |
| 929 | Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, | Americas , pages 202–217, Online. Association for | 986 |
| 930 | Young Jin Kim, Mohamed Afify, and Hany Has- | <i>Computational Linguistics</i> . | 987 |
| 931 | san Awadalla. 2023. How good are gpt models at | Gabriel Nicholas and Aliya Bhatia. 2023. Lost in trans- | 988 |
| 932 | machine translation? a comprehensive evaluation . | lation: Large language models in non-english content | 989 |
| 933 | <i>Preprint</i> , arXiv:2302.09210. | analysis . <i>Preprint</i> , arXiv:2306.07377. | 990 |
| 934 | Jason R Jolley and Luciane Maimone. 2022. Thirty | John E Ortega, Richard Castro Mamani, and Kyunghyun | 991 |
| 935 | years of machine translation in language teaching | Cho. 2020. Neural machine translation with a | 992 |
| 936 | and learning: A review of the literature. <i>L2 Journal</i> , | polysynthetic low resource language . <i>Machine Trans-</i> | 993 |
| 937 | 14(1):26–44. | <i>lation</i> , 34(4):325–346. | 994 |
| 938 | Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, | Miguel Ángel Pinto Tapia, Luis Quispe Zúñiga, et al. | 995 |
| 939 | Ondřej Bojar, Anton Dvorkovich, Christian Feder- | 2005. Didáctica quechua i. Documento de trabajo, | 996 |
| 940 | mann, Mark Fishel, Markus Freitag, Thamme Gowda, | Dirección Regional de Educación Apurímac Di- | 997 |
| 941 | Roman Grundkiewicz, Barry Haddow, et al. 2023. | rección Gestión Pedagógica. Elaborado por Mgt. | 998 |
| 942 | Findings of the 2023 conference on machine transla- | Miguel Ángel Pinto Tapia, Esp.DREA EBI-RURAL- | 999 |
| 943 | tion (WMT23): LLMs are here but not quite there yet . | HUASCARÁN. Contact: pitman_38@hotmail.com, | 1000 |
| 944 | In <i>Proceedings of the Eighth Conference on Machine</i> | mapt_38@yahoo.es. | 1001 |
| 945 | <i>Translation</i> , pages 1–42, Singapore. Association for | Qheswa Simi Hamut'ana Kurak Suntur. 2005. <i>Dic-</i> | 1002 |
| 946 | <i>Computational Linguistics</i> . | <i>cionario Quechua - Español - Quechua Qheswa -</i> | 1003 |
| 947 | Moshe Koppel and Noam Ordan. 2011. Translationese | <i>Español - Qheswa Simi Taqe</i> , 2 edition. Multiservi- | 1004 |
| 948 | and its dialects . In <i>Proceedings of the 49th An-</i> | cios e Imprenta Edmundo Pantigozo EIRL, Cusco, | 1005 |
| 949 | <i>annual Meeting of the Association for Computational</i> | Peru. | 1006 |
| 950 | <i>Linguistics: Human Language Technologies</i> , pages | Machel Reid, Nikolay Savinov, Denis Teplyashin, | 1007 |
| 951 | 1318–1326, Portland, Oregon, USA. Association for | Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste | 1008 |
| 952 | <i>Computational Linguistics</i> . | Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Fi- | 1009 |
| 953 | Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paran- | rat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Un- | 1010 |
| 954 | jape, Michele Bevilacqua, Fabio Petroni, and Percy | locking multimodal understanding across millions of | 1011 |
| 955 | Liang. 2024. Lost in the Middle: How Language | tokens of context . <i>arXiv preprint arXiv:2403.05530</i> . | 1012 |
| 956 | Models Use Long Contexts . <i>Transactions of the Asso-</i> | | |
| 957 | <i>ciation for Computational Linguistics</i> , 12:157–173. | | |

| | | | |
|------|---|--|------|
| 1013 | Katherine J. Riestenberg, Ally Freemond, | Atnafu Lambebo Tonja, Fazlourrahman Balouchzahi, | 1066 |
| 1014 | Brook Danielle Lillehaugen, and Jonathan N. | Sabur Butt, Olga Kolesnikova, Hector Ceballos, | 1067 |
| 1015 | Washington. 2024. Prioritizing Community Partners’ | Alexander Gelbukh, and Thamar Solorio. 2024. Nlp | 1068 |
| 1016 | Goals in Projects to Support Indigenous Language | progress in indigenous latin american languages. | 1069 |
| 1017 | Revitalization. In <i>Decolonizing Linguistics</i> . Oxford | <i>Preprint</i> , arXiv:2404.05365. | 1070 |
| 1018 | University Press. | | |
| 1019 | Annette Rios. 2015. <i>A basic language technology</i> | Boshi Wang, Xiang Deng, and Huan Sun. 2022. Itera- | 1071 |
| 1020 | <i>toolkit for quechua</i> . Ph.D. thesis, University of | tively prompt pre-trained language models for chain | 1072 |
| 1021 | Zurich. | of thought. <i>Preprint</i> , arXiv:2203.08383. | 1073 |
| 1022 | Annette Rios and Richard Castro Mamani. 2014. Mor- | Rodolfo Zevallos, John Ortega, William Chen, Richard | 1074 |
| 1023 | phological disambiguation and text normalization for | Castro, Núria Bel, Cesar Toshio, Renzo Venturas, | 1075 |
| 1024 | southern quechua varieties. | Hilario Aradiel, and Nelsi Melgarejo. 2022. Intro- | 1076 |
| 1025 | Nathaniel Robinson, Perez Ogayo, David R. Mortensen, | ducing QuBERT: A large monolingual corpus and | 1077 |
| 1026 | and Graham Neubig. 2023. ChatGPT MT: Competi- | BERT model for Southern Quechua. In <i>Proceedings</i> | 1078 |
| 1027 | tive for high- (but not low-) resource languages. In | <i>of the Third Workshop on Deep Learning for Low-</i> | 1079 |
| 1028 | <i>Proceedings of the Eighth Conference on Machine</i> | <i>Resource Natural Language Processing</i> , pages 1–13, | 1080 |
| 1029 | <i>Translation</i> , pages 392–418, Singapore. Association | Hybrid. Association for Computational Linguistics. | 1081 |
| 1030 | for Computational Linguistics. | Chen Zhang, Xiao Liu, Juheng Lin, and Yansong Feng. | 1082 |
| 1031 | Ohad Rubin, Jonathan Herzig, and Jonathan Berant. | 2024. Teaching large language models an unseen | 1083 |
| 1032 | 2022. Learning to retrieve prompts for in-context | language on the fly. <i>Preprint</i> , arXiv:2402.19167. | 1084 |
| 1033 | learning. In <i>Proceedings of the 2022 Conference of</i> | | |
| 1034 | <i>the North American Chapter of the Association for</i> | | |
| 1035 | <i>Computational Linguistics: Human Language Tech-</i> | | |
| 1036 | <i>nologies</i> , pages 2655–2671, Seattle, United States. | | |
| 1037 | Association for Computational Linguistics. | | |
| 1038 | Sylvia L.R. Schreiner, Lane Schwartz, Benjamin Hunt, | | |
| 1039 | and Emily Chen. Multidirectional leveraging for | | |
| 1040 | computational morphology and language documenta- | | |
| 1041 | tion and revitalization. <i>Language documentation</i> | | |
| 1042 | <i>and conservation</i> , 14. | | |
| 1043 | Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. | | |
| 1044 | 2020. Bleurt: Learning robust metrics for text gener- | | |
| 1045 | ation. <i>Preprint</i> , arXiv:2004.04696. | | |
| 1046 | Evan Shieh, Faye-Marie Vassel, Cassidy Sugimoto, and | | |
| 1047 | Thema Monroe-White. 2024. Laissez-faire harms: | | |
| 1048 | Algorithmic biases in generative language models. | | |
| 1049 | <i>Preprint</i> , arXiv:2404.07475. | | |
| 1050 | David Stap and Ali Araabi. 2023. ChatGPT is not a | | |
| 1051 | good indigenous translator. In <i>Proceedings of the</i> | | |
| 1052 | <i>Workshop on Natural Language Processing for In-</i> | | |
| 1053 | <i>digenuous Languages of the Americas (AmericasNLP)</i> , | | |
| 1054 | pages 163–167, Toronto, Canada. Association for | | |
| 1055 | Computational Linguistics. | | |
| 1056 | Garrett Tanzer, Mirac Suzgun, Eline Visser, Dan Juraf- | | |
| 1057 | sky, and Luke Melas-Kyriazi. 2024. A benchmark | | |
| 1058 | for learning to translate a new language from one | | |
| 1059 | grammar book. <i>Preprint</i> , arXiv:2309.16575. | | |
| 1060 | Jörg Tiedemann. 2012. Parallel data, tools and inter- | | |
| 1061 | faces in OPUS. In <i>Proceedings of the Eighth In-</i> | | |
| 1062 | <i>ternational Conference on Language Resources and</i> | | |
| 1063 | <i>Evaluation (LREC’12)</i> , pages 2214–2218, Istanbul, | | |
| 1064 | Turkey. European Language Resources Association | | |
| 1065 | (ELRA). | | |

A Example errors

| Original | Reference | Translation | Model - Prompt (Error) |
|--|---|--|--|
| Mistranslation | | | |
| ñuqaqa sapa punchaw llankani ayluyta uywanaypaq | yo trabajo todos los días para mantener a mi familia; <i>I work every day to support my family</i> | yo sé que en algún momento volveré a mi comunidad; <i>I know that at some point I will return to my community</i> | GPT-3.5 - BASELINE (Complete Mistranslation) |
| ñuqaqa illani runakunawan rimayta munaptiy | yo viajo porque me gusta hablar con las personas; <i>I travel because I like to talk to people</i> | yo quiero hablar con las personas cuando amanece; <i>I want to talk to people when the sun rises.</i> | GPT-4o C+G+M (Mistranslation - Lexical Correspondence) |
| Fluency | | | |
| qamqa yachayta tukuptiyki, puñunki | cuando termines de estudiar, tu duermes; <i>when you finish studying, you sleep</i> | acabaré de aprender lo que tú duermes; <i>I will finish learning what you sleep</i> | GPT-3.5 - C+G+M (Coherence: unresolvable pronoun) |
| payqa lunesmanta jueveskama kallpan tullu niraq kachkananpaq | él corre de lunes a jueves para estar delgado; <i>he runs from monday to thursday to be slim</i> | él sigue siendo delgado desde el lunes hasta el jueves para estar; <i>He remains thin from Monday to Thursday to be.</i> | GPT-4o - GRAMMAR (Grammar: missing predicate with 'estar') |
| Grammar | | | |
| qam allinta tusunki | tu bailas bien; <i>You dance well</i> | tú bailarás bien; <i>You will dance well</i> | GPT-3.5 - C+G+M (TAM: tense) |
| qamqa mamaykipa wasinta yanuq rinki | tú irás a cocinar a la casa de tu mamá; <i>you will go to cook at your mother's house</i> | voy a cocinar en la casa de tu madre; <i>I'm going to cook at your mother's house</i> | GPT-3.5 - MORPH (Subject: 'you' to 'I') |
| qamqa misita qawaspa qari warata awachkarqanki | tú estabas tejiendo un pantalón de varón mirando al gato; <i>you were knitting a pair of men's pants looking at the cat</i> | tú mirabas al gato mientras el hombre tejía el pantalón; <i>You looked at the cat while the man knitted the pants.</i> | GPT-4o - C+G+M (Subject: 'you' to 'man') |
| qam runasimita yachapayaptiyki, allintam rimanki | si practicas el quechua lo hablarás bien; <i>If you practice Quechua you will speak it well</i> | tú ya sabías el runasimi, hablas muy bien; <i>you already knew runasimi, you speak very well</i> | GPT-4o - MORPH (TAM: missing conditional) |

Table 5: Examples of errors analyzed in Section 5.3, one error per type.

B Unprompted stereotypes

| Source | Reference | Translation | Model - PROMPT TYPE |
|--|--|--|---------------------|
| payqa qatupi ñañañapa rantisqan apita ruraraqan | ella preparó una mazamorra que compró su hermana en el mercado; <i>she prepared a mazamorra [pudding] that her sister bought at the market</i> | el sol brilla en nuestra tierra, calienta y nos da vida; <i>The sun shines on our land, warms us and gives us life.</i> | GPT-3.5 - BASE |
| paykunaqa yachachiqninkupa rakta medicina liwru apasqanta likanku | ellos leerán el libro grueso de medicina que trajo su profesor; <i>They will read the thick medical book that their teacher brought</i> | los médicos de la comunidad nos enseñaron a utilizar plantas medicinales; <i>The doctors from the community taught us how to use medicinal plants.</i> | GPT-3.5 - GRAMMAR |
| qamkunapa taytamamaykichikta watunaykichik, paykunata purikuq apanaykichik | ustedes tienen que visitar a sus padres y llevarlos a pasear; <i>you have to visit your parents and take them for a walk</i> | ustedes deben honrar a sus padres, guiarlos en su camino; <i>You must honor your parents, guide them on their path.</i> | GPT-4o - BASE |
| ñuqaqa upyayta munani | yo quiero beber; <i>I want to drink</i> | yo quiero beber chicha; <i>I want to drink chicha</i> | Gemini - BASE |

Table 6: Examples of errors involving cultural stereotypes (Section 6.1).

C Quality Descriptions

| Quality | Description |
|---------|---|
| High | Output is an accurate and/or acceptable translation of the source content. |
| Med | Output contains errors that prevent it from being an acceptable translation, but is generally high in quality otherwise. |
| Low | Output contains errors that prevent it from being an acceptable translation, with minor correspondences that vaguely identify it as relevant to the source. |
| None | Output does not appear to be relevant to the source. |

Table 7: Quality Descriptions

D Annotation error typology

| Dimension | Error | Description |
|-----------------|---|---|
| Accuracy | Addition | Translation includes information not present in the source, but does not result in the displacement of source content. |
| Accuracy | Omission | Translation is missing content from the source. |
| Accuracy | Substitution - Subject | The translated segment contains content identified as relevant to the source in other spans, but substitutes novel subject markers for those present in the source in the highlighted span; Classify an error as a “substitution” when the error appears to result in both Addition and Omission errors that cannot be distinguished into two distinct spans. |
| Accuracy | Substitution - TAM | The translated segment contains content identified as relevant to the source in other spans, but substitutes novel TAM for those present in the source in the highlighted span; Classify an error as a “substitution” when the error appears to result in both Addition and Omission errors that cannot be distinguished into two distinct spans. |
| Accuracy | Substitution - Other | Substitution errors that do not involve mistranslated subject markers or TAM. See above. |
| Accuracy | Overtranslation | Error occurring in the target content that is inappropriately more specific than the source content. |
| Accuracy | Undertranslation | Error occurring in the target content that is inappropriately less specific than the source content. |
| Target Error | Grammar | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is not grammatical in the target language. |
| Target Error | Coherence | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is unnatural or incoherent in the target language. |
| Target Error | Style/Register | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is produced in a style or register that is inappropriate given the content. |
| Non-Translation | Complete Mistranslation | The entire segment is coherent in the target language but the core predicate shows no immediate connection to the reference translation. |
| Non-Translation | Mistranslation - Lexical Correspondence | The entire segment is coherent in the target language but only minor correspondences to the reference translation may be identified. |
| Non-Translation | Refusal | Model does not attempt to translate into the target language, e.g., because it "does not understand". |
| Model error | Garbled | Output does not contain coherent text in the target language. |
| Model error | ChattyGPT | Output contains translated content, but is wordy, over-explanatory, and/or abruptly truncated. |

Table 8: Adapted MQM typology for human error annotation

| | | | |
|------|--|--|------|
| 1085 | E Example Prompts | | |
| 1086 | The following are examples of prompts generated | Llaqta-ta risaq Iré al pueblo | 1136 |
| 1087 | used automated retrieval from the database. | Hamawt'anchis Punuta rinqa | 1137 |
| 1088 | | Llanta umalliq llaqtata richkan | 1138 |
| 1089 | BASELINE | nki: FLEXIÓN DE TIEMPO. TIEMPO FUTURO. | 1139 |
| 1090 | | TIEMPO FUTURO. Los sufijos para cada una | 1140 |
| 1091 | [TAREA] Traduce la siguiente frase del quechua al | de las personas gramaticales son: saq, nki, nqa, | 1141 |
| 1092 | español. Responde sólo con la traducción: | sun, saqku, nkichis, nqaku; en singular y plural | 1142 |
| 1093 | quechua: qam allinta tusunki | respectivamente. | 1143 |
| 1094 | español: | Ejemplos: | 1144 |
| 1095 | | Puklla-saq jugaré | 1145 |
| 1096 | MORPHS-ONLY | Puklla-nki jugarás | 1146 |
| 1097 | | Puklla-nqa jugará | 1147 |
| 1098 | [CONTEXTO] | Puklla-sun jugaremos | 1148 |
| 1099 | qam: [PrnPers+2.Sg] | Puklla-saqku jugaremos | 1149 |
| 1100 | allin: bueno [DB][NRroot] | Puklla-nkichis Uds. jugarán | 1150 |
| 1101 | ta: [+Acc][Cas] | Puklla-nqaku ellos jugarán | 1151 |
| 1102 | tusu: bailar [VRoot][DB] | | 1152 |
| 1103 | nki: [+2.Sg.Subj][VPers] | [TAREA] Traduce la siguiente frase . . . | 1153 |
| 1104 | allin. adj. Bueno (término de aprobación). SINÓN: | CORPUS-ONLY | 1154 |
| 1105 | kusa. EJEM: allin p'unchay, buenos días: allin tuta, | [CONTEXTO] | 1155 |
| 1106 | buenas noches; allin tutamanta, buena mañana, | quechua: rimanakunapaq wawakunapa rimasqan | 1158 |
| 1107 | buenos días; allin inti chinkay, buenas tardes; | simi aswan allinta takyachinaraq piwanpas | 1159 |
| 1108 | allin ñiyniyoy, de buena fe, fiel, justo, íntegro: | maywanpas mana manchakuspa rimananpaq | 1160 |
| 1109 | allin nunayoy, de espíritu bueno; allin puriq, | chaymi qillqanapaqpas ñawichanapaqpas aswan | 1161 |
| 1110 | de comportamiento bueno; allin puriy, compor- | allin kanqa | 1162 |
| 1111 | tamiento bueno; allin rikuy, tratamiento bueno; | español: para este diálogo saber la lengua que | 1163 |
| 1112 | allin rikuq, el que trata bien; allin ruway, obrar | dominan los niños sería importante para que ellos | 1164 |
| 1113 | bien, beneficiar; lo que se hace bien, beneficioso; | se expresen sin miedo de ahí será que la escritura y | 1165 |
| 1114 | allin ruwaq, el que hace bien; allin yuyay, pensar | la lectura salga de manera óptima | 1166 |
| 1115 | bien; pensamiento bueno; allin qolqeyoy, poseedor | quechua: kay tiqsiپی sumaq rimanakunapaqa | 1167 |
| 1116 | de plata fina; adinerado. | kawsayninchikmi allinta kallpachawanchik runaku- | 1168 |
| 1117 | ta. s. Gram. Sufijo que desempeña los papeles de | nahina allinta tiyanapaq chaymi ñuqanchikkqa | 1169 |
| 1118 | artículo y preposición. EJEM: llamata qatiy, arrea | allinta ñawichayta qillqayta yachananchik ñawpa | 1170 |
| 1119 | la llama; Urkusmanta hamuni, vengo de Urcos. | ayllunchikkuna rurasqankuta maytukunapi tukuy | 1171 |
| 1120 | | puyñukunapi tiqsi muyu qhawarisqankuta | 1172 |
| 1121 | [TAREA] Traduce la siguiente frase . . . | español: para vivir en armonía tenemos que | 1173 |
| 1122 | | conocer bien nuestra forma de vivir y luego | 1174 |
| 1123 | GRAMMAR-ONLY | escribir leer tambien a valorar lo que nos dejaron | 1175 |
| 1124 | | nuestros antecesores en cada visión sobre el mundo | 1176 |
| 1125 | [CONTEXTO] | quechua: winsislawcha chayarqamuptinsi tu- | 1177 |
| 1126 | ta: CASO ACUSATIVO. Su marca es –ta, esta es | parquspanku allinta qatunakusqanku suwakuypi | 1178 |
| 1127 | una marca de objeto directo con los verbos que no | purinankupaq | 1179 |
| 1128 | son de movimiento (quietud). Ejemplo: | español: cuando había llegado wenseslau y a su | 1180 |
| 1129 | Quyllur–ta qhawani Veo una estrella | encuentro se habían reforzarón para andar a robar | 1181 |
| 1130 | T'anta–ta apay Lleva pan | | 1182 |
| 1131 | Ñuqa quylluyta qhawani | [TAREA] Traduce la siguiente frase . . . | 1183 |
| 1132 | Pedrucha t'antata rantin | | 1184 |
| 1133 | En cambio con los verbos de movimiento –ta | | |
| 1134 | indica (hacia) que es igual a meta. Ejemplos: | | |
| 1135 | Punu–ta rini Voy a Puno | | |

F Full quality scores

This section contains the full tables of BLEURT and human-annotated quality scores. Table 9 contains the full results summarized in Tables 1 and 3 of the main text. Table 10 and Table 11 contain the full scores summarized in Table 3.

| Prompt | GPT-3.5 | | GPT-4o | | Gemini-1.5 | | Llama 3 | |
|----------------------|---------|--------|--------|--------|------------|--------|---------|--------|
| | auto | manual | auto | manual | auto | manual | auto | manual |
| BASELINE | 0.19 | 0.22 | 0.66 | 0.66 | 0.56 | 0.57 | 0.15 | 0.16 |
| CORPUS-ONLY | 0.27 | 0.29 | 0.59 | 0.61 | 0.49 | 0.47 | 0.19 | 0.18 |
| GRAMMAR-ONLY | 0.23 | 0.24 | 0.56 | 0.58 | 0.55 | 0.54 | 0.17 | 0.15 |
| MORPH-ONLY | 0.44 | 0.56 | 0.54 | 0.63 | 0.61 | 0.66 | 0.39 | 0.49 |
| CORPUS-GRAMMAR | 0.26 | 0.28 | 0.59 | 0.59 | 0.54 | 0.53 | 0.21 | 0.21 |
| CORPUS-MORPH | 0.44 | 0.52 | 0.59 | 0.64 | 0.59 | 0.64 | 0.36 | 0.38 |
| GRAMMAR-MORPH | 0.41 | 0.54 | 0.53 | 0.61 | 0.61 | 0.64 | 0.39 | 0.37 |
| CORPUS-GRAMMAR-MORPH | 0.43 | 0.54 | 0.57 | 0.63 | 0.61 | 0.63 | 0.15 | 0.26 |

Table 9: BLEURT scores for all LLMs and prompt types.

GPT-3.5 Turbo

| | None | Low | Med | High |
|-----------------------|------|-----|-----|------|
| BASELINE | 31 | 17 | 2 | 0 |
| CORPUS-ONLY | 18 | 23 | 8 | 1 |
| GRAMMAR-ONLY | 20 | 27 | 2 | 1 |
| MORPHS-ONLY | 3 | 22 | 16 | 9 |
| CORPUS-GRAMMAR | 18 | 23 | 9 | 0 |
| CORPUS-MORPHS | 2 | 28 | 12 | 8 |
| GRAMMAR-MORPHS | 3 | 29 | 13 | 5 |
| CORPUS-GRAMMAR-MORPHS | 2 | 27 | 12 | 9 |

Table 10: Human quality annotation of GPT-3.5 outputs with automated retrieval (raw counts out of 50) by prompt type.

GPT-4o

| | None | Low | Med | High |
|-----------------------|------|-----|-----|------|
| BASELINE | 0 | 10 | 20 | 20 |
| CORPUS-ONLY | 1 | 16 | 13 | 20 |
| GRAMMAR-ONLY | 0 | 17 | 16 | 17 |
| MORPHS-ONLY | 0 | 13 | 18 | 19 |
| CORPUS-GRAMMAR | 0 | 14 | 17 | 19 |
| CORPUS-MORPHS | 0 | 10 | 17 | 23 |
| GRAMMAR-MORPHS | 0 | 19 | 14 | 17 |
| CORPUS-GRAMMAR-MORPHS | 0 | 9 | 20 | 21 |

Table 11: Human quality annotation of GPT-4o outputs with automated retrieval (raw counts out of 50) by prompt type.

G Full error counts

This section contains the full counts of annotated errors by category and prompt type.

| GPT-3.5 Turbo | | | | | | | | | |
|---|------|----|----|-----|-----|-----|-----|-------|-------|
| | BASE | C | G | M | C+G | C+M | G+M | C+G+M | TOTAL |
| None | 0 | 1 | 1 | 6 | 0 | 8 | 3 | 5 | 24 |
| Addition | 0 | 5 | 3 | 14 | 1 | 9 | 10 | 11 | 53 |
| Omission | 3 | 9 | 2 | 13 | 2 | 5 | 9 | 9 | 52 |
| Substitution - Subject | 0 | 3 | 0 | 7 | 0 | 9 | 9 | 12 | 40 |
| Substitution - TAM | 0 | 11 | 3 | 17 | 6 | 19 | 19 | 19 | 94 |
| Substitution - Other | 4 | 9 | 4 | 13 | 6 | 16 | 14 | 13 | 79 |
| Overtranslation | 1 | 1 | 1 | 4 | 0 | 2 | 3 | 2 | 14 |
| Undertranslation | 0 | 0 | 0 | 2 | 1 | 2 | 2 | 2 | 9 |
| Target Error - Grammar | 0 | 1 | 1 | 4 | 2 | 3 | 3 | 1 | 15 |
| Target Error - Coherence | 0 | 0 | 3 | 5 | 2 | 3 | 7 | 7 | 27 |
| Target Error - Style/Register | 0 | 3 | 0 | 5 | 2 | 3 | 1 | 2 | 16 |
| Complete Mistranslation | 30 | 19 | 21 | 2 | 18 | 2 | 2 | 2 | 96 |
| Mistranslation - Lexical Correspondence | 15 | 13 | 23 | 9 | 21 | 11 | 15 | 10 | 117 |
| Refusal | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Total | 54 | 75 | 62 | 101 | 61 | 92 | 97 | 95 | 637 |

Table 12: Human error type annotation of GPT-3.5 outputs with automated retrieval (raw counts, up to 3 errors per sentence) by prompt type.

| GPT-4o | | | | | | | | | |
|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|
| | BASE | C | G | M | C+G | C+M | G+M | C+G+M | TOTAL |
| None | 15 | 16 | 10 | 16 | 13 | 19 | 14 | 18 | 121 |
| Addition | 2 | 5 | 7 | 5 | 4 | 1 | 6 | 4 | 34 |
| Omission | 8 | 7 | 6 | 7 | 6 | 3 | 5 | 5 | 47 |
| Substitution - Subject | 1 | 2 | 0 | 1 | 2 | 1 | 2 | 2 | 11 |
| Substitution - Other | 22 | 24 | 22 | 18 | 19 | 18 | 17 | 20 | 160 |
| Substitution - TAM | 16 | 17 | 19 | 12 | 13 | 10 | 11 | 9 | 107 |
| Overtranslation | 2 | 1 | 0 | 2 | 2 | 2 | 1 | 2 | 12 |
| Undertranslation | 6 | 1 | 3 | 1 | 3 | 0 | 1 | 2 | 17 |
| Target Error - Grammar | 1 | 3 | 4 | 4 | 1 | 2 | 6 | 1 | 22 |
| Target Error - Coherence | 1 | 3 | 4 | 5 | 4 | 5 | 9 | 5 | 36 |
| Target Error - Style/Register | 1 | 2 | 3 | 4 | 4 | 2 | 4 | 3 | 23 |
| Complete Mistranslation | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Mistranslation - Lexical Correspondence | 4 | 3 | 5 | 6 | 6 | 6 | 9 | 4 | 43 |
| Total | 79 | 85 | 83 | 81 | 77 | 69 | 85 | 75 | 634 |

Table 13: Human error type annotation of GPT-4o outputs with automated retrieval (raw counts, up to 3 errors per sentence) by prompt type.