

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 EMMA-500: ENHANCING MASSIVELY MULTILINGUAL ADAPTATION OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we introduce **EMMA-500**, a large-scale multilingual language model continue-trained on texts across 546 languages designed for enhanced multilingual performance, with a focus on improving language coverage for low-resource languages. To facilitate continual pre-training, we compile the **MaLA corpus**, a comprehensive multilingual dataset and enrich it with curated datasets across diverse domains. Leveraging this corpus, we conduct extensive continual pre-training of the Llama 2 7B model, resulting in EMMA-500, which demonstrates robust performance across a wide collection of benchmarks, including a comprehensive set of multilingual tasks and **PolyWrite**, an open-ended generation benchmark developed in this study. Our results highlight the effectiveness of continual pre-training in expanding large language models' language capacity, particularly for underrepresented languages, demonstrating significant gains in cross-lingual transfer, task generalization, and language adaptability.

1 INTRODUCTION

Multilingual language models (MLMs) are designed to process and generate text in multiple languages. These models have evolved rapidly over the past decade, fueled by advances in deep learning, e.g., Transformer networks (Vaswani et al., 2017), pre-training techniques, and the availability of large-scale multilingual corpora such as mC4 (Raffel et al., 2020) and ROOTS (Laurençon et al., 2022). The development of models like BERT (Devlin et al., 2019), GPT, and T5 (Raffel et al., 2020) opened the door for multilingual counterparts such as mBERT, XLM-R (Conneau et al., 2020), mGPT (Shliazhko et al., 2022), and mT5 (Xue et al., 2021). These models were trained on massive multilingual corpora, allowing text in dozens of languages to be processed with the same set of model weights. Multilingual language models have shown impressive performance across various tasks like text classification and machine translation by leveraging cross-lingual transfer from high-resource languages such as English and Chinese. However, many low-resource languages, with limited available data, remain underrepresented. While large corpora are abundant for high-resource languages like English, French, and Spanish, languages like Xhosa and Inuktitut suffer from scarce or fragmented data, leading to imbalanced training and models that prioritize high-resource languages.

Recent studies adopt continual pre-training to enhance the language coverage of large language models on low-resource languages. For example, Glot500 (Imani et al., 2023) and MaLA-500 (Lin et al., 2024) use continual pre-training and vocabulary extension using XLM-R and LLaMA, respectively, on the Glot500-c corpus covering 534 languages. xLLMs-100 (Lai et al., 2024) proceeds to multilingual instruction fine-tuning to improve the multilingual performance of LLaMA and BLOOM models on 100 languages, and Aya model (Üstün et al., 2024) applies continual training to the mT5 model (Xue et al., 2021) using their constructed instruction dataset. LLaMAX (Lu et al., 2024) pushes the envelope by focusing on translation tasks via continual pre-training of LLaMA in over 100 languages.

As the field of MLMs evolves, the role of data becomes increasingly critical in enhancing the performance of the models, particularly when it comes to low-resource languages. To address this need for more and better data, we extend existing work, such as MaLA-500, by expanding the corpus for continual pre-training, coupled with large-scale training methods. We emphasize the creation of a massively multilingual corpus that not only increases the quantity of data but also diversifies the types of texts (e.g., code, books, scientific papers, and instructions). This ensures better language adaptation and broader language coverage, thus improving the representation and performance of

054 multilingual language models, especially for underrepresented languages, ultimately creating more
 055 inclusive and versatile language models that cater to a broader linguistic diversity. Our contribution is
 056 summarized as follows: (1) We compile a massively multilingual corpus named MaLA to facilitate
 057 continual training of large language models for enhanced language adaptation across a wide range
 058 of linguistic contexts. (2) We extend the MaLA corpus by integrating multiple curated datasets,
 059 creating a comprehensive and diverse data mix specifically for continual pre-training. (3) We perform
 060 continual pre-training with the Llama 2 7B model¹ on the multilingual corpus with 546 languages,
 061 resulting in the new EMMA-500 model. This model is rigorously evaluated on a diverse set of tasks
 062 and benchmarks, including PolyWrite, a novel open-ended generation benchmark developed as part
 063 of this work.

064 **The MaLA Corpus** Our multilingual corpus features the following characteristics:
 065

- 066 • It contains 939 languages, 546 of which have more than 100k tokens and are used for
 067 training our EMMA-500 model, and 74 billion (B) whitespace delimited tokens in total.
 068
- 069 • It has more than 300 languages with over 1 million whitespace delimited tokens and 546
 070 languages with over 100k tokens.
 071
- 072 • It comes with four publicly available versions: (1) a noisy version after only basic pre-
 073 processing like extraction and harmonization; (2) a cleaned version after data cleaning; (3)
 074 a deduplicated version after approximate and exact deduplication; (4) a split version with
 075 train and valid splits.
 076
- 077 • Our augmentation to the MaLA corpus includes different types of texts such as code, books,
 078 scientific papers and instruction data, leading to a data mix with 100B+ whitespace delimited
 079 tokens.
 080

081 **Evaluation Results** In comparison to other decoder-only LLMs with parameter sizes from 4.5B to
 082 13B, including Llama 2-based multilingual models and the latest advanced models, our EMMA-500
 083 7B parameter model demonstrates strong performance. It achieves the lowest negative log-likelihood
 084 in intrinsic evaluation and significantly improves commonsense reasoning, machine translation, and
 085 open-ended generation tasks. It also outperforms all Llama 2-based models and multilingual LLMs in
 086 text classification and natural language inference. While math and machine reading comprehension
 087 remain challenging, our model still enhances the Llama 2 base model, and it surpasses the base model
 088 in code generation without regression in performance.

089 2 THE MALA CORPUS

090 The MaLA corpus—MaLA standing for **M**assive **L**anguage **A**daptation—is a diverse and extensive
 091 compilation of text data encompassing 939 languages sourced from a wide array of datasets. It is
 092 developed for continual training of multilingual large language models. The source datasets the
 093 MaLA corpus is compiled from exhibit a wide range of variance in various aspects. Examples of such
 094 elements include the data quality, the nature of the text content, how the data sources were organized
 095 into directory and file tree structures (if distributed as files rather than through an API) and the naming
 096 conventions used therein, the data formats and structures in the files or in-memory objects containing
 097 the text data, and the logic by which multilingual texts were aligned. This section introduces the
 098 efforts made in data extraction, harmonization, pre-processing, cleaning, and deduplication in order
 099 to build the corpus.

100 2.1 DATA PRE-PROCESSING

101 To develop the MaLA corpus for training our language model, we establish a processing workflow
 102 consisting of the following key steps: (1) loading and curating identified data sources, (2) extracting
 103 and harmonizing both textual and metadata from these diverse sources into a unified format—often
 104 with tailored filtering, and (3) performing deduplication and further filtering on the textual data. In
 105 step (1), source data are either organized and loaded into memory from a file tree structure or, if
 106

107 ¹Choosing Llama 2 allows us to compare our model with many other models derived from it using continual
 108 pre-training. We plan to continue training models based on Llama 3/3.1 in the future.

108 available, accessed through an API. In step (2), the output is designed to be JSON Lines (JSONL)
 109 files with extracted text data and other desired content. A JSONL file contains multiple JSON records
 110 for storing data, each separated by a newline character. We selectively process only data annotated as
 111 training or development (validation) data, deliberately excluding test data.
 112

113 2.1.1 EXTRACTION AND HARMONIZATION

114
 115 As mentioned, there is significant variability in the data quality of the source datasets used for
 116 compiling the corpus. Many of the datasets exhibit data quality challenges that would have adverse
 117 effects on model training if left unaddressed. These dataset-specific challenges are typically addressed
 118 during the pre-processing stage. For example, we identify one issue involving text records consisting
 119 solely of date and timestamp information in the dataset for Languages of Russia (Corpora and Tools,
 120 n.d.), likely resulting from a web scraper failing to differentiate between these elements and actual
 121 text. We address this by implementing a logic in the pre-processing script to detect and exclude such
 122 records. This issue is resolved by introducing a rule to identify and discard text containing consecutive
 123 repeating words. Due to the extensive volume of data, exhaustive examination of every source for
 124 data quality issues is impractical. Instead, we address issues only as we encounter them in our data
 125 exploration and pre-processing pipeline development efforts. This approach likely leaves some data
 126 quality problems undetected, since we do not go through the data systematically. Despite the need
 127 for customized handling of certain dataset-specific idiosyncrasies, the core logic and structure of
 128 the pre-processing workflow remain consistent across most datasets. We develop a standardized
 129 pre-processing script that can be adapted with minor adjustments to accommodate different datasets.
 130

131 2.1.2 LANGUAGE CODE NORMALIZATION

132 An essential component of our pre-processing pipeline involved converting language codes to the
 133 ISO 639-3 standard. This is crucial for ensuring consistent language identification across the source
 134 datasets. We rely on the declared language of each dataset and normalize it to the ISO standard without
 135 performing additional language identification at this step. This approach helps maintain uniformity
 136 while streamlining the pre-processing workflow. We primarily use the PyPI library `iso639-lang`²
 137 or `langcodes`³. While converting language codes to ISO 639-3, we encounter several challenges.
 138 One issue is that some languages in ISO 639-3 are divided into multiple subvarieties, but our source
 139 data does not specify which subvariety is present. Our solution is to retain the original language code
 140 from the dataset, even when it does not conform to the ISO 639-3 standard. Another issue arises
 141 when certain languages are merged into other languages in the ISO 639-3 standard. In these cases,
 142 we update the language code to reflect the merged language. Additionally, some language names or
 143 codes in the source data—referred to as “original language names” or “original language codes”—are
 144 not recognized by the conversion libraries. In some cases, the reason behind this is that the original
 145 code in fact represents language families or groups of dialects (e.g., the ISO 639-2 codes “ber” for
 146 Berber languages and “bih” for Bihari languages), rather than specific languages. If so, we then retain
 147 the original codes, despite their non-compliance with ISO 639-3. In other cases, the original language
 148 names are spelled differently from the standard recognized by the libraries. To address this mismatch,
 149 we implement a logic to detect and correct misspelled language names during pre-processing. All
 150 these “corner cases” require careful attention in the pre-processing stage to ensure correct language
 151 code identification.

152 2.1.3 WRITING SYSTEM RECOGNITION

153 In addition to normalizing language codes, we also identify the script or writing system used in the
 154 text data. We use the GlotScript library (Kargaran et al., 2024) to recognize writing systems
 155 accordant to the ISO 15924 standard. The process begins by sampling 100 random lines from each
 156 dataset (or the full dataset if it contains fewer than 100 lines). If GlotScript fails to identify a
 157 script from this sample, we attempt identification using just the first line of the sample. If this still
 158 does not yield a result, we set the script as “None”. It is worth noting that we choose not to classify a
 159 dataset into multiple scripts, even when code-mixing (i.e., the use of multiple scripts) is present.
 160

²<https://pypi.org/project/iso639-lang/>

³<https://pypi.org/project/langcodes/>

If script identification is unsuccessful after the initial steps, we assume the script matched a previously detected one for that language. In cases where no previous script information exists, we refer to a mapping of languages and their default scripts provided by the Glot500 corpora collection. Through this multi-step process, we are able to determine the script for every dataset without exceptions.

During script identification, we encounter several challenges. One issue is determining an appropriate length for the text chunk used for script recognition. A chunk that is too short could lead to incorrect identification if the text contains quotes or foreign language fragments using a different writing system. Conversely, using a chunk that is too long could result in excessive resource usage, slowing down processing or even causing memory exhaustion. Another consideration is whether to assume that a single file or dataset might contain multiple scripts. Such an assumption would require identifying the script at a more granular level, such as paragraph by paragraph or even sentence by sentence. Alternatively, we could assume that each file or dataset contained only one “main” script. This assumption would allow us to identify the script from a representative sample of the text for the whole file or dataset. We adopt the latter approach, recognizing a single dominant script for each dataset. The output of this process is a label in the format `language_Script`, e.g., `eng_Latn`, where “`Language`” represents the ISO 639-3 language code and “`Script`” represents the ISO 15924 script code.

2.2 DATA CLEANING

Most source data has already undergone data cleaning to different extents. Nonetheless, different cleaning processes have been adopted. We continue to clean the data to ensure consistency and accuracy for monolingual and bilingual texts. Following the pipeline used by BigScience’s pipeline for ROOTS corpus (Laurençon et al., 2022), we further adopt some necessary data cleaning to filter out text samples that might have undesirable quality. We first perform document modification for monolingual texts. The first step is whitespace standardization: all types of whitespace in a document are converted into a single, consistent space character. We split documents by newline characters, tabs, and spaces, strip words, and reconstruct the documents to remove very long words. However, these two steps do not apply to languages without whitespace word delimiters like Chinese, Japanese, Korean, Thai, Lao, Burmese, etc. We also remove words containing certain patterns, e.g., “`http`” and “`.com`”, which are likely to be links and page source code. We then perform document filtering, including word count filtering, character repetition filtering, word repetition filtering, special characters filtering, stop words filtering, and flag words filtering. We re-identify the languages that are supported by the pre-trained `fastText`-based language identification model (Joulin et al., 2016b;a). For other languages, we assume the language identification of the original data source and language code conversion are reliable.

As we collect data from different sources, we deduplicate the data to remove the overlap between different sources using MinHash and exact deduplication (Mou et al., 2023), with details described in Appendix A.3.

2.3 KEY STATISTICS

This section presents the final MaLA corpus obtained after data sourcing, pre-processing, cleaning, and deduplication. Table 1 first shows some basic data statistics and compares them with other multilingual corpora for pre-training language models or language adaptation. Additional statistics are presented in Appendix B in the appendix. The MaLA corpus harvests a wide range of datasets in multiple domains. Table 7 in the appendix lists the corpora and collections we used as monolingual data sources. The token counts are based on white-space delimitation, though it might not be accurate for languages like Chinese, Japanese and Korean since the entire clause is counted as one token. Glot500-c (Imani et al., 2023) has 534 languages in total, in which 454 languages are directly distributed on Huggingface ⁴. We also omit high-resource languages in the other three datasets, i.e., MADLAD (Kudugunta et al., 2024), CulturaX (Nguyen et al., 2023), and CC100 (Wenzek et al., 2020), as our main focus is continual pre-training for language adaptation. The final MaLA corpus consists of 939 languages, 546 of which have more than 100k tokens and are used for training our EMMA-500 model. Counting languages with more than 1 million tokens, the MaLA corpus and Glot500-c have more than 300, while MADLAD and CulturaX have 200 and 100 respectively. Compared with Glot500-c, the MaLA corpus contains documents with significantly higher sequence

⁴<https://huggingface.co/datasets/cis-lmu/Glot500>

216 Table 1: Data statistics of the MaLA corpus and comparison to other multilingual corpora. The
 217 number of documents and tokens is in millions.

Dataset	N Lang	N Lang Counted	N Docs	N Tokens	Avg Tokens/Doc
Glot500-c	534	454	1,815	35,449	19.53
MADLAD	419	414	1,043	645,111	618.51
CulturaX	167	161	2,141	1,029,810	480.99
CC100	116	101	2,557	52,201	20.41
MaLA	939	546	824	74,255	90.12

225 lengths, with an average token count of 90 versus 19. This higher sequence length is advantageous for
 226 continually training LLMs because it provides more context within each training example, allowing
 227 the model to better capture long-range dependencies and patterns in the data. As a result, MaLA is
 228 more effective for language adaptation.

230 3 DATA MIXING AND MODEL TRAINING

232 Incorporating a diverse data mix—spanning various languages, domains, document lengths and
 233 styles—is crucial for continual training of large language models to enhance their versatility, generalization
 234 ability, and robustness across a wide range of tasks and domains. We augment the
 235 MaLA corpus with diverse data to mitigate issues such as over-fitting to specific styles or topics or
 236 underperforming on tasks outside the training distribution.

238 3.1 DATA MIXING

240 **Curated Data** We enhance the corpus with high-quality curated data, specifically high-resource
 241 languages in the monolingual part. We use texts from scientific papers as these provide a structured,
 242 information-dense corpus that can improve the model’s ability to handle technical language and
 243 domain-specific content. They are (1) CSL (Li et al., 2022), a large-scale Chinese Scientific Liter-
 244 ature dataset, that contains titles, abstracts, keywords and academic fields of 396,209 papers; (2)
 245 pes2o (Soldaini & Lo, 2023), a collection of full-text open-access academic papers derived from the
 246 Semantic Scholar Open Research Corpus (S2ORC) (Lo et al., 2020). We further add free e-books
 247 from the Gutenberg project⁵ compiled by Faysse (2023). These texts enhance the range of literary
 248 styles and narrative forms, thus enhancing the model’s versatility. Adding high-resource languages
 249 into the pre-training corpora also mitigates the forgetting in model training.

250 **Instruction Data** We further augmented the training corpus by incorporating instruction-based
 251 datasets, inspired by Li et al. (2023); Taylor et al. (2022); Nakamura et al. (2024). We mix two
 252 instruction data into our training corpus. They are: (1) xp3x (Crosslingual Public Pool of Prompts
 253 eXtended)⁶, a multitask instruction collection in 277 languages (Muennighoff et al., 2022); (2)
 254 the Aya collection⁷ that contains both human-curated and machine translated instructions in 101
 255 languages (Singh et al., 2024). For both instruction datasets, we use their training set.

256 **Code** We additionally enrich the training corpus by sourcing code data from The Stack (Kocetkov
 257 et al., 2023). This is done following existing work that demonstrates the value of code data in
 258 improving the reasoning ability of language models (Zhang et al., 2024b; Ma et al., 2024) while also
 259 mitigating any catastrophic forgetting of the base model’s programming knowledge.

261 We subsample The Stack at an effective rate of 15.2%, prioritizing high-quality source files and data
 262 science code⁸. We retain the 32 most important non-data programming languages by prevalence
 263 while also adding in all LLVM code following prior work detailing its importance in multi-lingual
 264 code generation (Szafraniec et al., 2023; Paul et al., 2024). We also source from data-heavy formats

265 ⁵<https://www.gutenberg.org/>

266 ⁶<https://huggingface.co/datasets/CohereForAI/xP3x>

267 ⁷https://huggingface.co/datasets/CohereForAI/aya_collection_language_split

268 ⁸https://huggingface.co/datasets/AlgorithmicResearchGroup/arxiv_research_code

270 Table 2: Data mix for continual training. Code and reasoning-related data are counted by Llama 2
 271 tokenizer and others are counted as whitespace delimited; ‘inst’ stands for instruction and ‘mono’
 272 stands for monolingual texts.

Data	Original Counts	Sample Rate	Final Counts	Percentage
inst high	42,121,055,562	0.1	4,212,105,556	3.08%
inst medium-high+	6,486,592,274	0.2	1,297,318,455	0.95%
inst medium-high	30,651,187,534	0.5	15,325,593,767	11.21%
inst medium	1,444,764,863	1.0	1,444,764,863	1.06%
inst medium-low	47,691,495	5.0	238,457,475	0.17%
inst low	3,064,796	20.0	61,295,920	0.04%
inst code/reasoning	612,208,775	1.0	612,208,775	0.45%
code	221,003,976,266	0.1	20,786,882,764	15.20%
curated (EN pes2o)	56,297,354,921	0.2	11,241,574,489	8.22%
curated (ZH CSL & wiki)	61,787,372	1.0	61,787,372	0.05%
curated (Gutenberg)	5,173,357,710	1.0	5,173,357,710	3.78%
mono high EN	3,002,029,817	0.1	300,202,982	0.22%
mono high	40,411,201,964	0.5	20,205,600,982	14.78%
mono medium-high	27,515,227,962	1.0	27,515,227,962	20.12%
mono medium	2,747,484,380	5.0	13,737,421,900	10.05%
mono medium-low	481,935,633	20.0	9,638,712,660	7.05%
mono low	97,535,696	50.0	4,876,784,800	3.57%

287 but follow precedent (Lozhkov et al., 2024) and subsample them more aggressively. For a more
 288 detailed read on filtering heuristics, we direct the reader to Appendix A.2.

290 **Data Mix** Our final data mix for continual training is listed in Table 2. The resource categorization
 291 refers to Appendix B.1 in the appendix and `inst medium-high+` is a separate category with
 292 languages with more than 500 million but less than 1B tokens. For monolingual text, we also have
 293 a separate category for English. In continual learning, where new data is introduced to an existing
 294 model, there is a risk of “catastrophic forgetting”, where the model loses knowledge from earlier
 295 training stages. Although our work’s primary focus is in a low-resource regime, we enhance the
 296 training corpus with a wide range of data types, including books and scientific papers in high-resource
 297 languages, code, and instruction data in our data mix. We downsample texts in high-resource
 298 languages and upsample text in low-resource languages using different sample rates according to
 299 how resourceful the language is. We make our data mix diverse and balanced towards different
 300 resource groups of languages in order to retain the prior knowledge of the model while learning new
 301 information, especially in medium- and low-resource languages, thus maintaining high performance
 302 across both previously seen and new languages. The final data mix has around 136B tokens.

303 3.2 MODEL TRAINING

305 We employ continual training using the causal language modelling objective for the decoder-only
 306 Llama model and exposing the pre-trained model to new data to develop our EMMA-500 model.
 307 We adopt efficient training strategies combining optimization, memory management, precision
 308 handling, and distributed training techniques. Our EMMA-500 model is trained on the Leonardo
 309 supercomputer⁹, occupying 256 Nvidia A100 GPUs, using the GPT-NeoX framework (Andonian
 310 et al., 2023). During training, we set a global batch size of 4096 and worked with sequences of 4096
 311 tokens. The training process ran for 12,000 steps, resulting in a total of 200 billion Llama 2 tokens.
 312 We use the Adam optimizer (Kingma & Ba, 2015) with a learning rate of 0.0001, betas of [0.9, 0.95],
 313 and an epsilon of 1e-8. We use a cosine learning rate scheduler with a warm-up of 500 iterations. To
 314 reduce memory consumption, activation checkpointing is employed. Precision is managed through
 315 mixed-precision techniques, using bfloat16 for computational efficiency and maintaining FP32 for
 316 gradient accumulation.

317 4 EVALUATION

319 4.1 TASKS, BENCHMARKS, AND BASELINES

321 **Tasks and Benchmarks** We conduct a comprehensive evaluation to validate the usability of our
 322 processed data and data mixing for massively multilingual language adaptation. We perform the

323 ⁹<https://leonardo-supercomputer.cineca.eu>

324 Table 3: Evaluation statistics. Sample/Lang: average number of test samples per language; N Lang:
 325 number of languages covered; NLL: negative log-likelihood; ACC: accuracy.

Tasks	Dataset	Metric	Samples/Lang	N Lang	Domain
Intrinsic Evaluation (Appendix E.1)	Glot500-c test (Imani et al., 2023)	NLL	1000	534	Misc
	PBC (Mayer & Cysouw, 2014)	NLL	500	370	Bible
Text Classification (Section 4.3)	SIB200 (Adelani et al., 2023)	ACC	204	205	Misc
	Taxi1500 (Ma et al., 2023)	ACC	111	1507	Bible
Commonsense Reasoning (Appendix E.3)	XCOPA (Ponti et al., 2020)	ACC	600	11	Misc
	XStoryCloze (Lin et al., 2022)	ACC	1870	11	Misc
	XWinograd (Tikhonov & Ryabinin, 2021)	ACC	741.5	6	Misc
Natural Language Inference (Appendix E.5)	XNLI (Conneau et al., 2018)	ACC	2490	15	Misc
Machine Translation (Section 4.2)	FLORES-200 (Costa-jussà et al., 2022)	BLEU, chrF++	1012	204	Misc
Open-Ended Generation (Section 4.4)	Aya (Singh et al., 2024)	BLEU, Self-BLEU	215	119	Misc
	PolyWrite (Ours)	Self-BLEU	149	240	Misc
Summarization (Appendix E.2)	XL-Sum (Hasan et al., 2021)	ROUGE-L, BERTScore	2537	44	News
Math (Appendix E.6)	MGSM direct (Shi et al., 2022)	ACC	250	10	Misc
	MGSM CoT (Shi et al., 2022)	ACC	250	10	Misc
Machine Comprehension (Appendix E.7)	BELEBELE (Bandarkar et al., 2023)	ACC	900	122	Misc
	ARC multilingual (Lai et al., 2023)	ACC	1170	31	Misc
Code Generation (Appendix E.8)	Multipl-E (Cassano et al., 2022)	Pass@k	164	7	Misc

343 intrinsic evaluation of the models’ performance on next-word prediction and evaluate the model’s
 344 performance on downstream tasks. Table 3 lists the datasets we used as downstream evaluation
 345 datasets in this work.

347 We present the evaluation results of intrinsic evaluation in Appendix E.1, commonsense reasoning in
 348 Appendix E.3, text summarization in Appendix E.2, math tasks in Appendix E.6, machine reading
 349 comprehension in Appendix E.7, and coding generation in Appendix E.8.

350 **Baselines** We compare our model with three groups of decoder-only models. They are (1) Llama 2
 351 models (Touvron et al., 2023) and continual pre-trained models based on Llama 2, such as CodeL-
 352 llama (Roziere et al., 2023), MaLA-500 (Lin et al., 2024), LLaMAX (Lu et al., 2024), Tower (Alves
 353 et al., 2024), and YaYi¹⁰; (2) other LLMs and continual pre-trained LLMs designed to be massively
 354 multilingual, including BLOOM (Scao et al., 2022), mGPT (Shliazhko et al., 2022), XGLM (Lin
 355 et al., 2022), and Occiglot¹¹; and (3) recent LLMs with superior English capabilities like Llama
 356 3 (Dubey et al., 2024), Llama 3.1¹², Qwen 2 (Yang et al., 2024a), and Gemma 2 (Team et al., 2024).
 357 There are also some other LLMs such as OpenAI’s API models¹³ and xLLMs-100 (Lai et al., 2024).
 358 However, they do not release the model weights or they limit access to them through commercial API,
 359 so we did not include them. The MADLAD model (Kudugunta et al., 2024) that uses the decoder-only
 360 T5 architecture is not supported by inference engines such as the HuggingFace transformers (Wolf
 361 et al., 2019). We do not compare them in this work.

363 4.2 MACHINE TRANSLATION

364 FLORES-200 is an evaluation benchmark for translation tasks with 204 language pairs involving
 365 English and thus 408 translation directions, with a particular focus on low-resource languages. We
 366 assess all language models by adopting a 3-shot evaluation approach with the prompt in Appendix D.1.

368 The performance is measured by BLEU (Papineni et al., 2002) and chrF++ (Popović, 2015) imple-
 369 mented in `sacrebleu` (Post, 2018). The BLEU score is calculated with the `flores200` tokenizer
 370 applied to the texts and chrF++ uses word order 2. The choice of `flores200` tokenization ensures
 371 that languages that do not have a whitespace delimiter can be evaluated at the (sub-)word level. For
 372 reproducibility, we attach the BLEU and chrF++ signatures.^{14,15}

373 ¹⁰<https://huggingface.co/wenge-research/yayi-7b-llama2>

374 ¹¹<https://huggingface.co/occiglot/occiglot-7b-eu5>

375 ¹²https://llama.meta.com/docs/model-cards-and-prompt-formats/llama3_1

376 ¹³<https://platform.openai.com/docs/models>

377 ¹⁴BLEU: nrefs:1—case:mixed—eff:no—tok:flores200—smooth:exp—version:2.4.2

378 ¹⁵chrF++: nrefs:1—case:mixed—eff:yes—nc:6—nw:2—space:no—version:2.4.2

378 Table 4: 3-shot results on FLORES-200 (X-Eng, BLEU/chrF++). EMMA-500 Llama 2 7B has better
 379 average performance than all baselines.
 380

Model	Avg	High	Medium-High	Medium	Medium-Low	Low
Llama 2 7B	12.93/ 30.32	19.91/ 39.04	17.56/ 35.84	12.49/ 29.81	8.27/ 24.35	6.96/ 23.36
Llama 2 7B Chat	12.28/ 31.72	18.98/ 39.65	17.06/ 37.03	11.74/ 31.1	7.79/ 26.34	6.18/ 25.03
CodeLlama 2 7B	10.82/ 28.57	17.39/ 37.43	15.27/ 33.94	10.39/ 28.05	6.45/ 22.85	5.04/ 21.48
LLaMAX Llama 2 7B	1.99/ 13.66	3.68/ 22.18	2.95/ 18.15	1.83/ 12.84	0.67/ 7.2	1.01/ 9.04
LLaMAX Llama 2 7B Alpaca	22.29/ 42.27	32.83/ 54.56	30.04/ 51.25	21.7/ 41.94	13.06/ 31.32	14.24/ 32.88
MaLA-500 Llama 2 10B v1 [‡]	2.29/ 13.6	4.64/ 15.95	3.18/ 14.64	2.68/ 14.23	1.24/ 12.58	0.33/ 11.18
MaLA-500 Llama 2 10B v2 [‡]	2.87/ 15.44	5.58/ 18.65	3.81/ 16.33	3.55/ 16.29	1.63/ 14.2	0.55/ 12.76
Yayi Llama 2 7B	12.98/ 31.38	19.48/ 39.58	17.55/ 36.71	12.47/ 30.79	8.54/ 25.63	7.22/ 24.84
TowerBase Llama 2 7B	13.74/ 31.47	21.76/ 40.96	18.92/ 37.27	13.15/ 30.9	8.3/ 25.05	7.21/ 24.1
TowerInstruct Llama 2 7B	4.81/ 25.43	9.18/ 34.4	6.66/ 30.01	4.62/ 25.22	2.64/ 20.24	1.8/ 18.69
Occiglot Mistral 7B v0.1	13.12/ 31.13	19.53/ 38.93	17.57/ 36.27	13.07/ 31.2	9.03/ 26.15	6.86/ 23.83
Occiglot Mistral 7B v0.1 Instruct	11.61/ 31.65	16.72/ 39.28	15.06/ 36.48	11.7/ 31.73	8.48/ 26.88	6.54/ 24.7
BLOOM 7B	9.57/ 27.84	15.75/ 36.65	9.65/ 28.19	9.42/ 27.81	6.81/ 23.95	8.61/ 25.89
BLOOMZ 7B [†]	20.22/ 34.74	32.23/ 47.03	19.2/ 34.08	20.09/ 34.49	16.25/ 30.58	18.54/ 32.63
mGPT	5.29/ 20.69	9.37/ 26.64	8.28/ 25.29	3.41/ 17.87	2.43/ 16.07	2.84/ 17.28
mGPT-13B	7.42/ 24.58	12.61/ 31.95	11.11/ 30.16	5.72/ 22.49	3.57/ 18.16	4.11/ 20.04
Yayi 7B	4.82/ 21.36	5.69/ 25.18	4.53/ 19.97	4.41/ 21.52	3.71/ 19.18	6.13/ 23.12
Llama 3 8B	23.78/ 43.72	33.71/ 55.36	30.31/ 51.3	24.75/ 44.91	15.18/ 33.65	16.01/ 34.65
Llama 3.1 8B	24.19/ 44.1	34.15/ 55.7	30.79/ 51.7	24.98/ 45.26	15.89/ 34.24	16.13/ 34.85
Gemma 2 9B	23.15/ 38.87	33.11/ 51.36	30.81/ 48.53	25.58/ 41.23	15.37/ 30.03	11.73/ 24.15
Gemma 7B	23.79/ 43.68	34.23/ 55.77	29.87/ 50.95	24.0/ 44.25	16.16/ 34.36	16.03/ 34.58
Qwen 1.5 7B	15.58/ 35.87	24.07/ 46.29	19.92/ 40.74	15.76/ 36.27	9.74/ 28.81	9.77/ 29.13
Qwen 2 7B	17.39/ 37.61	27.63/ 50.06	22.48/ 43.28	18.13/ 38.63	9.89/ 28.54	10.64/ 29.99
EMMA-500 Llama 2 7B	25.37/ 45.78	32.24/ 53.74	31.39/ 52.85	25.72/ 46.16	20.32/ 39.96	17.18/ 36.15

400 Table 4 presents the average X-to-English (X-Eng) translation results.¹⁶ Our EMMA-500 model
 401 outperforms all other models on average. We achieve the best performance across all language
 402 settings, except for high-resource languages where our model slightly lags behind Llama 3/3.1,
 403 Gemma 7B, and LLaMAX 7B Alpaca. In the English-to-X (Eng-X) translation direction, as shown in
 404 Table 16 in Appendix E.4, the advantage of EMMA-500 is even more pronounced. We outperform
 405 all other models even in high-resource languages, and the advantage becomes more significant in
 406 lower-resource languages. Overall, we note that our model outperforms Tower models which are
 407 explicitly adjusted to perform translation tasks in high-resource languages. Further, the much larger
 408 margin between EMMA-500 and other models in Eng-X compared with X-Eng indicates that our
 409 EMMA-500 model is particularly good at generating non-English texts.

4.3 TEXT CLASSIFICATION

412 SIB-200 (Adelani et al., 2023) and Taxi1500 (Ma et al., 2023) are two prominent topic classification
 413 datasets. SIB-200 encompasses seven categories: science/technology, travel, politics, sports, health,
 414 entertainment, and geography. Taxi1500 spans 1507 languages, involving six classes: Recommendation,
 415 Faith, Description, Sin, Grace, and Violence. We use 3-shot prompting with prompts in
 416 Appendix D.2, drawing demonstrations from the development set and testing models on the test split.
 417 The outcomes on SIB-200 and TAXI-1500 are tabulated in Table 5. For SIB-200, our EMMA-500
 418 model outperforms all Llama2-based models, with particularly notable gains in languages with
 419 medium or fewer resources—seeing an average improvement of 47.5%. Taxi-1500 could be a more
 420 challenging task since it is in the religious domain, but our model still surpasses all Llama2-based
 421 models except for MaLA-500. However, despite these improvements in both classification tasks, our
 422 models lag behind the latest models such as Llama3 and 3.1, especially in high-resource languages.

4.4 OPEN-ENDED GENERATION

425 **Aya Evaluation** We choose the two subsets aya-human-annotated and
 426 dolly-machine-translated from the Aya evaluation suite (Singh et al., 2024), which have
 427 both inputs and targets for subsequent evaluation. To quantitatively assess the quality of the generated
 428 text by the models, we employ two metrics: BLEU (Papineni et al., 2002) and Self-BLEU (Zhu et al.,

429 ¹⁶We mark BLOOMZ with a † because it has used FLORES in its instruction tuning data; we mark MaLA-500
 430 with a ‡ because it has used FLORES in its training data but with source and target sides split. Besides, as
 431 a remark, Tower, LLaMAX, and our EMMA-500 have intentionally used parallel data (not FLORES) in the
 432 training stage.

432 Table 5: 3-shot results on SIB-200 and Taxi-1500 (ACC). EMMA-500 Llama 2 7B has better average
 433 performance than Llama 2 models and comparable performance with multilingual LLMs, and has
 434 comparable performance with the compared LLMs.

436 Model	SIB-200						Taxi-1500					
	Avg	High	Med-High	Medium	Med-Low	Low	Avg	High	Med-High	Medium	Med-Low	Low
Llama 2 7B	0.2241	0.2664	0.2469	0.2205	0.1968	0.1900	0.1754	0.1950	0.1949	0.1847	0.1746	0.1737
Llama 2 7B Chat	0.2558	0.2972	0.2811	0.2501	0.2303	0.2191	0.1544	0.1873	0.1766	0.1661	0.1559	0.1522
CodeLlama 2 7B	0.2335	0.2606	0.2542	0.2310	0.2142	0.2037	0.1703	0.1745	0.1741	0.1720	0.1705	0.1700
LLaMAX Llama 2 7B	0.1061	0.1242	0.1160	0.1001	0.0945	0.0954	0.2352	0.2320	0.2340	0.2376	0.2356	0.2352
LLaMAX Llama 2 7B Alpaca	0.2789	0.3309	0.3212	0.2716	0.2338	0.2282	0.1509	0.1870	0.1688	0.1599	0.1500	0.1491
MaLA-500 Llama 2 10B v1	0.2325	0.2330	0.2364	0.2288	0.2276	0.2358	0.2527	0.2390	0.2402	0.2476	0.2457	0.2543
MaLA-500 Llama 2 10B v2	0.1930	0.1893	0.2105	0.1949	0.1755	0.1846	0.2339	0.2136	0.2230	0.2132	0.2172	0.2367
Yayi Llama 2 7B	0.2457	0.2904	0.2717	0.2442	0.2144	0.2069	0.1773	0.1874	0.1846	0.1819	0.1789	0.1765
TowerBase Llama 2 7B	0.1934	0.2200	0.2092	0.1874	0.1790	0.1693	0.1773	0.1849	0.1881	0.1867	0.1810	0.1761
TowerInstruct Llama 2 7B	0.2053	0.2321	0.2196	0.2026	0.1915	0.1804	0.1729	0.2017	0.1960	0.1808	0.1740	0.1709
Occiglot Mistral 7B v0.1	0.3269	0.3880	0.3582	0.3174	0.2892	0.2836	0.2226	0.2464	0.2291	0.2299	0.2233	0.2215
Occiglot Mistral 7B v0.1 Instruct	0.3431	0.3948	0.3716	0.3336	0.3147	0.3008	0.1876	0.2430	0.2090	0.1941	0.1918	0.1848
BLOOM 7B	0.1781	0.2313	0.1805	0.1717	0.1576	0.1702	0.1476	0.1558	0.1489	0.1456	0.1511	0.1473
BLOOMZ 7B	0.2973	0.3039	0.2963	0.2980	0.2953	0.2970	0.1696	0.1693	0.1698	0.1696	0.1699	0.1695
mGPT	0.2711	0.2858	0.2799	0.2673	0.2589	0.2648	0.1072	0.0867	0.0844	0.0992	0.1029	0.1093
mGPT-13B	0.3320	0.3669	0.3427	0.3448	0.2939	0.3226	0.1723	0.1798	0.1644	0.1588	0.1610	0.1738
XGLM 7.5B	0.3181	0.3528	0.3512	0.3169	0.2696	0.2996	0.2041	0.2421	0.2369	0.2125	0.2105	0.2010
Yayi 7B	0.3576	0.4057	0.3620	0.3563	0.3472	0.3318	0.1612	0.1665	0.1638	0.1645	0.1583	0.1611
Llama 3 8B	0.6369	0.7345	0.7025	0.6462	0.5316	0.5696	0.2173	0.3184	0.2708	0.2560	0.2261	0.2105
Llama 3.1 8B	0.6142	0.7070	0.6790	0.6199	0.5146	0.5475	0.2020	0.2751	0.2521	0.2443	0.2097	0.1959
Gemma 7B	0.5821	0.6806	0.6455	0.5816	0.4886	0.5112	0.1805	0.2868	0.2855	0.2499	0.2028	0.1689
Gemma 2 9B	0.4625	0.5177	0.4900	0.4692	0.4304	0.4045	0.1383	0.2429	0.2413	0.1874	0.1497	0.1283
Qwen 1.5 7B	0.4795	0.5600	0.5181	0.4825	0.4156	0.4286	0.0729	0.1265	0.1145	0.0878	0.0730	0.0692
Qwen 2 7B	0.5495	0.6637	0.6014	0.5517	0.4519	0.4925	0.2187	0.2737	0.2557	0.2401	0.2233	0.2147
EMMA-500 Llama 2 7B	0.3127	0.3275	0.3328	0.3099	0.3083	0.2760	0.1982	0.2366	0.2333	0.2277	0.2200	0.1930

456 2018). BLEU is crucial for assessing the linguistic accuracy and relevance of the generated text in
 457 comparison to the expected human-like text present in the dataset. Self-BLEU is used to evaluate the
 458 diversity of the text generated by a model. It measures how similar different texts from the same
 459 model are to each other by treating one generated text as the “candidate” and others as the “reference”
 460 texts. This metric is useful in scenarios where high degrees of variation are desirable, as it helps
 461 identify models that might be overfitting to particular styles or patterns of text.

462 The results are presented in Table 6. The BLEU metric is an indicator that measures how generated
 463 texts are close to the references. However, in open-ended generation, LLMs cannot generate identical
 464 texts to references, leading to low BLEU scores. The EMMA-500 model obtains remarkably high
 465 BLEU scores in both high-resource and low-resource settings when compared with baselines. Its
 466 performance in low-resource settings is particularly noteworthy, as it not only sustains high BLEU
 467 scores but also exhibits a Self-BLEU score of 5.09, the highest among all models evaluated. This
 468 high Self-BLEU score indicates less diversity in the generated text, suggesting that while EMMA-500
 469 maintains consistency, it may produce less varied outputs. When compared to other high-performing
 470 models like Qwen 2 7B and Llama 3.1 8B, the EMMA-500 model exhibits a superior balance between
 471 accuracy and linguistic creativity. Unlike Qwen 2 7B, which shows a spike in performance primarily
 472 in medium-low resource settings, EMMA-500 maintains a consistently high performance across
 473 varying levels of resource availability.

474 **PolyWrite** This is a novel multilingual benchmark composed in this work for evaluating open-
 475 ended generation in 240 languages. We use ChatGPT to generate different prompts in English and use
 476 Google Translate to translate them into different languages for models to generate creative content.
 477 This benchmark consists of 31 writing tasks, such as storytelling and email writing, and 155 prompts
 478 in total. We back-translate the multilingual prompts to English, calculate the BLEU scores between
 479 original English prompts and back-translation, and filter out translated prompts with BLEU scores
 480 below 20, and the entire dataset contains a total of 35,751 prompts. The details of PolyWrite are
 481 described in Appendix C.1.

482 We use Self-BLEU (Zhu et al., 2018) to evaluate the diversity of generated texts in the PolyWrite
 483 benchmark, as presented in Table 6. A lower Self-BLEU score indicates more diverse generation,
 484 but does not mean a better generation quality. Our EMMA-500 model demonstrates comparable
 485 performance across various languages. Compared to other models like Llama 3/3.1 and Qwen 1.5/2,

particularly in medium-low and low-resource languages, EMMA-500 has higher Self-BLEU scores, indicating lower diversity in its generated content.

Evaluating open-ended generation poses significant challenges, as it goes beyond simply measuring accuracy or correctness. Metrics like BLEU or Self-BLEU, while useful for assessing similarity to reference texts or the diversity of given texts, often fail to capture more nuanced aspects. Subjective factors like cultural relevance and the appropriateness of responses in low-resource languages are difficult to quantify. This makes it challenging to create evaluation benchmarks and metrics that fully capture the strengths and weaknesses of models like EMMA-500 in diverse, real-world scenarios.

Table 6: Results on Aya (BLEU/Self-BLEU) and PolyWrite (Self-BLEU). EMMA-500 Llama 2 7B has higher average BLEU scores than all baselines on Aya.

Model	Aya						PolyWrite					
	Avg	High	Med-High	Medium	Med-Low	Low	Avg	High	Med-High	Medium	Med-Low	Low
Llama 2 7B	1.24/0.74	1.27/0.57	1.47/0.60	0.86/0.37	1.17/0.45	0.77/1.87	0.5358	0.4282	0.6545	0.3769	0.4766	0.6587
Llama 2 7B Chat	1.17/1.29	1.46/1.15	1.36/1.15	0.69/1.03	1.14/1.14	0.54/2.23	1.1550	0.8640	0.8902	1.1877	1.4435	1.4167
CodeLlama 2 7B	1.22/1.21	1.31/1.19	1.40/1.00	0.85/0.84	1.23/0.56	0.78/2.57	1.0313	1.2052	1.2883	0.9191	0.9092	0.6798
LLaMAX Llama 2 7B	1.72/1.70	1.80/1.37	2.03/1.48	1.30/1.24	1.65/0.89	0.98/3.74	0.9564	1.0066	1.1655	0.8786	0.8363	0.7709
LLaMAX Llama 2 7B Alpaca	1.68/1.67	1.82/1.22	1.96/1.42	1.28/1.14	1.55/1.04	1.00/3.94	0.8086	0.7321	0.9517	0.8322	0.8351	0.4981
MaLA-500 Llama 2 10B v1	0.40/2.29	0.42/2.53	0.49/2.79	0.32/1.22	0.31/2.42	0.18/1.16	3.7079	3.7902	3.5066	2.4541	4.8052	3.5163
MaLA-500 Llama 2 10B v2	0.41/2.31	0.42/2.01	0.50/2.65	0.32/1.42	0.33/1.02	0.18/3.09	4.0059	3.1737	4.0148	3.0476	5.0652	3.8916
Yayi Llama 2 7B	1.65/0.61	1.84/0.82	1.88/0.62	1.26/0.62	1.53/0.45	1.02/0.41	0.6274	0.6207	0.6921	0.4169	0.6813	0.6352
TowerBase Llama 2 7B	1.44/0.83	1.45/0.64	1.67/0.56	1.19/0.49	1.46/0.38	0.88/2.54	0.4938	0.7268	0.4736	0.3945	0.4417	0.5396
TowerInstruct Llama 2 7B	1.55/0.93	1.80/0.85	1.82/0.78	1.15/0.65	1.21/0.54	0.85/1.97	0.7124	0.9565	0.7651	0.6492	0.5998	0.6615
Occiglot Mistral 7B v0.1	1.53/2.43	1.57/0.91	1.78/2.89	1.17/1.73	1.61/1.31	0.95/4.27	0.9647	0.8818	0.7975	0.9543	0.9563	1.4157
Occiglot Mistral 7B v0.1 Instruct	0.75/2.81	0.82/2.38	0.89/3.10	0.43/2.72	0.79/1.17	0.45/3.50	3.9033	5.3884	4.7140	3.7555	2.8444	2.9568
BLOOM 7B	0.85/1.17	0.92/1.20	0.96/1.32	0.73/1.05	0.78/1.27	0.55/0.72	1.3892	1.2845	1.6705	1.9685	1.1513	0.6600
BLOOMZ 7B	0.12/0.61	0.07/0.33	0.17/0.62	0.08/1.00	0.06/0.89	0.08/0.51	0.0024	0.0000	0.0005	0.0093	0.0000	0.0049
mGPT	1.24/0.55	1.22/0.64	1.47/0.59	0.91/0.48	1.21/0.60	0.84/0.29	0.7560	0.9222	0.6291	0.4534	0.8156	1.1134
mGPT-13B	1.42/0.57	1.42/0.80	1.63/0.58	1.00/0.48	1.53/0.44	1.03/0.36	0.7479	0.8483	0.7091	0.4230	0.7962	1.0201
Yayi 7B	1.05/0.39	1.22/0.38	1.18/0.41	0.76/0.42	1.01/0.54	0.67/0.22	0.5151	0.4266	0.3574	0.3283	0.6706	0.8574
Llama 3 8B	1.59/0.94	1.03/0.60	1.31/0.53	2.96/0.54	1.11/0.28	2.43/3.47	0.5796	0.5753	0.5921	0.7141	0.5586	0.4440
Llama 3.1 8B	1.85/1.08	1.41/0.95	1.60/0.64	3.11/0.52	1.33/0.39	2.52/3.52	0.7995	0.6805	0.8253	0.5044	0.6965	1.3585
Gemma 2 9B	1.55/0.82	1.59/0.94	1.73/0.93	1.38/0.70	1.33/0.47	1.21/0.65	1.1736	1.2347	1.2913	1.1616	0.9261	1.3240
Gemma 7B	1.29/0.57	1.41/0.62	1.39/0.66	1.16/0.40	1.26/0.50	0.93/0.40	1.0541	0.9629	1.1222	0.9275	1.1112	1.0284
Qwen 1.5 7B	1.93/1.13	1.66/0.85	1.64/0.65	3.14/0.79	1.51/0.58	2.45/3.69	0.6441	0.9398	0.7457	0.3797	0.4164	0.8738
Qwen 2 7B	1.99/1.13	1.84/0.94	1.69/0.70	3.29/0.72	1.36/0.44	2.47/3.53	0.5709	0.7763	0.6135	0.3695	0.4186	0.7989
EMMA-500 Llama 2 7B	2.93/1.54	2.90/1.29	2.75/0.83	3.79/0.82	2.86/0.95	2.87/5.09	0.9879	0.9833	0.7058	0.8465	1.3130	1.1702

5 CONCLUSION AND OUTLOOKS

This paper addresses critical advancements and challenges in adapting language models to more than 500 languages, focusing on enhancing their performance across diverse languages. We compile the MaLA corpus, a multilingual dataset for continual pre-training of multilingual language models. By expanding and augmenting existing corpora, we train the EMMA-500 model. It demonstrates notable improvements in a range of tasks such as next-token prediction, commonsense reasoning, machine translation, text classification, and open-ended generation, showing remarkable improvements in low-resource languages. Our results show that a well-curated, massively multilingual corpora can advance model capabilities. This work sets a new benchmark for inclusive and effective multilingual language models and paves the way for future research to address the disparities between high-resource and low-resource languages.

Multilingual language models are designed to cater to diverse linguistic and cultural backgrounds, yet many existing multilingual benchmarks rely on human or machine translations, primarily reflecting English-speaking communities and introducing imperfections that affect evaluation integrity. We emphasize the need for natively-created multilingual test sets to provide more accurate assessments. While our EMMA-500 model shows enhanced multilingual performance compared to Llama 2 7B and other variants, it falls short on certain benchmarks compared to newer models like Llama 3 and Gemma 2. Future efforts will explore multilingual extension with these models and direct pre-training using the released corpus, along with multilingual instruction tuning to improve task performance and interactions.

540 REFERENCES
541

- 542 Solomon Teferra Abate, Michael Melese, Martha Yifiru Tachbelie, Million Meshesha, Solomon
543 Atinifu, Wondwossen Mulugeta, Yaregal Assabie, Hafte Abera, Binyam Ephrem, Tewodros
544 Abebe, Wondimagegnhue Tsegaye, Amanuel Lemma, Tsegaye Andargie, and Seifedin Shifaw.
545 Parallel corpora for bi-lingual English-Ethiopian languages statistical machine translation. In
546 *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 3102–3111,
547 Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics. URL
548 <https://aclanthology.org/C18-1262>.
- 549 Ahmed Abdelali, Hamdy Mubarak, Younes Samih, Sabit Hassan, and Kareem Darwish. QADI:
550 Arabic dialect identification in the wild. In *Proceedings of the Sixth Arabic Natural Language
551 Processing Workshop*, pp. 1–10, Kyiv, Ukraine (Virtual), April 2021. Association for Computational
552 Linguistics. URL <https://aclanthology.org/2021.wanlp-1.1>.
- 553 Kathrein Abu Kwaik, Motaz Saad, Stergios Chatzikyriakidis, and Simon Dobnik. Shami: A corpus
554 of Levantine Arabic dialects. In *Proceedings of the Eleventh International Conference on Lan-
555 guage Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language
556 Resources Association (ELRA). URL <https://aclanthology.org/L18-1576>.
- 557 David Adelani, Dana Ruiter, Jesujoba Alabi, Damilola Adebomojo, Adesina Ayeni, Mofe Adeyemi,
558 Ayodele Esther Awokoya, and Cristina España-Bonet. The effect of domain and diacritics in
559 Yoruba–English neural machine translation. In *Proceedings of Machine Translation Summit XVIII:
560 Research Track*, pp. 61–75, Virtual, August 2021. Association for Machine Translation in the
561 Americas. URL <https://aclanthology.org/2021.mtsummit-research.6>.
- 562 David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter,
563 Dietrich Klakow, Peter Nabende, Ernie Chang, Tajudeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou,
564 Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad,
565 Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala,
566 Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott,
567 Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi,
568 Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire
569 Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulkumin, Ayodele
570 Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. A few thou-
571 sand translations go a long way! leveraging pre-trained models for African news translation.
572 In *Proceedings of the 2022 Conference of the North American Chapter of the Association for
573 Computational Linguistics: Human Language Technologies*, pp. 3053–3070, Seattle, United States,
574 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.223.
575 URL <https://aclanthology.org/2022.naacl-main.223>.
- 576 David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba O. Alabi, Yanke
577 Mao, Haonan Gao, and En-Shiun Annie Lee. SIB-200: A simple, inclusive, and big evaluation
578 dataset for topic classification in 200+ languages and dialects. *Corr*, abs/2309.07445, 2023. doi:
579 10.48550/arXiv.2309.07445. URL <https://doi.org/10.48550/arXiv.2309.07445>.
- 580 Rodrigo Agerri, Xavier Gómez Guinovart, German Rigau, and Miguel Anxo Solla Portela. De-
581 veloping new linguistic resources and tools for the Galician language. In *Proceedings of
582 the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*,
583 Miyazaki, Japan, May 2018. European Language Resources Association (ELRA). URL <https://aclanthology.org/L18-1367>.
- 584 Israa Alsarsour, Esraa Mohamed, Reem Suwaileh, and Tamer Elsayed. DART: A large dataset of
585 dialectal Arabic tweets. In *Proceedings of the Eleventh International Conference on Language Re-
586 sources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language Resources
587 Association (ELRA). URL <https://aclanthology.org/L18-1579>.
- 588 Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian,
589 Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. Tower: An open multilingual
590 large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*, 2024. URL
591 <https://arxiv.org/abs/2402.17733>.

- 594 Antonios Anastasopoulos, Alessandro Cattelan, Zi-Yi Dou, Marcello Federico, Christian Federmann,
 595 Dmitriy Genzel, Franscisco Guzmán, Junjie Hu, Macduff Hughes, Philipp Koehn, Rosie Lazar,
 596 Will Lewis, Graham Neubig, Mengmeng Niu, Alp Öktem, Eric Paquin, Grace Tang, and Sylwia
 597 Tur. TICO-19: the translation initiative for COvid-19. In *Proceedings of the 1st Workshop on NLP*
 598 *for COVID-19 (Part 2) at EMNLP 2020*, Online, December 2020. Association for Computational
 599 Linguistics. doi: 10.18653/v1/2020.nlp covid19-2.5. URL <https://aclanthology.org/2020.nlpcovid19-2.5>.
- 600
- 601 Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Halla-
 602 han, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Jason Phang,
 603 Shivanshu Purohit, Hailey Schoelkopf, Dashiell Stander, Tri Songz, Curt Tigges, Benjamin Thérien,
 604 Phil Wang, and Samuel Weinbach. GPT-NeoX: Large Scale Autoregressive Language Modeling in
 605 PyTorch, 9 2023. URL <https://www.github.com/eleutherai/gpt-neox>.
- 606
- 607 Mikko Aulamo, Nikolay Bogoychev, Shaoxiong Ji, Graeme Nail, Gema Ramírez-Sánchez, Jörg Tiede-
 608 mann, Jelmer Van Der Linde, and Jaume Zaragoza. Hplt: High performance language technologies.
 609 In *Proceedings of the 24th Annual Conference of the European Association for Machine Transla-*
 610 *tion*, pp. 517–518, 2023. URL <https://aclanthology.org/2023.eamt-1.61>.
- 611 Niyati Bafna. Empirical models for an indic language continuum, 2022.
- 612 Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa,
 613 Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. The BELEBELE
 614 benchmark: a parallel reading comprehension dataset in 122 language variants. *arXiv preprint*
 615 *arXiv:2308.16884*, 2023. URL <https://aclanthology.org/2024.acl-long.44>.
- 616
- 617 Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis,
 618 Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo
 619 Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarriás, Marek Strelec, Brian Thompson, William
 620 Waites, Dion Wiggins, and Jaume Zaragoza. ParaCrawl: Web-scale acquisition of parallel corpora.
 621 In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp.
 622 4555–4567, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.
 623 acl-main.417. URL <https://aclanthology.org/2020.acl-main.417>.
- 624
- 625 Marta Bañón, Miquel Esplà-Gomis, Mikel L. Forcada, Cristian García-Romero, Taja Kuzman, Nikola
 626 Ljubetic, Rik van Noord, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Peter Rupnik, Vít
 627 Suchomel, Antonio Toral, Tobias van der Werff, and Jaume Zaragoza. Macocu: Massive collection
 628 and curation of monolingual and bilingual data: focus on under-resourced languages. In Helena
 629 Moniz, Lieve Macken, Andrew Rufener, Loïc Barrault, Marta R. Costa-jussà, Christophe Declercq,
 630 Maarit Koponen, Ellie Kemp, Spyridon Pilos, Mikel L. Forcada, Carolina Scarton, Joachim Van
 631 den Bogaert, Joke Daems, Arda Tezcan, Bram Vanroy, and Margot Fonteyne (eds.), *Proceedings of*
 632 *the 23rd Annual Conference of the European Association for Machine Translation, EAMT 2022,*
 633 *Ghent, Belgium, June 1-3, 2022*, pp. 301–302. European Association for Machine Translation,
 634 2022. URL <https://aclanthology.org/2022.eamt-1.41>.
- 635
- 636 Andrei Z Broder. On the resemblance and containment of documents. In *Proceedings. Compression*
 637 *and Complexity of SEQUENCES 1997 (Cat. No. 97TB100171)*, pp. 21–29. IEEE, 1997.
- 638
- 639 José Camacho-Collados, Claudio Delli Bovi, Alessandro Raganato, and Roberto Navigli. A large-
 640 scale multilingual disambiguation of glosses. In *Proceedings of the Tenth International Conference*
 641 *on Language Resources and Evaluation (LREC'16)*, pp. 1701–1708, Portorož, Slovenia, May 2016.
 642 European Language Resources Association (ELRA). URL <https://aclanthology.org/L16-1269>.
- 643
- 644 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald
 645 Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. Multipl-e:
 646 A scalable and extensible approach to benchmarking neural code generation. *arXiv preprint*
 647 *arXiv:2208.08227*, 2022.
- 648
- 649 Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and Benjamin K. Bergen. When is multilinguality a
 650 curse? Language modeling for 250 high- and low-resource languages. *arXiv preprint*, 2023. URL
 651 <https://arxiv.org/abs/2311.09205>.

- 648 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
 649 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
 650 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 651
- 652 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
 653 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
 654 *ArXiv*, abs/1803.05457, 2018.
- 655 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher
 656 Hesse, and John Schulman. Training verifiers to solve math word problems, 2021.
- 657
- 658 Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger
 659 Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. In
 660 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
 Association for Computational Linguistics, 2018.
- 661
- 662 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Fran-
 663 cisco Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised
 664 cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the*
 665 *Association for Computational Linguistics*, pp. 8440–8451, 2020.
- 666 Corpora and Tools. Languages of Russia: Collections of texts in small languages, n.d. URL
 667 <http://web-corpora.net/wsg3/minorlangs/download>.
- 668
- 669 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan,
 670 Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling
 671 human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.
- 672
- 673 Ona de Gibert, Graeme Nail, Nikolay Arefyev, Marta Bañón, Jelmer van der Linde, Shaoxiong Ji,
 674 Jaume Zaragoza-Bernabeu, Mikko Aulamo, Gema Ramírez-Sánchez, Andrey Kutuzov, et al. A
 675 new massive multilingual dataset for high-performance language technologies. In *Proceedings of*
LREC-COLING, 2024.
- 676
- 677 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
 678 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
the North American Chapter of the Association for Computational Linguistics: Human Language
Technologies, 2019.
- 679
- 680 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 681 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 682 *arXiv preprint arXiv:2407.21783*, 2024.
- 683
- 684 Jonathan Dunn. Mapping languages: the corpus of global language use. *Lang. Resour. Evaluation*,
 685 54(4):999–1018, 2020. doi: 10.1007/s10579-020-09489-2. URL <https://doi.org/10.1007/s10579-020-09489-2>.
- 686
- 687 EdTeKLA. Indigenous languages corpora, 2022. URL https://github.com/EdTeKLA/IndigenousLanguages_Corpora.
- 688
- 689 Mahmoud El-Haj. Habibi - a multi dialect multi national Arabic song lyrics corpus. In *Proceedings*
 690 *of the Twelfth Language Resources and Evaluation Conference*, pp. 1318–1326, Marseille, France,
 691 May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.165>.
- 692
- 693 Mahmoud El-Haj, Paul Rayson, and Mariam Aboelezz. Arabic dialect identification in the context of
 694 bivalency and code-switching. In *Proceedings of the Eleventh International Conference on Lan-*
695 guage Resources and Evaluation (LREC 2018), Miyazaki, Japan, May 2018. European Language
 696 Resources Association (ELRA). URL <https://aclanthology.org/L18-1573>.
- 697
- 698 Manuel Faysse. Dataset card for "project gutenberg", 2023. URL https://huggingface.co/datasets/manu/project_gutenberg.
- 699
- 700 Fitsum Gaim, Wonsuk Yang, and Jong C. Park. Tlmd: Tigrinya language modeling dataset (1.0.0),
 701 2021. URL <https://doi.org/10.5281/zenodo.5139094>. Dataset.

- 702 Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence
 703 Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot
 704 language model evaluation. Zenodo, 2023.
- 705
- 706 Javier García Gilabert, Carlos Escolano, Aleix Sant Savall, Francesca De Luca Fornaciari, Audrey
 707 Mash, Xixian Liao, and Maite Melero. Investigating the translation capabilities of large language
 708 models trained on parallel data only, 2024. URL <https://arxiv.org/abs/2406.09140>.
- 709
- 710 Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. Building large monolingual dictionaries
 711 at the Leipzig corpora collection: From 100 to 200 languages. In Nicoletta Calzolari, Khalid
 712 Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion
 713 Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Eighth International Conference
 714 on Language Resources and Evaluation (LREC'12)*, pp. 759–765, Istanbul, Turkey, May 2012.
 European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2012/pdf/327_Paper.pdf.
- 715
- 716 Santiago Góngora, Nicolás Giossa, and Luis Chiruzzo. Experiments on a Guarani corpus of news and
 717 social media. In *Proceedings of the First Workshop on Natural Language Processing for Indigenous
 718 Languages of the Americas*, pp. 153–158, Online, June 2021. Association for Computational
 719 Linguistics. doi: 10.18653/v1/2021.americasnlp-1.16. URL <https://aclanthology.org/2021.americasnlp-1.16>.
- 720
- 721
- 722 Santiago Góngora, Nicolás Giossa, and Luis Chiruzzo. Can we use word embeddings for enhancing
 723 Guarani-Spanish machine translation? In *Proceedings of the Fifth Workshop on the Use of
 724 Computational Methods in the Study of Endangered Languages*, pp. 127–132, Dublin, Ireland,
 725 May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.computel-1.16.
 726 URL <https://aclanthology.org/2022.computel-1.16>.
- 727
- 728 Thamme Gowda, Zhao Zhang, Chris Mattmann, and Jonathan May. Many-to-English machine
 729 translation tools, data, and pretrained models. In *Proceedings of the 59th Annual Meeting
 730 of the Association for Computational Linguistics and the 11th International Joint Conference
 731 on Natural Language Processing: System Demonstrations*, pp. 306–316, Online, August
 732 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-demo.37. URL
<https://aclanthology.org/2021.acl-demo.37>.
- 733
- 734 Grupo de Inteligencia Artificial PUCP. Monolingual and parallel corpora of peruvian languages, n.d.
 735 URL <https://github.com/iapucp/multilingual-data-peru>.
- 736
- 737 Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubashir, Yuan-Fang Li, Yong-Bin
 738 Kang, M. Sohel Rahman, and Rifat Shahriyar. XL-sum: Large-scale multilingual abstractive
 739 summarization for 44 languages. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli
 740 (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 4693–
 4703, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
 findings-acl.413. URL <https://aclanthology.org/2021.findings-acl.413>.
- 741
- 742 Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L. Richter, Quentin Anthony, Timothée
 743 Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train
 744 large language models, 2024. URL <https://arxiv.org/abs/2403.08763>.
- 745
- 746 David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D'souza, Julia Kreutzer, Constantine
 747 Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, et al.
 748 Masakhaner: Named entity recognition for african languages. *arXiv e-prints*, pp. arXiv–2103,
 749 2021.
- 750
- 751 Ayyoob Imani, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora
 752 Kassner, Chunlan Ma, Helmut Schmid, André F. T. Martins, François Yvon, and Hinrich Schütze.
 753 Glot500: Scaling multilingual corpora and language models to 500 languages. *arXiv preprint*,
 2023. URL <https://aclanthology.org/2023.acl-long.61/>.
- 754
- 755 Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hervé Jégou, and Tomas Mikolov.
 Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016a.

- 756 Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient
 757 text classification. *arXiv preprint arXiv:1607.01759*, 2016b.
- 758
- 759 Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M.
 760 Khapra, and Pratyush Kumar. IndicNLPSuite: Monolingual corpora, evaluation benchmarks and
 761 pre-trained multilingual language models for Indian languages. In *Findings of the Association for
 762 Computational Linguistics: EMNLP 2020*, pp. 4948–4961, Online, November 2020. Association
 763 for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.445. URL <https://aclanthology.org/2020.findings-emnlp.445>.
- 764
- 765 Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schütze. GlotLID: Lan-
 766 guage identification for low-resource languages. In *The 2023 Conference on Empirical Methods
 767 in Natural Language Processing*, 2023. URL <https://openreview.net/forum?id=d14e3EBz5j>.
- 768
- 769 Amir Hossein Kargaran, François Yvon, and Hinrich Schütze. Glotscript: A resource and tool for low
 770 resource writing system identification. In *Proceedings of the 2024 Joint International Conference
 771 on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp.
 772 7774–7784, 2024.
- 773
- 774 Najoung Kim, Sebastian Schuster, and Shubham Toshniwal. Code pretraining improves entity tracking
 775 abilities of language models. *CoRR*, abs/2405.21068, 2024. doi: 10.48550/ARXIV.2405.21068.
 776 URL <https://doi.org/10.48550/arXiv.2405.21068>.
- 777
- 778 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International
 779 Conference for Learning Representations*, 2015.
- 780
- 781 Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Yacine Jernite, Margaret
 782 Mitchell, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von
 783 Werra, and Harm de Vries. The stack: 3 TB of permissively licensed source code. *Trans. Mach.
 784 Learn. Res.*, 2023, 2023. URL <https://openreview.net/forum?id=pxpbTdUEpD>.
- 785
- 786 Minato Kondo, Takehito Utsuro, and Masaaki Nagata. Enhancing translation accuracy of large
 787 language models through continual pre-training on parallel data, 2024. URL <https://arxiv.org/abs/2407.03145>.
- 788
- 789 Fajri Koto and Ikhwan Koto. Towards computational linguistics in Minangkabau language: Studies on
 790 sentiment analysis and machine translation. In *Proceedings of the 34th Pacific Asia Conference on
 791 Language, Information and Computation*, pp. 138–148, Hanoi, Vietnam, October 2020. Association
 792 for Computational Linguistics. URL <https://aclanthology.org/2020.pacific-1.17>.
- 793
- 794 Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi
 795 Stella, Ankur Bapna, and Orhan Firat. Madlad-400: A multilingual and document-level large
 796 audited dataset. *Advances in Neural Information Processing Systems*, 2024.
- 797
- 798 Anoop Kunchukuttan, Pratik Mehta, and Pushpak Bhattacharyya. The IIT Bombay English-Hindi par-
 799 allel corpus. In *Proceedings of the Eleventh International Conference on Language Resources and
 800 Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language Resources Association
 (ELRA). URL <https://aclanthology.org/L18-1548>.
- 801
- 802 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 803 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
 804 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
 805 Systems Principles*, 2023.
- 806
- 807 Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A.
 808 Rossi, and Thien Huu Nguyen. Okapi: Instruction-tuned large language models in multiple lan-
 809 guages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference
 810 on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 318–327,
 2023.

- 810 Wen Lai, Mohsen Mesgar, and Alexander Fraser. LLMs beyond English: Scaling the multilingual
 811 capability of LLMs with cross-lingual feedback. In Lun-Wei Ku, Andre Martins, and Vivek
 812 Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 8186–
 813 8213, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational
 814 Linguistics. URL <https://aclanthology.org/2024.findings-acl.488>.
- 815 Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral,
 816 Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen,
 817 et al. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. In *Thirty-sixth*
 818 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- 819
- 820 Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel
 821 Whitenack. Bloom library: Multimodal datasets in 300+ languages for a variety of downstream
 822 tasks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022*
 823 *Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi,*
 824 *United Arab Emirates, December 7-11, 2022*, pp. 8608–8621. Association for Computational
 825 Linguistics, 2022. URL <https://aclanthology.org/2022.emnlp-main.590>.
- 826 Hector Levesque, Ernest Davis, and Leora Morgenstern. The Winograd schema challenge. In
 827 *Thirteenth international conference on the principles of knowledge representation and reasoning*,
 828 2012.
- 829 Shenggui Li, Hongxin Liu, Zhengda Bian, Jiarui Fang, Haichen Huang, Yuliang Liu, Boxiang Wang,
 830 and Yang You. Colossal-ai: A unified deep learning system for large-scale parallel training. In
 831 *Proceedings of the 52nd International Conference on Parallel Processing*, pp. 766–775, 2023.
- 832
- 833 Yudong Li, Yuqing Zhang, Zhe Zhao, Linlin Shen, Weijie Liu, Weiquan Mao, and Hui Zhang. CSL: A
 834 large-scale chinese scientific literature dataset. In *Proceedings of the 29th International Conference*
 835 *on Computational Linguistics*, pp. 3917–3923, 2022.
- 836 Evenki Life. Evenki life newspaper, 2014. URL https://drive.google.com/file/d/1he2q6RncA_NKHP1IJjSzlkK-2qgEFTiCG/view.
- 837
- 838 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization*
 839 *Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
 840 URL <https://aclanthology.org/W04-1013>.
- 841
- 842 Peiqin Lin, Shaoxiong Ji, Jörg Tiedemann, André FT Martins, and Hinrich Schütze. MaLA-500:
 843 Massive language adaptation of large language models. *arXiv preprint arXiv:2401.13303*, 2024.
- 844
- 845 Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuhui Chen, Daniel Simig, Myle Ott,
 846 Naman Goyal, Shruti Bhosale, Jingfei Du, et al. Few-shot learning with multilingual generative
 847 language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural*
 848 *Language Processing*, pp. 9019–9052, 2022.
- 849
- 850 Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. S2ORC: The semantic
 851 scholar open research corpus. In *Proceedings of the 58th Annual Meeting of the Association*
 852 *for Computational Linguistics*, pp. 4969–4983, Online, July 2020. Association for Computational
 853 Linguistics. doi: 10.18653/v1/2020.acl-main.447. URL <https://www.aclweb.org/anthology/2020.acl-main.447>.
- 854
- 855 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane
 856 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The
 857 next generation. *arXiv preprint arXiv:2402.19173*, 2024.
- 858
- 859 Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. LLaMAX: Scaling linguistic hori-
 860 zons of LLM by enhancing translation capabilities beyond 100 languages. *arXiv preprint*
 861 *arXiv:2407.05975*, 2024.
- 862
- 863 Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catas-
 864 troptic forgetting in large language models during continual fine-tuning. *CoRR*, abs/2308.08747,
 865 2023. doi: 10.48550/ARXIV.2308.08747. URL <https://doi.org/10.48550/arXiv.2308.08747>.

- 864 Chunlan Ma, Ayyoob ImaniGooghari, Haotian Ye, Ehsaneddin Asgari, and Hinrich Schütze. Taxi1500:
 865 A multilingual dataset for text classification in 1500 languages, 2023.
 866
- 867 Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. At
 868 which training stage does code data help llms reasoning? In *The Twelfth International Conference*
 869 *on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
 870 URL <https://openreview.net/forum?id=KIPJKST4gw>.
- 871 Martin Majliš. W2C – web to corpus – corpora, 2011. URL [http://hdl.handle.net/11858/
 872 00-097C-0000-0022-6133-9](http://hdl.handle.net/11858/00-097C-0000-0022-6133-9). LINDAT/CLARIAH-CZ digital library at the Institute of
 873 Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
 874
- 875 Masakhane. Lacuna project, 2023. URL https://github.com/masakhane-io/lacuna_pos_ner.
 876
- 877 Thomas Mayer and Michael Cysouw. Creating a massively parallel bible corpus. In Nicoletta
 878 Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani,
 879 Asunción Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Ninth International*
 880 *Conference on Language Resources and Evaluation, LREC 2014, Reykjavik, Iceland, May 26-*
 881 *31, 2014*, pp. 3158–3163. European Language Resources Association (ELRA), 2014. URL
 882 <http://www.lrec-conf.org/proceedings/lrec2014/summaries/220.html>.
 883
- 884 Jamshidbek Mirzakhalov, Anoop Babu, Duygu Ataman, Sherzod Kariev, Francis Tyers, Otabek Abduraufov, Mammad Hajili, Sardana Ivanova, Abror Khaytbaev, Antonio Laverghetta Jr., Bekhzodbek Moydinboyev, Esra Onal, Shaxnoza Pulatova, Ahsan Wahab, Orhan Firat, and Sriram Chellappan. A large-scale study of machine translation in Turkic languages. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5876–5890, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.475. URL <https://aclanthology.org/2021.emnlp-main.475>.
 885
- 886 Steven Moran, Christian Bentz, Ximena Gutierrez-Vasques, Olga Pelloni, and Tanja Samardzic. TeDDi sample: Text data diversity sample for language comparison and multilingual NLP. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 1150–1158, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.123>.
 887
- 888 Makoto Morishita, Jun Suzuki, and Masaaki Nagata. JParaCrawl: A large scale web-based English-Japanese parallel corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 3603–3609, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.443>.
 889
- 890 Nasrin Mostafazadeh, Michael Roth, Annie Louis, Nathanael Chambers, and James Allen. Lsdsem 2017 shared task: The story cloze test. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*, 2017.
 891
- 892 Chenghao Mou, Chris Ha, Kenneth Enevoldsen, and Peiyuan Liu. Chenghaomou/text-dedup: Reference snapshot, 2023. URL <https://doi.org/10.5281/zenodo.8364980>.
 893
- 894 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. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*, 2022.
 895
- 896 Jonathan Mukibi, Andrew Katumba, Joyce Nakatumba-Nabende, Ali Hussein, and Joshua Meyer. The makerere radio speech corpus: A luganda radio corpus for automatic speech recognition. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 1945–1954, 2022.
 897
- 898 Taishi Nakamura, Mayank Mishra, Simone Tedeschi, Yekun Chai, Jason T Stillerman, Felix Friedrich, Prateek Yadav, Tanmay Laud, Vu Minh Chien, Terry Yue Zhuo, et al. Aurora-m: The first open source multilingual language model red-teamed according to the us executive order. *arXiv preprint arXiv:2404.00399*, 2024.
 899

- 918 Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya
 919 Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, Ondřej Bojar, Chenhui Chu,
 920 Akiko Eriguchi, Kaori Abe, Yusuke Oda, and Sadao Kurohashi. Overview of the 8th workshop on
 921 Asian translation. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pp. 1–45,
 922 Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.wat-1.1.
 923 URL <https://aclanthology.org/2021.wat-1.1>.
- 924 Toshiaki Nakazawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida,
 925 Anoop Kunchukuttan, Makoto Morishita, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori
 926 Abe, Yusuke Oda, and Sadao Kurohashi. Overview of the 9th workshop on Asian translation.
 927 In *Proceedings of the 9th Workshop on Asian Translation*, pp. 1–36, Gyeongju, Republic of
 928 Korea, October 2022. International Conference on Computational Linguistics. URL <https://aclanthology.org/2022.wat-1.1>.
- 930 Graham Neubig. The Kyoto free translation task. <http://www.phontron.com/kftt>, 2011.
- 931 Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu Man, Nghia Trung Ngo, Franck Dernoncourt,
 932 Ryan A. Rossi, and Thien Huu Nguyen. Culturax: A cleaned, enormous, and multilingual dataset
 933 for large language models in 167 languages, 2023.
- 934 Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. Small data? no problem! exploring the viability
 935 of pretrained multilingual language models for low-resourced languages. In *Proceedings
 936 of the 1st Workshop on Multilingual Representation Learning*, pp. 116–126, Punta Cana, Do-
 937 minican Republic, November 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.mrl-1.11>.
- 938 OSCAR. Oscar (open super-large crawled aggregated corpus) 2301, 2023. URL <https://huggingface.co/datasets/oscar-corpus/OSCAR-2301>.
- 939 Chester Palen-Michel, June Kim, and Constantine Lignos. Multilingual open text release 1: Public
 940 domain news in 44 languages. In *Proceedings of the Thirteenth Language Resources and Evalua-
 941 tion Conference*, pp. 2080–2089, Marseille, France, June 2022. European Language Resources
 942 Association. URL <https://aclanthology.org/2022.lrec-1.224>.
- 943 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
 944 evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin (eds.),
 945 *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp.
 946 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.
 947 doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 948 Indraneil Paul, Goran Glavas, and Iryna Gurevych. Ircoder: Intermediate representations make
 949 language models robust multilingual code generators. In Lun-Wei Ku, Andre Martins, and
 950 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Com-
 951 putational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11–
 952 16, 2024, pp. 15023–15041. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.acl-long.802>.
- 953 Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen.
 954 XCOPA: A multilingual dataset for causal commonsense reasoning. In *Proceedings of the 2020
 955 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- 956 Maja Popović. chrF: character n-gram F-score for automatic MT evaluation. In Ondřej Bojar,
 957 Rajan Chatterjee, Christian Federmann, Barry Haddow, Chris Hokamp, Matthias Huck, Varvara
 958 Logacheva, and Pavel Pecina (eds.), *Proceedings of the Tenth Workshop on Statistical Machine
 959 Translation*, pp. 392–395, Lisbon, Portugal, September 2015. Association for Computational Lin-
 960 guistics. doi: 10.18653/v1/W15-3049. URL <https://aclanthology.org/W15-3049>.
- 961 Matt Post. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference
 962 on Machine Translation: Research Papers*, pp. 186–191, Belgium, Brussels, October 2018. As-
 963 sociation for Computational Linguistics. URL <https://www.aclweb.org/anthology/W18-6319>.

- 972 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
 973 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
 974 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 975
- 976 Anand Rajaraman and Jeffrey D Ullman. *Mining of massive datasets*. Autoedicion, 2011.
- 977
- 978 Ronald Rivest. The md5 message-digest algorithm, 1992.
- 979
- 980 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
 981 Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. Code llama: Open foundation models for code.
 982 *arXiv preprint arXiv:2308.12950*, 2023.
- 983
- 984 Roberts Rozis and Raivis Skadiņš. Tilde MODEL - multilingual open data for EU languages.
 985 In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pp. 263–265,
 986 Gothenburg, Sweden, May 2017. Association for Computational Linguistics. URL <https://aclanthology.org/W17-0235>.
- 987
- 988 Hassan Sajjad, Ahmed Abdelali, Nadir Durrani, and Fahim Dalvi. AraBench: Benchmarking dialec-
 989 tal Arabic-English machine translation. In *Proceedings of the 28th International Conference on*
 990 *Computational Linguistics*, pp. 5094–5107, Barcelona, Spain (Online), December 2020. Interna-
 991 tional Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.447. URL
 992 <https://aclanthology.org/2020.coling-main.447>.
- 993
- 994 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman
 995 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. BLOOM: A 176B-
 996 parameter open-access multilingual language model. *arXiv preprint*, 2022.
- 997
- 998 Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. WikiMatrix:
 999 Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In *Proceedings of the*
 1000 *16th Conference of the European Chapter of the Association for Computational Linguistics: Main*
 1001 *Volume*, pp. 1351–1361, Online, April 2021a. Association for Computational Linguistics. doi: 10.
 1002 18653/v1/2021.eacl-main.115. URL <https://aclanthology.org/2021.eacl-main.115>.
- 1003
- 1004 Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela
 1005 Fan. CCMATRIX: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the*
 1006 *59th Annual Meeting of the Association for Computational Linguistics and the 11th International*
 1007 *Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 6490–6500,
 1008 Online, 2021b. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.507.
 1009 URL <https://aclanthology.org/2021.acl-long.507>.
- 1010
- 1011 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
 1012 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language
 1013 models are multilingual chain-of-thought reasoners, 2022.
- 1014
- 1015 Oleh Shliazhko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and
 1016 Tatiana Shavrina. mGPT: Few-shot learners go multilingual. *arXiv preprint arXiv:2204.07580*,
 1017 2022.
- 1018
- 1019 Anil Kumar Singh. Named entity recognition for south and south East Asian languages: Taking stock.
 1020 In *Proceedings of the IJCNLP-08 Workshop on Named Entity Recognition for South and South*
 1021 *East Asian Languages*, 2008. URL <https://aclanthology.org/I08-5003>.
- 1022
- 1023 Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F Karlsson, Abinaya Mahendiran, Wei-Yin
 1024 Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, et al. Aya dataset: An
 1025 open-access collection for multilingual instruction tuning. In *Proceedings of the 62nd Annual*
 1026 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11521–
 1027 11567, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL
 1028 <https://aclanthology.org/2024.acl-long.620>.
- 1029
- 1030 Luca Soldaini and Kyle Lo. peS2o (Pretraining Efficiently on S2ORC) Dataset. Technical report,
 1031 Allen Institute for AI, 2023. ODC-By, <https://github.com/allenai/pes2o>.

- 1026 Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. Asynchronous pipeline for processing
 1027 huge corpora on medium to low resource infrastructures. In *7th Workshop on the Challenges in the*
 1028 *Management of Large Corpora (CMLC-7)*. Leibniz-Institut für Deutsche Sprache, 2019.
- 1029
- 1030 Marc Szafraniec, Baptiste Rozière, Hugh Leather, Patrick Labatut, François Charton, and Gabriel
 1031 Synnaeve. Code translation with compiler representations. In *The Eleventh International Conference*
 1032 *on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
 1033 2023. URL <https://openreview.net/forum?id=XomEU3eNeSQ>.
- 1034 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia,
 1035 Andrew Poult, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science.
 1036 *arXiv preprint arXiv:2211.09085*, 2022.
- 1037
- 1038 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya
 1039 Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al.
 1040 Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*,
 1041 2024.
- 1042 Atula Tejaswi, Nilesh Gupta, and Eunsol Choi. Exploring design choices for building language-
 1043 specific llms, 2024. URL <https://arxiv.org/abs/2406.14670>.
- 1044
- 1045 Huu Nguyen Thuat Nguyen and Thien Nguyen. Culturay: A large cleaned multilingual dataset of 75
 1046 languages, 2024.
- 1047 Jörg Tiedemann. Parallel data, tools and interfaces in OPUS. In Nicoletta Calzolari, Khalid Choukri,
 1048 Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion Moreno,
 1049 Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Eighth International Conference on*
 1050 *Language Resources and Evaluation (LREC'12)*, pp. 2214–2218, Istanbul, Turkey, May 2012.
 1051 European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2012/pdf/463_Paper.pdf.
- 1052
- 1053 Jörg Tiedemann. The Tatoeba Translation Challenge – Realistic data sets for low resource and
 1054 multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pp. 1174–
 1055 1182, Online, November 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.wmt-1.139>.
- 1056
- 1057 Alexey Tikhonov and Max Ryabinin. It's All in the Heads: Using Attention Heads as a Baseline
 1058 for Cross-Lingual Transfer in Commonsense Reasoning. In *Findings of the Association for*
 1059 *Computational Linguistics: ACL-IJCNLP 2021*, 2021.
- 1060
- 1061 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
 1062 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
 1063 and fine-tuned chat models. *arXiv preprint*, 2023.
- 1064
- 1065 Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude,
 1066 Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. Aya model: An instruction
 1067 finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*, 2024.
- 1068
- 1069 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
 1070 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
 1071 *processing systems*, pp. 5998–6008, 2017.
- 1072
- 1073 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán,
 1074 Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets
 1075 from web crawl data. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri,
 1076 Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani,
 1077 Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the*
 1078 *Twelfth Language Resources and Evaluation Conference*, pp. 4003–4012, Marseille, France, May
 1079 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.494>.
- Wikimedia Foundation. Wikimedia downloads, n.d. URL <https://dumps.wikimedia.org>.

- 1080 Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, X. Li, Zhi Yuan Lim,
 1081 S. Soleman, R. Mahendra, Pascale Fung, Syafri Bahar, and A. Purwarianti. Indonlu: Benchmark
 1082 and resources for evaluating indonesian natural language understanding. In *Proceedings of the 1st*
 1083 *Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the*
 1084 *10th International Joint Conference on Natural Language Processing*, 2020.
- 1085 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
 1086 Pierrick Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:
 1087 State-of-the-art natural language processing. *arXiv preprint*, 2019.
- 1088 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya
 1089 Barua, and Colin Raffel. mT5: A massively multilingual pre-trained text-to-text transformer.
 1090 In *Proceedings of the 2021 Conference of the North American Chapter of the Association for*
 1091 *Computational Linguistics: Human Language Technologies*, pp. 483–498, 2021.
- 1092 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
 1093 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
 1094 *arXiv:2407.10671*, 2024a.
- 1095 Ke Yang, Jiateng Liu, John Wu, Chaoqi Yang, Yi R. Fung, Sha Li, Zixuan Huang, Xu Cao, Xingyao
 1096 Wang, Yiquan Wang, Heng Ji, and Chengxiang Zhai. If LLM is the wizard, then code is the wand:
 1097 A survey on how code empowers large language models to serve as intelligent agents. *CoRR*,
 1098 abs/2401.00812, 2024b. doi: 10.48550/ARXIV.2401.00812. URL <https://doi.org/10.48550/arXiv.2401.00812>.
- 1099 Rodolfo Zevallos, John Ortega, William Chen, Richard Castro, Núria Bel, Cesar Toshio, Renzo
 1100 Venturas, Hilario Aradiel, and Nelsi Melgarejo. Introducing QuBERT: A large monolingual
 1101 corpus and BERT model for Southern Quechua. In *Proceedings of the Third Workshop on*
 1102 *Deep Learning for Low-Resource Natural Language Processing*, pp. 1–13, Hybrid, July 2022.
 1103 Association for Computational Linguistics. doi: 10.18653/v1/2022.deeplo-1.1. URL <https://aclanthology.org/2022.deeplo-1.1>.
- 1104 Chen Zhang, Mingxu Tao, Quzhe Huang, Jiaheng Lin, Zhibin Chen, and Yansong Feng. Mc²:
 1105 Towards transparent and culturally-aware nlp for minority languages in china. *arXiv preprint*
 1106 *arXiv:2311.08348*, 2024a.
- 1107 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating
 1108 text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- 1109 Xinlu Zhang, Zhiyu Zoey Chen, Xi Ye, Xianjun Yang, Lichang Chen, William Yang Wang, and
 1110 Linda Ruth Petzold. Unveiling the impact of coding data instruction fine-tuning on large language
 1111 models reasoning. *CoRR*, abs/2405.20535, 2024b. doi: 10.48550/ARXIV.2405.20535. URL
 1112 <https://doi.org/10.48550/arXiv.2405.20535>.
- 1113 Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Texxygen:
 1114 A benchmarking platform for text generation models. In *The 41st international ACM SIGIR*
 1115 *conference on research & development in information retrieval*, pp. 1097–1100, 2018.
- 1116
- ## A DATA SOURCES
- 1117
- ### A.1 MONOLINGUAL DATA
- 1118
- Table 7 lists the corpora and collections we use as monolingual data sources in this work. Monolingual
 1119 data sources simply contain text data in a single language.
- 1120
- Metadata** For the monolingual data sources, we define the contents of the JSONL output file
 1121 from the pre-processing workflow to consist of the fields url, text, collection, source, and
 1122 original_code. The contents of these fields are as follows. The field text contains the language
 1123 data in a granularity specific to the given corpus. If the granularity was sentence-level, then we could
 1124 expect the sentences in the corpus to generally be independent of each other, while only parts of the

sentences—such as phrases, clauses, and words—exhibit serial dependence. If the granularity was paragraph-level, then we could expect the sentences within paragraphs to have serial dependence, while paragraphs to largely be independent of each other. The field `url` contains a URL indicating the web address from which the text data has been extracted, if available. The field `collection` contains the name of the collection, i.e., a corpus or a set of corpora, which the text is extracted from, whereas the field `source` contains the name of a more specific part of the collection, such as the name of an individual corpus or a file in the collection the text was extracted from. Lastly, the field `original_code` contains the language code of the text data as it is designated in the data source, e.g., in the directory structure, the filenames, or the data object returned by an API call.

Table 7: Datasets used as monolingual source data.

Name	Languages	Domains	URL	Year
AfriBERTa (Ogueji et al., 2021)	10	news	https://huggingface.co/datasets/castorini/afriberita-corpus	2021
Bloom library (Leong et al., 2022)	363	religious, books	https://huggingface.co/datasets/bloom/bloom-1m	2022
CC100 (Conneau et al., 2020)	100	web	https://huggingface.co/datasets/c0100	2020
CulturaX (Nguyen et al., 2023)	167	web	https://huggingface.co/datasets/unipr/CulturaX	2023
CulturaY (Thuat Nguyen & Nguyen, 2024)	75	web	https://huggingface.co/datasets/onotcord/CulturaY	2024
curse-of-multilinguality (Chang et al., 2023)	200	misc	https://github.com/tylerachang/curse-of-multilinguality	2023
Evenki Life (Life, 2014)	1	newspapers	https://drive.google.com/file/d/1he2qRncA_NKHP1JjSz1kk-2qgEFTicG/view	2014
Glot500 (Imani et al., 2023)	511	misc	https://huggingface.co/datasets/cis-lmu/Glot500	2023
GlotSparse (Kargaran et al., 2023)	10	news	https://huggingface.co/datasets/cis-lmu/GlotSparse	2023
HPLT v1.2 (de Gibert et al., 2024)	75	web	https://hplt-project.org/datasets/v1.2	2024
Indigenous Languages Corpora (EdTeKLA, 2022)	1	UNK	https://github.com/EdTeKLA/IndigenousLanguages_Corpora	2022
Indo4B (Wilie et al., 2020)	1	misc	https://github.com/IndoNLB/indonlu?tab=readme-ov-file	2020
Lacuna Project (Masakhane, 2023)	20	UNK	https://github.com/masakhane-io/lacuna_pos_ner	2023
Languages of Russia (Corpora and Tools, n.d.)	46	social media, web	https://zenodo.org/record/585101lang/download	UNK
MADLAD-400 (Kudugunta et al., 2024)	419	web	https://huggingface.co/datasets/allenai/MADLAD-400	2023
Makerere Radio Speech Corpus (Mukiibi et al., 2022)	1	transcription	https://zenodo.org/record/585017	2022
masakhane-ner1.0 (Ifeoluwa Adelani et al., 2021)	12	UNK	https://github.com/masakhane-io/masakhane-ner	2021
MC2 (Zhang et al., 2024a)	4	web	https://huggingface.co/datasets/pkupie/mc2_corpus	2024
mC4 (Raffel et al., 2020)	101	web	https://huggingface.co/datasets/allenai/c4	2020
multilingual-data-peru (Grupo de Inteligencia Artificial PUCP, n.d.)	4	UNK	https://github.com/laipucp/multilingual-data-peru	2020
OSCAR 2301 (OSCAR, 2023)	152	web	https://huggingface.co/datasets/oscar-corpus/OSCAR-2301	2023
The Leipzig Corpora (Goldhahn et al., 2012)	136	newspapers, wikipedia	https://huggingface.co/datasets/imvladion/leipzig_corpora_collection	2012
Tigrinya Language Modeling (Gaim et al., 2021)	1	news, blogs, books	https://zenodo.org/record/5139094	2021
Wikipedia 20231101 (Wikimedia Foundation, n.d.)	323	wikipedia	https://huggingface.co/datasets/wikimedia/wikipedia	2023
Wikisource 20231201 (Wikimedia Foundation, n.d.)	73	books	https://huggingface.co/datasets/wikimedia/wikisource	2023
Tatoeba challenge monolingual (Tiedemann, 2020)	280	wikimedia	https://github.com/Helsinki-NLP/Tatoeba-Challenge/blob/master/data/MonolingualData-v2020-07-28.md	2020

A.1.1 LIST OF DATA SOURCES

The major ones include AfriBERTa (Ogueji et al., 2021), Bloom library (Leong et al., 2022), CC100 (Conneau et al., 2020; Wenzek et al., 2020) CulturaX (Nguyen et al., 2023), CulturaY (Thuat Nguyen & Nguyen, 2024), the Curse of Multilinguality (Chang et al., 2023), Evenki Life (Life, 2014), Glot500 (Imani et al., 2023), GlotSparse (Kargaran et al., 2023), monoHPLT of HPLT v1.2 (de Gibert et al., 2024) from the HPLT project (Aulamo et al., 2023), Indigenous Languages Corpora (EdTeKLA, 2022), Indo4B (Wilie et al., 2020), Lacuna Project (Masakhane, 2023), Languages of Russia (Corpora and Tools, n.d.), MADLAD-400 (Kudugunta et al., 2024), Makerere Radio Speech Corpus (Mukiibi et al., 2022), masakhane-ner1.0 (Ifeoluwa Adelani et al., 2021), MC2 (Zhang et al., 2024a), mC4 (Raffel et al., 2020), multilingual-data-peru (Grupo de Inteligencia Artificial PUCP, n.d.), OSCAR 2301 (OSCAR, 2023) from the OSCAR project¹⁷, Tatoeba challenge monolingual collection (Tiedemann, 2020), The Leipzig Corpora (Goldhahn et al., 2012), Tigrinya Language Modeling (Gaim et al., 2021), Wikipedia 20231101 (Wikimedia Foundation, n.d.) and Wikisource 20231201 (Wikimedia Foundation, n.d.). We exclude high-resource languages in CulturaX, HPLT, MADLAD-400, CC100, mC4, and OSCAR 2301. We exclude Gahuza and Pidgin in the AfriBERTa dataset. We filter out texts that mainly contain a date or timestamp in the Languages of Russia dataset. For Glot500-c, we filter out texts, which may come from train or test tests from datasets for machine translation, such as Flores200, Tatoeba, and mtdata. Despite the translation data being split into source and target languages in the Glot500-c corpus, we decide to filter them to avoid potential data leakage, especially since we use Flores200 as the evaluation benchmark.

Glot500-c uses the following datasets: AI4Bharat,¹⁸ AIFORTHAI-LotusCorpus,¹⁹ Add (El-Haj et al., 2018), AfriBERTa (Ogueji et al., 2021), AfroMAFT (Adelani et al., 2022; Xue et al., 2021), Anuvaad,²⁰ AraBench (Sajjad et al., 2020), AUTSHUMATO,²¹ Bloom (Leong et al., 2022), CC100

¹⁷<https://oscar-project.org/>

¹⁸<https://ai4bharat.org/>

¹⁹<https://github.com/korakot/corpus/releases/download/v1.0/AIFORTHAI-LotusCorpus.zip>

²⁰<https://github.com/project-anuvaad/anuvaad-parallel-corpus>

²¹<https://autshumato.sourceforge.net/>

(Conneau et al., 2020), CCNet (Wenzek et al., 2020), CMU_Haitian_Creole,²² CORP.NCHLT,²³ Clarin,²⁴ DART (Alsarsour et al., 2018), Earthlings (Dunn, 2020), FFR,²⁵ Flores200 (Costa-jussà et al., 2022), GiossaMedia (Góngora et al., 2022; 2021), Glosses (Camacho-Collados et al., 2016), Habibi (El-Haj, 2020), HinDialect (Bafna, 2022), HornMT,²⁶ IITB (Kunchukuttan et al., 2018), IndicNLP (Nakazawa et al., 2021), Indiccorp (Kakwani et al., 2020), isiZulu,²⁷ JParaCrawl (Morishita et al., 2020), KinyaSMT,²⁸ LeipzigData (Goldhahn et al., 2012), Lindat,²⁹ Lingala_Song_Lyrics,³⁰ Lyrics,³¹ MC4 (Raffel et al., 2020), MTData (Gowda et al., 2021), MaCoCu (Bañón et al., 2022), Makerere MT Corpus,³² Masakhane community,³³ Mburisano_Covid,³⁴ Menyo20K (Adelani et al., 2021), Minangkabau corpora (Koto & Koto, 2020), MoT (Palen-Michel et al., 2022), NLLB_seed (Costa-jussà et al., 2022), Nart/abkhaz,³⁵ OPUS (Tiedemann, 2012), OSCAR (Suárez et al., 2019), ParaCrawl (Bañón et al., 2020), Parallel Corpora for Ethiopian Languages (Abate et al., 2018), Phontron (Neubig, 2011), QADI (Abdelali et al., 2021), Quechua-IIC (Zevallos et al., 2022), SLI_GalWeb.1.0 (Agerri et al., 2018), Shami (Abu Kwaik et al., 2018), Stanford NLP,³⁶ StatMT,³⁷ TICO (Anastasopoulos et al., 2020), TIL (Mirzakhalov et al., 2021), Tatoeba,³⁸ TeDDi (Moran et al., 2022), Tilde (Rozis & Skadiqš, 2017), W2C (Majliš, 2011), WAT (Nakazawa et al., 2022), WikiMatrix (Schwenk et al., 2021a), Wikipedia,³⁹ Workshop on NER for South and South East Asian Languages (Singh, 2008), XLSum (Hasan et al., 2021). We filter out Flores200 when processing the Glot500-c dataset. Glot500-c includes texts in languages, i.e., Azerbaijani, Gujarati, Igbo, Oromo, Rundi, Tigrinya and Yoruba, from the XLSum dataset.

A.2 CODE

All the programming language splits are filtered using the following conditions:

- For files forked more than 25 times, we retain them if the average line length is less than 120, the maximum line length is less than 300, and the alphanumeric fraction is more than 30%.
- For files forked between 15 and 25 times, we retain them if the average line length is less than 90, the maximum line length is less than 150, and the alphanumeric fraction is more than 40%.
- For files forked less than 15 times, we retain them if the average line length is less than 80, the maximum line length is less than 120, and the alphanumeric fraction is more than 45%.

Subsequently, an aggressive MinHash deduplication pipeline with a threshold of 0.5 and a shingle size of 20 is applied. Finally, the resultant language splits are then capped at 5 million samples each.

A.3 DATA DEDUPLICATION

As we collect data from different sources, we deduplicate the data to remove the overlap between different sources.

²²<http://www.speech.cs.cmu.edu/haitian/text/>

²³<https://repo.sadilar.org/handle/20.500.12185/7>

²⁴<https://www.clarin.si/>

²⁵<https://github.com/bonaventuredossou/ffr-v1/tree/master/FFR-Dataset>

²⁶<https://github.com/asmelashteka/HornMT>

²⁷<https://zenodo.org/record/5035171>

²⁸<https://github.com/pniyongabo/kinyarwandaSMT>

²⁹<https://lindat.cz/faq-repository>

³⁰https://github.com/espoirMur/songs_lyrics_webscrap

³¹<https://lyricstranslate.com/>

³²<https://zenodo.org/record/5089560>

³³<https://github.com/masakhane-io/masakhane-community>

³⁴<https://repo.sadilar.org/handle/20.500.12185/536>

³⁵https://huggingface.co/datasets/Nart/abkhaz_text

³⁶<https://nlp.stanford.edu/>

³⁷<https://statmt.org/>

³⁸<https://tatoeba.org/en/>

³⁹<https://huggingface.co/datasets/wikipedia>

1242 **MinHash Deduplication** For each language’s dataset in each writing system, we start by using the
 1243 MinHashLSH algorithm (Rajaraman & Ullman, 2011) to filter out similar documents. It is a near-
 1244 deduplication technique that builds on MinHash (Broder, 1997), utilizing multiple hash functions for
 1245 n-grams and the Jaccard similarity, and incorporates Locality-Sensitive Hashing to enhance efficiency.
 1246 We use the implementation by `text-dedup` repository (Mou et al., 2023), applying 5-grams and a
 1247 similarity threshold of 0.7 to identify similar documents based on the Jaccard index.

1248
 1249 **Exact Deduplication** We further deploy exact deduplication using the `text-dedup` repository
 1250 again with precise matching. It takes each document through a hash function, i.e., MD5 (Rivest,
 1251 1992) in our choice, and the hash values of all documents are compared to identify duplicates.

1252

1253 B ADDITIONAL STATISTICS OF MALA CORPUS

1254

1255 B.1 SUPPORTED LANGUAGES

1256 Table 8 shows the languages codes of MaLA corpus, where “unseen” means the languages are not used
 1257 for training EMMA-500. The classification system for token counts categorizes language resources
 1258 based on their size into five distinct tiers: “high” for resources exceeding 1 billion tokens, indicating
 1259 a vast amount of data; “medium-high” for those with more than 100 million tokens, reflecting a
 1260 substantial dataset; “medium” for resources that contain over 10 million tokens, representing a
 1261 moderate size; “medium-low” for datasets with over “1 million tokens”, indicating a smaller yet
 1262 significant amount of data; and finally, “low” for resources containing less than 1 million tokens,
 1263 which suggests a minimal data presence. This hierarchy helps in understanding the scale and potential
 1264 utility of the language resources available. Figure 1 shows the number of texts and tokens in different
 1265 resource groups.

1266

1267 B.2 DATA ANALYSIS

1268

1269 We examine the Unicode block distribution of each language, which counts the percentage of tokens
 1270 falling into the Unicode block of each language. This aims to check whether language code conversion
 1271 and writing system recognition are reasonably good. Figure 2 shows the Unicode block distribution.
 1272 The result aligns with our observation of the presence of code-mixing, but in general, the majority of
 1273 languages have tokens falling in their own Unicode blocks.

1274

1275 We also check the data source distribution as shown in Figure 3 to see where the texts in the MaLA
 1276 corpus come from. The main source is common crawl, e.g., CC 2018, CC, OSCAR, and CC 20220801.
 1277 A large number of documents come from Earthlings which comes from Glot500-c (Imani et al., 2023).
 1278 For web-crawled data with a URL in the original metadata of corpora like CulturaX (Nguyen et al.,
 1279 2023) and HPLT (de Gibert et al., 2024), we extract the domain. Thus, the final corpus has many
 1280 sources with a small portion.

1281

1282 C EVALUATION BENCHMARKS

1283

1284 C.1 POLYWRITE

1285

1286 The PolyWrite dataset has 51 writing tasks with the number of prompts per task shown in Figure 4.
 1287 We use Google Translate to translate the English prompts into 240 languages and back-translate for
 1288 translation quality assessment.

1289

1290 The metadata of PolyWrite includes several key fields. The `category` specifies the task type. The
 1291 name field typically holds the specific identifier for each prompt, while `prompt_en` contains the
 1292 English version of the prompt. `lang_script` identifies the language and script used, ensuring
 1293 correct language processing. The `prompt_translated` field holds the translated prompt in the
 1294 target language, and `prompt_backtranslated` contains the back-translated version to assess
 1295 translation quality. Both `bleu` and `chrf++` fields provide numeric evaluation metrics, with BLEU
 1296 and chrf++ scores measuring the quality of the generated text. Finally, the `uuid` ensures a unique
 1297 identifier for each dataset entry, allowing for precise reference and tracking of individual prompts. To

Table 8: Languages by resource groups

Category	Languages	Language Codes
high	27	fra.Latn, mon.Cyril, kat.Geor, tgk.Cyril, kaz.Cyril, glg.Latn, hbs.Latn, kan.Knda, mal.Mlym, rus.Cyril, cat.Latn, hye.Armen, guj.Gujr, slv.Latn, fil.Latn, bel.Cyril, isl.Latn, nep.Deva, mlt.Latn, pan.Guru, afr.Latn, urd.Arab, mkd.Cyril, aze.Latn, deu.Latn, eng.Latn, ind.Latn
low	210	prs.Arab, nqo.Nkoo, emp.Latn, pfl.Latn, teo.Latn, gpe.Latn, izz.Latn, shin.Mymr, hak.Latn, pls.Latn, evn.Cyril, djk.Latn, toj.Latn, nog.Cyril, ctu.Latn, tea.Latn, jiv.Latn, ach.Latn, mrj.Latn, aip.Arab, apc.Arab, tab.Cyril, hvn.Latn, ts.Latn, bak.Latn, ndc.Latn, trv.Latn, top.Latn, kjh.Cyril, guh.Latn, mnj.Mtei, csy.Latn, noa.Latn, dov.Latn, bho.Deva, kon.Latn, hne.Deva, keg.Latn, mni.Beng, hus.Latn, pau.Latn, jbo.Latn, dtp.Latn, kmb.Latn, hau.Arab, pdc.Latn, nch.Latn, acf.Latn, bim.Latn, ix.Latn, dyt.Deva, kas.Arabs, irc.Arabs, alz.Latn, lez.Cyril, lld.Latn, tdt.Latn, acm.Arabs, bih.Deva, mzh.Latn, guw.Latn, rop.Latn, rwo.Latn, ahk.Latn, qub.Latn, kri.Latn, gub.Latn, laj.Latn, sxn.Latn, luo.Latn, tly.Latn, pwn.Latn, mag.Deva, xav.Latn, bum.Latn, ubu.Latn, roa.Latn, math.Latn, tsq.Latn, ger.Latn, arn.Latn, csh.Latn, guc.Latn, bat.Latn, knj.Latn, cre.Latn, bus.Latn, anp.Deva, aln.Latn, nab.Latn, zai.Latn, kpv.Cyril, eng.Latn, gvl.Latn, wal.Latn, fiu.Latn, swh.Latn, crh.Latn, nia.Latn, bqc.Latn, map.Latn, atj.Latn, npi.Deva, bru.Latn, din.Latn, pis.Latn, gur.Latn, cuk.Latn, zne.Latn, cdo.Latn, lhu.Latn, pcd.Latn, mas.Latn, bis.Latn, ncj.Latn, ibb.Latn, tay.Latn, bts.Latn, tzj.Latn, bzz.Latn, cce.Cyril, jvn.Latn, ndo.Latn, rug.Latn, koi.Cyril, mco.Latn, fat.Latn, olo.Latn, inb.Latn, mkn.Latn, qvi.Latn, mak.Latn, itn.Latn, nrm.Latn, kua.Latn, san.San, nbl.Latn, kik.Latn, dyu.Latn, sgs.Latn, msm.Latn, mnw.Latn, zha.Latn, sja.Latn, xal.Cyril, rmc.Latn, ami.Latn, sda.Latn, tdx.Latn, yap.Latn, tzh.Latn, sus.Latn, ikk.Latn, bas.Latn, nde.Latn, dsb.Latn, seh.Latn, knv.Latn, amu.Latn, dwr.Latn, cyr.Latn, uig.Latn, bsr.Cyril, tcy.Knda, mau.Latn, aoj.Latn, gor.Latn, cha.Latn, fip.Latn, chr.Cher, mdv.Cyril, arb.Arab, quw.Latn, shp.Latn, spp.Latn, frp.Latn, ape.Latn, cbk.Latn, mnw.Mymr, mfe.Latn, jam.Latn, lad.Latn, awa.Deva, mad.Latn, ote.Latn, shi.Latn, btx.Latn, maz.Latn, ppk.Latn, smn.Latn, twu.Latn, blk.Mymr, msi.Latn, naq.Latn, tly.Arabs, wuu.Hani, mos.Latn, cab.Latn, zlm.Latn, gag.Latn, suz.Deva, ksw.Mymr, gug.Latn, nij.Latn, nov.Latn, srm.Latn, jae.Latn, nyu.Latn, yom.Latn, gui.Latn
medium	68	tha.Thai, kat.Latn, lim.Latn, tgk.Arabs, che.Cyril, lav.Latn, xho.Latn, war.Latn, nan.Latn, grc.Grek, orm.Latn, zsm.Latn, cnh.Latn, yor.Latn, arg.Latn, tgk.Latn, azj.Latn, tel.Latn, slk.Latn, pap.Latn, zho.Hani, sme.Latn, tgl.Latn, uzn.Cyril, als.Latn, san.Deva, azb.Arab, ory.Orya, lmo.Latn, bre.Latn, mvf.Mong, fao.Latn, oci.Latn, sah.Cyril, sco.Latn, tuk.Latn, aze.Arabs, hin.Deva, haw.Latn, glk.Arabs, oss.Cyril, lug.Latn, tet.Latn, tsn.Latn, hrv.Latn, gsw.Latn, arz.Arabs, vec.Latn, mon.Latn, ilo.Latn, ctd.Latn, ben.Beng, roh.Latn, kal.Latn, asm.Beng, srp.Latn, bod.Tib, hil.Latn, rus.Latn, nds.Latn, lus.Latn, ido.Latn, lao.Lao, tir.Ethi, chv.Cyril, wln.Latn, kaa.Latn, pnb.Arabs, div.Thaa, som.Latn, jpn.Japan, hat.Latn, sna.Latn, heb.Hebr, bak.Cyril, nil.Latn, tel.Tel, kin.Latn, msa.Latn, gla.Latn, bos.Latn, dan.Latn, smo.Latn, ita.Latn, mar.Deva, pus.Arabs, spr.Cyril, spa.Latn, lat.Latn, hmnn.Latn, sin.Sinh, zul.Latn, bul.Cyril, amh.Ethi, ron.Latn, tam.Taml, khm.Khmr, nno.Latn, cos.Latn, fin.Latn, ori.Orya, uig.Orya, wig.Latn, hbs.Cyril, gle.Latn, cym.Latn, vie.Latn, kor.Hang, lit.Latn, yid.Hebr, ara.Arabs, sqi.Latn, pol.Latn, tur.Latn, swa.Latn, hau.Latn, ceb.Latn, eus.Latn, kir.Cyril, mlg.Latn, jav.Latn, snd.Arabs, sot.Latn, por.Latn, uzb.Cyril, fas.Arabs, nor.Latn, est.Latn, hun.Latn, ibo.Latn, itz.Latn, swe.Latn, tat.Cyril, ast.Latn, mya.Mymr, uzb.Latn, sun.Latn, ell.Grek, ces.Latn, mri.Latn, ckb.Arabs, kur.Latn, kaa.Cyril, nob.Latn, ukr.Cyril, fry.Latn, epo.Latn, nya.Latn
medium-high	79	div.Thaa, som.Latn, jpn.Japan, hat.Latn, sna.Latn, heb.Hebr, bak.Cyril, nil.Latn, tel.Tel, kin.Latn, msa.Latn, gla.Latn, bos.Latn, dan.Latn, smo.Latn, ita.Latn, mar.Deva, pus.Arabs, spr.Cyril, spa.Latn, lat.Latn, hmnn.Latn, sin.Sinh, zul.Latn, bul.Cyril, amh.Ethi, ron.Latn, tam.Taml, khm.Khmr, nno.Latn, cos.Latn, fin.Latn, ori.Orya, uig.Orya, wig.Latn, hbs.Cyril, gle.Latn, cym.Latn, vie.Latn, kor.Hang, lit.Latn, yid.Hebr, ara.Arabs, sqi.Latn, pol.Latn, tur.Latn, swa.Latn, hau.Latn, ceb.Latn, eus.Latn, kir.Cyril, mlg.Latn, jav.Latn, snd.Arabs, sot.Latn, por.Latn, uzb.Cyril, fas.Arabs, nor.Latn, est.Latn, hun.Latn, ibo.Latn, itz.Latn, swe.Latn, tat.Cyril, ast.Latn, mya.Mymr, uzb.Latn, sun.Latn, ell.Grek, ces.Latn, mri.Latn, ckb.Arabs, kur.Latn, kaa.Cyril, nob.Latn, ukr.Cyril, fry.Latn, epo.Latn, nya.Latn
medium-low	162	aym.Latn, rue.Cyril, rom.Latn, dzo.Tib, poh.Latn, sat.Olck, ary.Arabs, fur.Latn, mbt.Latn, bpy.Beng, iso.Latn, pon.Latn, glv.Latn, new.Deva, gym.Latn, bgp.Latn, kac.Latn, abt.Latn, que.Cyril, otq.Latn, sag.Latn, cak.Latn, avk.Latn, pan.Latn, meo.Latn, tum.Latn, bam.Latn, kha.Latn, syr.Syrc, kom.Cyril, nhe.Latn, abr.Arabs, nso.Arabs, nsl.Patn, pck.Latn, crs.Latn, acr.Latn, tat.Latn, acr.Arabs, uzs.Arabs, hil.Latn, mgh.Latn, tpi.Latn, bbe.Latn, meu.Latn, zza.Latn, ext.Latn, yue.Hani, ekk.Latn, xmf.Geor, nap.Latn, mzn.Arabs, pcm.Latn, lij.Latn, myv.Cyril, scn.Latn, dag.Latn, ban.Latn, twi.Latn, udm.Cyril, arb.Nso, nsl.Patn, pck.Latn, crs.Latn, acr.Latn, tat.Latn, acr.Arabs, uzs.Arabs, hil.Latn, mgh.Latn, tpi.Latn, ady.Cyril, pag.Latn, kiu.Latn, ber.Latn, iba.Latn, ksh.Latn, plt.Latn, lin.Latn, chk.Latn, tzo.Latn, th.Latn, ile.Latn, lub.Latn, hui.Latn, min.Latn, bjin.Latn, szl.Latn, kbp.Latn, inh.Cyril, que.Latn, ven.Latn, vls.Latn, kbd.Cyril, run.Latn, wol.Latn, ace.Latn, ada.Latn, kek.Latn, yua.Latn, tbz.Latn, gom.Latn, ful.Latn, mrj.Cyril, abk.Cyril, tuc.Latn, stq.Latn, mwl.Latn, tv.Latn, qwh.Latn, gom.Deva, mhr.Cyril, fij.Latn, grn.Latn, zap.Latn, mam.Latn, mps.Latn, tiv.Latn, ksd.Latn, ton.Latn, bik.Latn, vol.Latn, ava.Cyril, tsr.Latn, sz.y.Latn, ngu.Latn, hye.Armen, fon.Latn, skr.Arabs, kos.Latn, tyz.Latn, kur.Arabs, smr.Latn, tyv.Cyril, bci.Latn, vepr.Latn, erh.Cyril, kpg.Latn, hsb.Latn, ssw.Latn, zeia.Latn, ewe.Latn, ium.diq.Latn, lg.Latn, nzi.Latn, guj.Deva, ina.Latn, pms.Latn, bua.Cyril, lvs.Latn, eml.Latn, hmo.Latn, kum.Cyril, kab.Latn, chm.Cyril, cor.Latn, cfm.Latn, alt.Cyril, bcl.Latn, ang.Latn, fr.Latn, mai.Deva
unseen	393	rap.Latn, pmf.Latn, lsi.Latn, dje.Latn, blkx.Latn, ipk.Latn, syw.Deva, ann.Latn, bag.Latn, bat.Cyril, chu.Cyril, gwc.Arabs, adh.Latn, szy.Hani, shi.Arabs, njy.Latn, pdu.Latn, buo.Latn, cuv.Latn, udg.Mlym, bac.Latn, tio.Latn, jbj.Latn, taj.Deva, lez.Latn, olo.Cyril, rnl.Latn, bri.Latn, inh.Latn, kas.Cyril, wni.Latn, amp.Latn, tsc.Latn, mng.Latn, ldi.Cyril, mdf.Latn, arn.Latn, xty.Latn, llg.Latn, nge.Latn, gan.Latn, tuv.Latn, stk.Latn, nut.Latn, thy.Thai, lgr.Latn, hnj.Latn, dar.Cyril, aia.Latn, lwl.Thai, tnl.Latn, tvs.Latn, jra.Khmr, tay.Hani, gal.Latn, ybi.Deva, sink.Arabs, gag.Cyril, tuk.Cyril, trv.Hani, ydd.Hebr, kea.Latn, gbm.Deva, kwi.Latn, hro.Latn, rki.Latn, quy.Latn, tgd.Deva, zha.Hani, pcg.Mlym, tom.Latn, nsn.Latn, qul.Latn, jmx.Latn, kqr.Latn, mrr.Latn, bxa.Latn, abc.Latn, mve.Arabs, lfa.Latn, qup.Latn, yin.Latn, roo.Latn, mrw.Latn, nxa.Latn, yrk.Cyril, bem.Latn, kvt.Latn, csw.Cans, bjr.Latn, mgm.Latn, ngn.Latn, pib.Latn, quiz.Latn, awb.Latn, myk.Latn, otq.Arabs, ino.Latn, tkd.Latn, bef.Latn, bug.Bugs, aeu.Latn, nlv.Latn, dyt.Latn, bkc.Latn, mmf.Latn, hak.Hani, sea.Latn, sef.Latn, mlk.Latn, cbr.Latn, Imp.Latn, tmn.Latn, qzv.Latn, pbt.Arabs, cjs.Cyril, mlv.Cyril, bfm.Latn, dig.Latn, thk.Latn, zxx.Latn, lkb.Latn, chr.Latn, pnt.Latn, vif.Latn, fl.Latn, got.Latn, hbb.Latn, tll.Latn, bug.Latn, kxp.Arabs, qaa.Latn, krr.Khmr, kjj.Lao, isu.Latn, kmu.Latn, gof.Latn, sdk.Latn, mne.Latn, baw.Latn, idt.Latn, xkg.Latn, mgo.Latn, dtr.Latn, kins.Latn, ffn.Latn, hna.Latn, nxl.Latn, bld.Latn, odk.Arabs, miq.Latn, mhx.Latn, kam.Latn, yao.Latn, pnt.Grek, kby.Cyril, kpv.Cyril, kbx.Latn, cim.Latn, qvo.Latn, pih.Latn, nog.Latn, nco.Cyril, ryo.Cyril, clo.Latn, dmg.Latn, aaa.Latn, rel.Latn, ben.Latn, loh.Latn, thi.Deva, chd.Latn, cni.Latn, cjs.Latn, lbe.Latn, ybh.Deva, zxx.Zyyy, awa.Latn, gou.Latn, xmm.Latn, nqo.Latn, rut.Cyril, kbq.Latn, tkr.Latn, dvr.Ethi, ckt.Cyril, ady.Latn, ify.Latn, xal.Latn, bra.Deva, cgc.Latn, bhs.Latn, pwg.Latn, ang.Runn, oon.Runn, kri.Latn, qwe.Latn, qvm.Latn, bkm.Latn, bkh.Latn, niv.Latn, zuh.Latn, myr.Latn, fiu.Cyril, ssn.Latn, rki.Mymr, sox.Latn, yav.Latn, nyo.Latn, dag.Arabs, qkh.Latn, bze.Latn, myx.Latn, zav.Latn, ddg.Latn, wnk.Latn, bwx.Latn, mqy.Latn, lad.Hebr, boz.Latn, lue.Latn, ded.Latn, ph.Latn, avk.Cyril, wms.Latn, sgd.Latn, azn.Latn, ajz.Latn, psp.Latn, jra.Latn, smt.Latn, ags.Latn, csw.Latn, wtq.Latn, bfn.Latn, pli.Deva, snl.Latn, kwd.Latn, lgg.Latn, nza.Latn, wbr.Deva, lan.Latn, kmz.Latn, bzi.Thai, hao.Latn, nla.Latn, qxr.Latn, ken.Latn, tbj.Latn, blk.Latn, ybb.Latn, nwe.Latn, gan.Hani, snk.Latn, kak.Latn, tpi.Latn, hla.Latn, tks.Arabs, pea.Latn, bya.Latn, enc.Latn, jgo.Latn, tnp.Latn, aph.Deva, bfg.Latn, brv.Lao, nod.Thai, niq.Latn, nwi.Latn, xmd.Latn, gbj.Orya, syr.Latn, ify.Latn, xal.Latn, bra.Deva, cgc.Latn, bhs.Latn, pwg.Latn, ang.Runn, oon.Runn, kri.Latn, qwe.Latn, qvm.Latn, bkm.Latn, bkh.Latn, niv.Latn, zuh.Latn, myr.Latn, fiu.Cyril, ssn.Latn, rki.Mymr, sox.Latn, yav.Latn, nyo.Latn, dag.Arabs, qkh.Latn, bze.Latn, myx.Latn, zav.Latn, ekm.Latn, msb.Latn, unr.Orya, cac.Latn, chp.Cans, ckt.Latn, bss.Latn, lts.Latn, bbj.Latn, ttt.Cyril, kwu.Latn, smn.Cyril, kpy.Cyril, tod.Latn, wbm.Latn, tec.Latn, arc.Syrc, nst.Latn, tuz.Latn, bob.Latn, bfn.Latn, pli.Deva, snl.Latn, kwd.Latn, lgg.Latn, nza.Latn, wbr.Deva, lan.Latn, kmz.Latn, bzi.Thai, hao.Latn, nla.Latn, qxr.Latn, ken.Latn, tbj.Latn, blk.Latn, ybb.Latn, nwe.Latn, gan.Hani, snk.Latn, kak.Latn, tpi.Latn, hla.Latn, tks.Arabs, pea.Latn, bya.Latn, enc.Latn, jgo.Latn, tnp.Latn, aph.Deva, bfg.Latn, brv.Lao, nod.Thai, niq.Latn, nwi.Latn, xmd.Latn, gbj.Orya, syr.Latn, ify.Latn, xal.Latn, bra.Deva, cgc.Latn, bhs.Latn, pwg.Latn, ang.Runn, oon.Runn, kri.Latn, qwe.Latn, qvm.Latn, bkm.Latn, bkh.Latn, niv.Latn, zuh.Latn, myr.Latn, fiu.Cyril, ssn.Latn, rki.Mymr, sox.Latn, yav.Latn, nyo.Latn, dag.Arabs, qkh.Latn, bze.Latn, myx.Latn, zav.Latn, ekm.Latn, msb.Latn, unr.Orya, cac.Latn, chp.Cans, ckt.Latn, bss.Latn, lts.Latn, bbj.Latn, ttt.Cyril, kwu.Latn, smn.Cyril, kpy.Cyril, tod.Latn, wbm.Latn, tec.Latn, arc.Syrc, nst.Latn, tuz.Latn, bob.Latn, bfn.Latn, pli.Deva, snl.Latn, kwd.Latn, lgg.Latn, nza.Latn, wbr.Deva, lan.Latn, kmz.Latn, bzi.Thai, hao.Latn, nla.Latn, qxr.Latn, ken.Latn, tbj.Latn, blk.Latn, ybb.Latn, nwe.Latn, gan.Hani, snk.Latn, kak.Latn, tpi.Latn, hla.Latn, tks.Arabs, pea.Latn, bya.Latn, enc.Latn, jgo.Latn, tnp.Latn, aph.Deva, bfg.Latn, brv.Lao, nod.Thai, niq.Latn, nwi.Latn, xmd.Latn, gbj.Orya, syr.Latn, ify.Latn, xal.Latn, bra.Deva, cgc.Latn, bhs.Latn, pwg.Latn, ang.Runn, oon.Runn, kri.Latn, qwe.Latn, qvm.Latn, bkm.Latn, bkh.Latn, niv.Latn, zuh.Latn, myr.Latn, fiu.Cyril, ssn.Latn, rki.Mymr, sox.Latn, yav.Latn, nyo.Latn, dag.Arabs, qkh.Latn, bze.Latn, myx.Latn, zav.Latn, ddb.Latn, wnk.Latn, bwx.Latn, mqy.Latn, lad.Hebr, boz.Latn, lue.Latn, ded.Latn, ph.Latn, avk.Cyril, wms.Latn, sgd.Latn, azn.Latn, ajz.Latn, psp.Latn, jra.Latn, smt.Latn, ags.Latn, csw.Latn, wtq.Latn, bfn.Latn, pli.Deva, snl.Latn, kwd.Latn, lgg.Latn, nza.Latn, wbr.Deva, omw.Latn, khb.Tah, doi.Deva, gld.Cyril, ava.Latn, chu.Latn, dwr.Latn, azo.Latn, dug.Latn, bce.Latn, kmr.Latn, kpy.Armen, abq.Cyril, trp.Latn, ewo.Latn, the.Deva, hig.Latn, pkb.Latn, mxu.Latn, oji.Latn, tnt.Latn, mzm.Latn, mns.Cyril, lbe.Cyril, qvh.Latn, kmg.Latn, sps.Latn, brb.Khmr, tab.Latn, sxb.Latn, mkz.Latn, mgq.Latn, got.Goth, ins.Latn, arc.Latn, akb.Latn, skr.Latn, nsk.Cans, smi.Latn, pee.Mymr, eee.Thai, ihm.Deva, yux.Cyril, bqm.Latn, bcc.Arabs, nas.Latn, agg.Latn, xog.Latn, tsb.Latn, fub.Latn, mij.Latn, nsk.Latn, bsr.Cyril, dln.Latn, ozm.Latn, rmy.Latn, cre.Cans, kim.Cyril, cuh.Latn, ngl.Latn, yas.Latn, bud.Latn, miy.Latn, ame.Latn, pnz.Latn, raj.Deva, enb.Latn, cmo.Khmr, saq.Latn, tpu.Khmr, eve.Cyril, doc.Hani

mitigate errors introduced by the machine translation process, we filter out prompts with a BLEU score of less than 20. Figure 5 shows the average BLEU score of each language in the final dataset.

D PROMPTS

D.1 MACHINE TRANSLATION

We use the prompt below for machine translation.

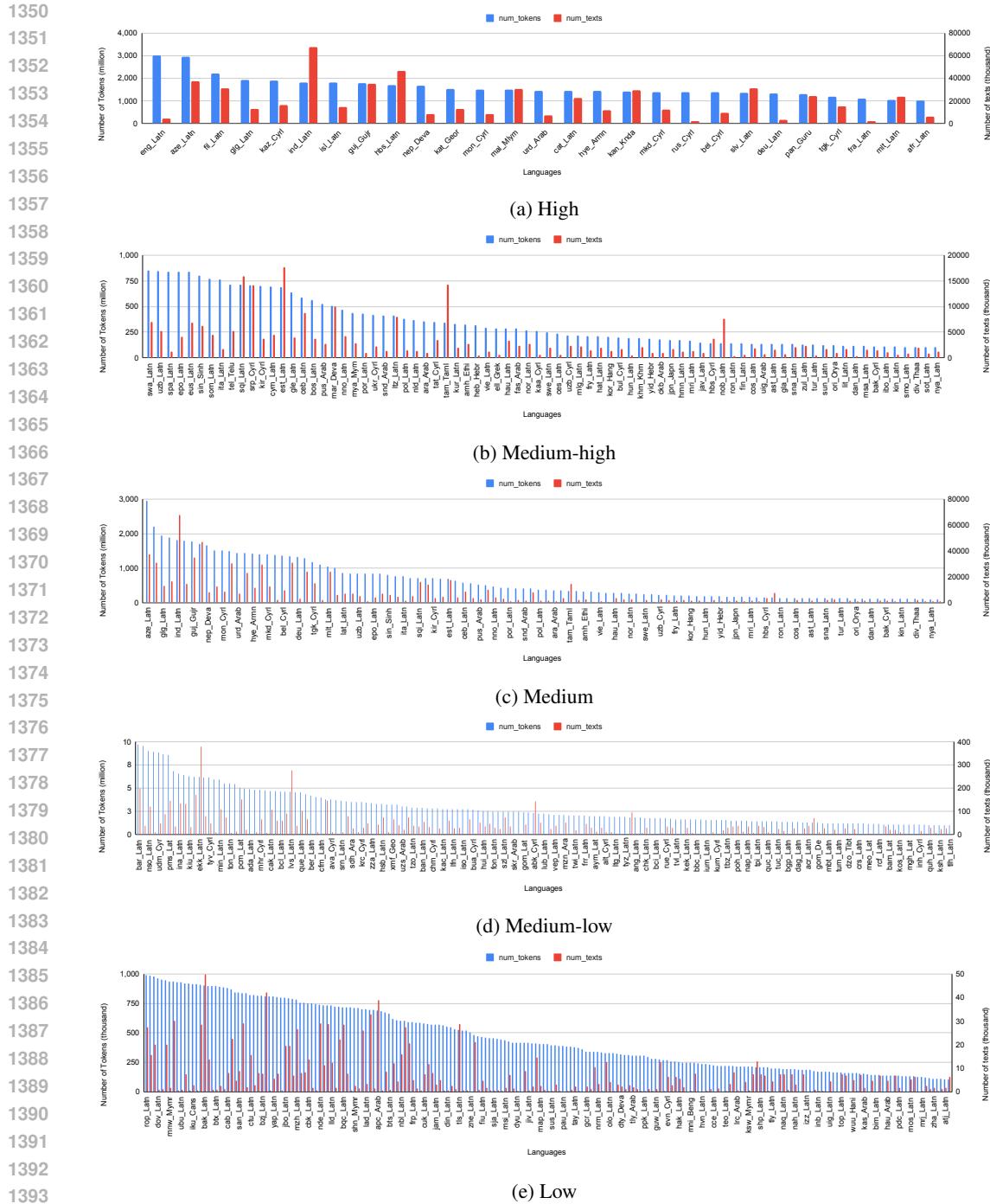


Figure 1: The number of texts and tokens of MaLA corpus in different resource groups.

Translate the following sentence from {src_lang} to {tgt_lang}
[{src_lang}]: {src_sent}
[{tgt_lang}]:

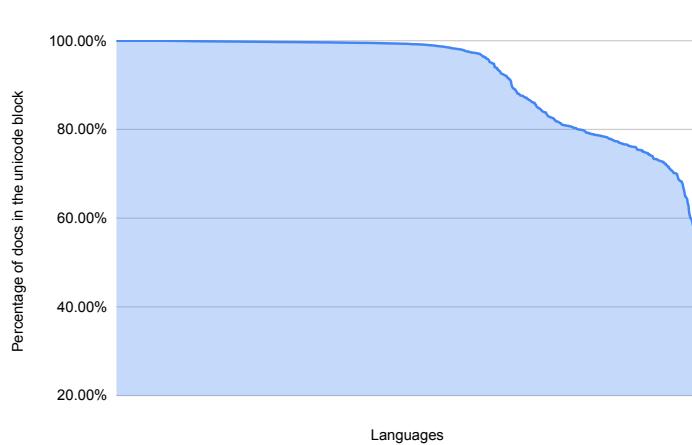


Figure 2: Unicode block distribution that measures the percentage of token counts falling into the Unicode block of each language

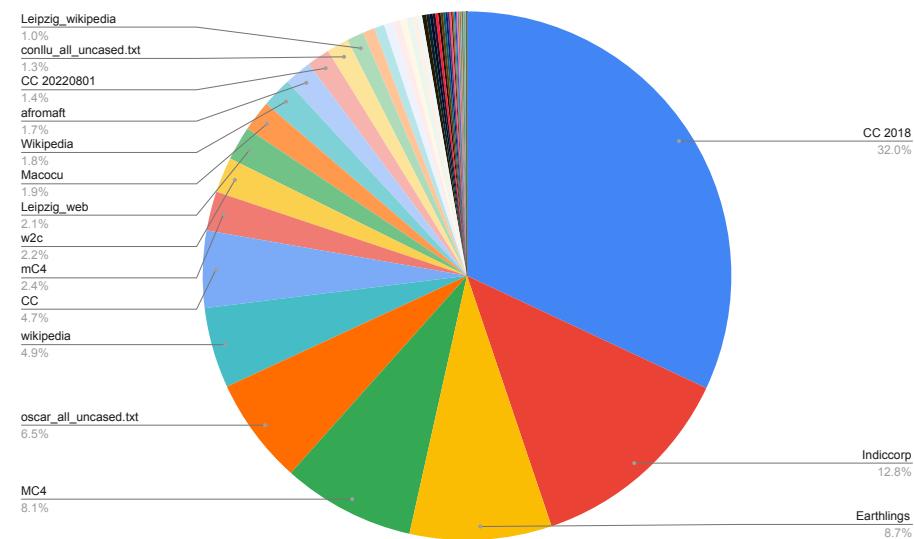


Figure 3: Data source distribution of MaLA corpus calculated by the number of documents

D.2 TEXT CLASSIFICATION

The prompt template for SIB-200 is as follows:

```
Topic Classification: science/technology, travel, politics, sports,
    health, entertainment, geography.
{examples}
The topic of the news "${text}" is
```

For Taxi-1500, the prompt template is as follows:

```
Topic Classification: Recommendation, Faith, Description, Sin, Grace
    , Violence.
{examples}
The topic of the verse "${text}" is
```

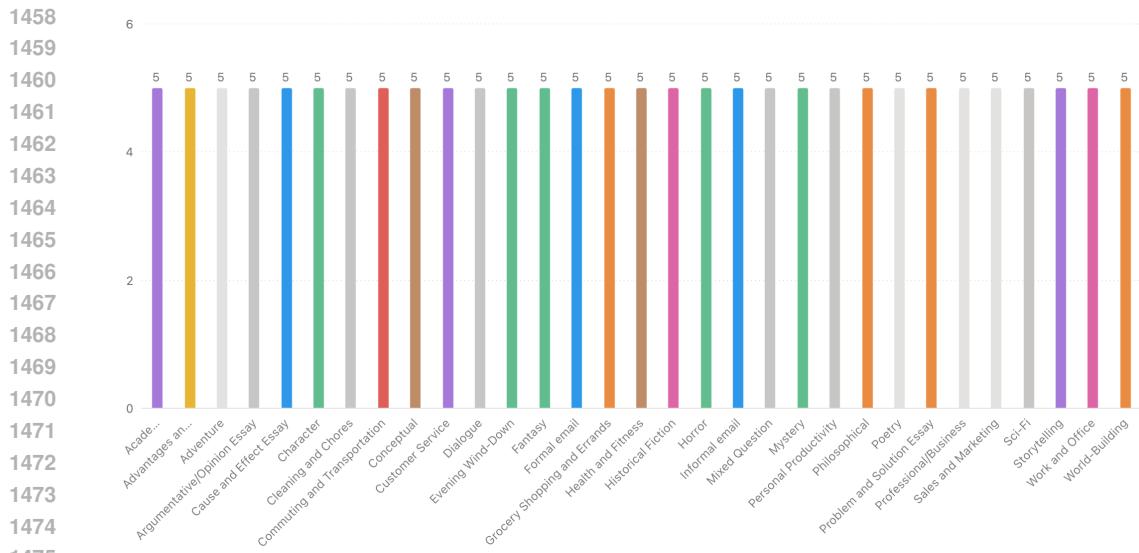


Figure 4: Writing tasks in the PolyWrite dataset.

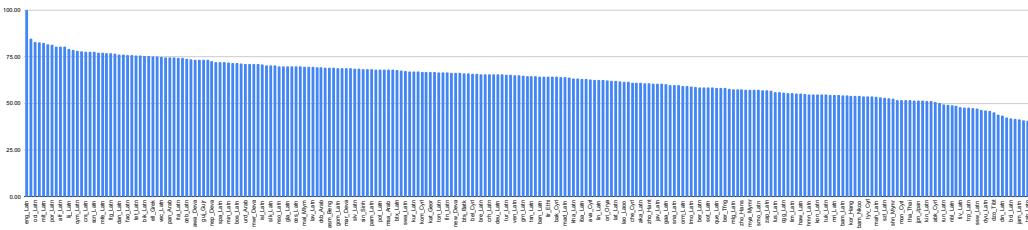


Figure 5: Mean BLEU scores per language in the PolyWrite dataset.

E DETAILED EVALUATION

This section presents detailed results of the evaluation. For benchmarks consisting of multiple languages, we host the results on GitHub⁴⁰. Those benchmarks include ARC multilingual, BELEBELE, Glot500-c test set, PBC, SIB-200, Taxi-1500, FLORES200, XLSum, Aya evaluation suite, and PolyWrite. For FLORES200, XLSum, Aya evaluation suite, and PolyWrite, we also release the generated texts of all compared models.

Evaluation Software We use the Language Model Evaluation Harness (lm-evaluation-harness) framework (Gao et al., 2023) for benchmarking test sets that are already ingested in the framework. For other benchmarks, we use in-house developed evaluation scripts and other open-source implementations. In text classification tasks such as SIB-200 and Taxi-1500, the evaluation protocol involves calculating the probability of the next output token for each candidate category. These probabilities are then sorted in descending order, with the category having the highest probability being selected as the model’s prediction. This process is implemented using the Transformers (Wolf et al., 2019) library. In tasks like Machine Translation and Open-Ended Generation, the vLLM library (Kwon et al., 2023) is used to accelerate inference, providing significant speed improvements. Similarly, for code generation tasks, we use a VLLM-enabled evaluation harness package (vllm-code-harness)⁴¹ for the execution-tested evaluations. For massively multilingual benchmarks with more than 100 languages, we categorize languages into five groups as listed in Table 8 in Appendix B.1.

⁴⁰https://github.com/MaLA-LM/emma-500/tree/main/evaluation_results

⁴¹<https://github.com/iNeil77/vllm-code-harness>

1512 Table 9: NLL on Glot500-c test. EMMA-500 Llama 2 7B has better average performance than all
 1513 baselines.

Model	Avg	High	Medium-High	Medium	Medium-Low	Low
Llama 2 7B	190.58	146.43	176.30	205.86	210.41	196.45
Llama 2 7B Chat	218.87	173.67	204.78	239.18	240.97	223.86
CodeLlama 2 7B	193.96	146.43	180.20	210.09	212.70	200.60
LLaMAX Llama 2 7B	187.37	108.58	142.84	197.74	212.22	203.83
LLaMAX Llama 2 7B Alpaca	169.12	94.84	123.90	173.54	193.83	187.34
MaLA-500 Llama 2 10B v1	155.62	127.51	153.93	173.02	166.01	158.12
MaLA-500 Llama 2 10B v2	151.25	123.82	147.69	167.46	161.20	155.13
TowerBase Llama 2 7B	192.98	150.41	180.19	209.12	212.18	198.89
TowerInstruct Llama 2 7B	199.44	157.33	186.42	216.92	218.82	204.93
Occiglot Mistral 7B v0.1	191.11	159.48	185.53	209.26	207.67	194.07
Occiglot Mistral 7B v0.1 Instruct	193.83	162.31	188.20	212.41	210.81	196.58
BLOOM 7B	202.95	160.33	195.01	216.15	220.46	206.89
BLOOMZ 7B	217.32	178.12	210.99	235.53	235.60	220.65
mGPT	340.37	311.14	337.29	388.97	367.45	343.14
XGLM 7.5B	205.07	199.08	210.22	225.53	214.92	205.67
Yayi 7B	226.67	181.48	217.34	243.42	246.58	231.65
Llama 3 8B	156.36	102.78	129.36	153.11	167.26	173.56
Llama 3.1 8B	154.59	101.23	127.57	150.87	164.81	172.05
Gemma 7B	692.25	583.39	721.77	817.40	729.61	689.60
Gemma 2 9B	320.81	348.26	351.08	380.49	338.00	303.62
Qwen 2 7B	188.55	132.47	171.50	200.26	210.21	196.78
Qwen 1.5 7B	195.52	141.37	181.41	212.10	217.16	202.46
EMMA-500 Llama 2 7B	112.20	81.78	100.89	122.53	99.28	109.25

E.1 INTRINSIC EVALUATION

For intrinsic evaluation, we compute the negative log-likelihood (NLL) of the test text given by the tested LLMs instead of using length-normalized perplexity due to different text tokenization schemes across models. We concatenate the test set as a single input and run a sliding-window approach.⁴² To make a comparison across models, we use the Glot500-c test set, which covers 534 evaluated languages. We also test on the Parallel Bible Corpus (PBC) which could yield NLL more comparable across languages. Table 9 and Table 10 show the intrinsic evaluation results on Glot500-c test and PBC, respectively. As shown, our model attains lower NLL compared to all other models on both test sets and languages with different resource availability. This suggests that the test sets are more similar to the underlying training data of our model and it can be interpreted that EMMA-500 has learned to compress massive multilingual text more efficiently.

E.2 SUMMARIZATION

XL-Sum (Hasan et al., 2021) is a large-scale multilingual abstractive summarization dataset that covers 44 languages. To evaluate the quality of the generated summaries, we use ROUGE-L (Lin, 2004) and BERTScore (Zhang et al., 2019) as evaluation metrics. ROUGE-L measures the longest common subsequence (LCS) between the reference and generated summaries. Recall and precision are calculated by dividing the LCS length by the reference length and the generated summary length, respectively. An F-score is then used to combine these two aspects into a single metric. BERTScore computes the semantic similarity between summaries by leveraging contextual embeddings from pre-trained language models. Specifically, we use the bert-base-multilingual-cased⁴³ model for BERTScore to accommodate multiple languages.

The evaluation results are presented in Table 11, we observe that our model, EMMA-500, performs comparably to other Llama2-based models like TowerInstruct Llama 2 7B. While EMMA-500 shows strength in certain language categories, it does not significantly outperform these models overall. This suggests that, although EMMA-500 is effective in multilingual summarization tasks, there is room for improvement to achieve more consistent performance across all languages. It is also important to note that this test set includes only 44 languages, limiting the evaluation of EMMA-500’s capabilities in low-resource languages.

⁴²<https://huggingface.co/docs/transformers/en/perplexity>

⁴³<https://huggingface.co/google-bert/bert-base-multilingual-cased>

Table 10: NLL on PBC. EMMA-500 Llama 2 7B has better average performance than all baselines.

Model	Avg	High	Medium-High	Medium	Medium-Low	Low
Llama 2 7B	122.10	91.30	99.41	112.31	133.08	135.34
Llama 2 7B Chat	139.14	108.40	115.78	129.79	149.54	152.82
CodeLlama 2 7B	123.98	93.27	101.83	113.47	134.65	137.52
LLaMAX Llama 2 7B	117.39	69.41	79.06	103.74	131.90	138.03
LLaMAX Llama 2 7B Alpaca	107.81	60.05	69.36	93.39	122.44	128.77
MaLA-500 Llama 2 10B v1	103.20	94.04	98.60	100.53	105.04	107.65
MaLA-500 Llama 2 10B v2	101.67	92.42	96.30	98.88	103.34	106.62
TowerBase Llama 2 7B	123.70	93.64	101.44	114.72	134.41	136.77
TowerInstruct Llama 2 7B	127.30	98.21	105.17	118.65	137.53	140.28
Occiglot Mistral 7B v0.1	121.64	95.15	101.86	114.37	131.49	132.93
Occiglot Mistral 7B v0.1 Instruct	123.18	96.86	103.41	115.88	132.98	134.48
BLOOM 7B	129.55	96.62	111.22	115.33	138.19	143.03
BLOOMZ 7B	137.72	107.03	119.89	125.85	145.27	150.95
mGPT	225.14	211.35	203.84	219.75	229.91	239.82
XGLM 7.5B	131.31	116.86	117.69	125.15	136.47	140.53
Yayi 7B	143.80	108.79	123.20	130.60	152.37	158.73
Llama 3 8B	102.55	64.20	71.82	85.22	114.79	121.29
Llama 3.1 8B	101.43	62.98	70.68	83.83	113.36	120.33
Gemma 7B	460.86	399.22	427.14	468.66	463.11	483.29
Gemma 2 9B	197.41	200.76	192.07	197.88	196.56	202.69
Qwen 2 7B	120.44	83.34	94.87	107.84	133.38	136.52
Qwen 1.5 7B	124.02	89.36	100.55	113.08	135.54	139.13
EMMA-500 Llama 2 7B	68.11	50.12	55.62	64.78	60.53	65.68

Table 11: Results on XL-Sum (ROUGE-L/BERTScore).

Model	Avg	High	Medium-High	Medium	Medium-Low	Low
Llama 2 7B	0.07/0.67	0.07/0.67	0.07/0.66	0.06/0.65	0.10/0.62	0.08/0.71
Llama 2 7B Chat	0.09/0.68	0.09/0.68	0.08/0.68	0.08/0.67	0.12/0.64	0.10/0.74
CodeLlama 2 7B	0.07/0.66	0.07/0.65	0.07/0.66	0.06/0.63	0.10/0.64	0.08/0.71
LLaMAX Llama 2 7B	0.05/0.65	0.05/0.65	0.05/0.65	0.05/0.61	0.07/0.62	0.06/0.69
LLaMAX Llama 2 7B Alpaca	0.10/0.69	0.10/0.70	0.09/0.69	0.09/0.68	0.13/0.65	0.12/0.74
MaLA-500 Llama 2 10B v1	0.05/0.64	0.06/0.64	0.05/0.64	0.05/0.61	0.07/0.64	0.06/0.68
MaLA-500 Llama 2 10B v2	0.05/0.64	0.06/0.64	0.05/0.64	0.05/0.61	0.07/0.64	0.06/0.69
Yayi Llama 2 7B	0.08/0.67	0.09/0.67	0.07/0.67	0.07/0.66	0.11/0.61	0.09/0.71
TowerBase Llama 2 7B	0.08/0.67	0.08/0.67	0.07/0.67	0.07/0.65	0.10/0.63	0.09/0.71
TowerInstruct Llama 2 7B	0.09/0.68	0.09/0.69	0.09/0.69	0.08/0.66	0.12/0.64	0.09/0.73
Occiglot Mistral 7B v0.1	0.07/0.66	0.08/0.66	0.07/0.66	0.07/0.64	0.10/0.63	0.09/0.71
Occiglot Mistral 7B v0.1 Instruct	0.08/0.67	0.09/0.67	0.08/0.67	0.08/0.64	0.10/0.64	0.09/0.73
BLOOM 7B	0.07/0.65	0.07/0.65	0.07/0.65	0.06/0.62	0.10/0.58	0.08/0.71
BLOOMZ 7B	0.11/0.70	0.14/0.69	0.09/0.70	0.10/0.69	0.12/0.66	0.15/0.74
mGPT	0.04/0.60	0.05/0.60	0.04/0.60	0.02/0.56	0.07/0.59	0.04/0.63
mGPT-13B	0.05/0.62	0.06/0.62	0.05/0.63	0.04/0.59	0.08/0.59	0.05/0.67
Yayi 7B	0.12/0.70	0.14/0.69	0.10/0.69	0.11/0.69	0.12/0.66	0.15/0.75
Llama 3 8B	0.08/0.67	0.08/0.67	0.08/0.67	0.08/0.65	0.10/0.65	0.11/0.72
Llama 3.1 8B	0.09/0.67	0.08/0.66	0.08/0.67	0.08/0.65	0.09/0.66	0.10/0.72
Gemma 2 9B	0.07/0.65	0.07/0.64	0.07/0.66	0.07/0.63	0.09/0.64	0.09/0.71
Gemma 7B	0.07/0.65	0.07/0.63	0.06/0.65	0.07/0.67	0.07/0.60	0.08/0.63
Qwen 1.5 7B	0.10/0.69	0.11/0.69	0.09/0.69	0.10/0.68	0.11/0.65	0.11/0.74
Qwen 2 7B	0.10/0.69	0.12/0.69	0.09/0.69	0.10/0.68	0.11/0.65	0.11/0.74
EMMA-500 Llama 2 7B	0.09/0.67	0.08/0.66	0.08/0.67	0.08/0.69	0.11/0.66	0.10/0.67

1620 We assess the models’ commonsense reasoning ability using three multilingual benchmarks. XCOPA
 1621 (Ponti et al., 2020) is a dataset covering 11 languages that focuses on causal commonsense
 1622 reasoning across multiple languages. XStoryCloze (Lin et al., 2022) is derived from the English
 1623 StoryCloze dataset (Mostafazadeh et al., 2017) and translated into 10 non-English languages, testing
 1624 commonsense reasoning within a story. In this task, the system must choose the correct ending
 1625 for a four-sentence narrative. XWinograd (Tikhonov & Ryabinin, 2021) is a multilingual collec-
 1626 tion of Winograd Schemas (Levesque et al., 2012) available in six languages, aimed at evaluating
 1627 cross-lingual commonsense reasoning.

1629 E.3 COMMONSENSE REASONING

1630 We perform zero-shot evaluations. Accuracy is used as the evaluation metric. We categorize languages
 1631 into different resource groups based on language availability, accessibility, and possible corpus size.
 1632 For XCOPA, we have three groups, i.e., high-resource (Italian, Turkish, Vietnamese, and Chinese),
 1633 medium-resource (Swahili due to its regional influence, Estonian, Haitian Creole, Indonesian, Thai,
 1634 and Tamil), and low-resource languages (Cusco-Collao Quechua).⁴⁴ For XStoryCloze, the resource
 1635 group is high-resource (Arabic, English, Spanish, Russian, and Chinese), medium (Hindi, Indonesian,
 1636 Swahili, and Telugu), and low (Basque and Burmese). For XWinograd, there is only one group for
 1637 high-resource languages. Table 12 shows the evaluation results of zero-shot commonsense reasoning.

1638 Compared with Llama 2-based models on XCOPA, our model improves the average performance by
 1639 a large margin—up to a 5% increase when compared with the best-performing TowerInstruct based
 1640 on Llama 2 7B.⁴⁵ Our model also outperforms all the multilingual LLMs. Recent LLMs such as
 1641 Gemma and Llama 3 have stronger performance than Llama 2 models, and Gemma 2 9B performs the
 1642 best. However, our model achieves better performance than Qwen, Llama 3, and 3.1. We gain similar
 1643 results on XStoryCloze, outperforming all the models except Gemma 2 9B. As for XWinograd, a
 1644 multilingual benchmark with high-resource only, our model achieves improved performance than
 1645 Llama 2, despite not being as good as Tower models, which target high-resource languages. However,
 1646 our model is comparable to other multilingual LLMs. For low-resource languages, our model
 1647 outperforms all the compared LLMs except LLaMAX Llama 2 7B Alpaca on XCOPA, where we
 1648 achieve the same accuracy as it.

1649 E.5 NATURAL LANGUAGE INFERENCE

1650 We evaluate on the XNLI (Cross-lingual Natural Language Inference) benchmark (Conneau et al.,
 1651 2018) where sentence pairs in different languages need to be classified as entailment, contradiction,
 1652 or neutral. We categorize the languages in XNLI into 3 groups, i.e., high-resource (German, English,
 1653 Spanish, French, Russian, and Chinese), medium-resource (Arabic, Bulgarian, Greek, Hindi, Turkish,
 1654 and Vietnamese), and low-resource (Swahili, Thai, and Urdu).⁴⁶ Table 17 shows the aggregated
 1655 accuracy. Our model outperforms most baselines including Llama 2-based models and multilingual
 1656 LLMs. We achieve the second-best average accuracy, slightly behind the Llama 3.1 8B model. On
 1657 the low-resource end, we perform the second-best, slightly behind the Gemma 2 9B model.

1659 E.6 MATH

1660 We evaluate the ability of LLMs to solve grade-school math problems across multiple languages on
 1661 the MGSM (Multilingual Grade School Math) benchmark (Shi et al., 2022). MGSM is an extension
 1662 of the GSM8K (Grade School Math 8K) dataset (Cobbe et al., 2021) by translating 250 of the original
 1663 GSM8K problems into ten languages. We also split these ten languages into three groups, i.e.,
 1664 high-resource (Spanish, French, German, Russian, Chinese, and Japanese), medium-resource (Thai,
 1665 Swahili, and Bengali), and low-resource (Telugu). Table 19 shows the results for 3-shot prompting
 1666 with a flexible match to obtain the answers in model generation. We evaluate all the models by
 1667 directly prompting the questions (denoted as direct) and the questions with answers followed by
 1668 Chain-of-Thoughts prompt in the same languages as the subset being evaluated (denoted as CoT)

1669 ⁴⁴Note that there is no perfect categorization for language resource groups.
 1670

1671 ⁴⁵The biggest improvements are on medium-resource languages. However, we note that the resource categor-
 1672 ization is not perfect.

1673 ⁴⁶Again, the categorization is not perfect.

Table 12: 0-shot results (ACC) on commonsense reasoning: XCOPA, XStoryCloze, and XWinograd. EMMA-500 Llama 2 7B has better average performance than Llama 2 models and multilingual LLMs on XCOPA, XStoryCloze, and comparable performance on XWinograd.

Model	XCOPA				XStoryCloze				XWinograd
	Avg	High	Medium	Low	Avg	High	Medium	Low	Avg
Llama 2 7B	0.5667	0.6210	0.5390	0.5160	0.5755	0.6338	0.5445	0.4921	0.7247
Llama 2 7B Chat	0.5585	0.6125	0.5313	0.5060	0.5841	0.6480	0.5477	0.4974	0.6945
CodeLlama 2 7B	0.5469	0.5870	0.5253	0.5160	0.5568	0.6068	0.5233	0.4990	0.7092
LLaMAX Llama 2 7B	0.5438	0.5550	0.5413	0.5140	0.6036	0.6434	0.5882	0.5347	0.6749
LLaMAX Llama 2 7B Alpaca	0.5660	0.5980	0.5517	0.5240	0.6383	0.6908	0.6201	0.5433	0.6986
MaLA-500 Llama 2 10B v1	0.5309	0.5355	0.5327	0.5020	0.5307	0.5815	0.4922	0.4808	0.6589
MaLA-500 Llama 2 10B v2	0.5309	0.5355	0.5327	0.5020	0.5307	0.5815	0.4922	0.4808	0.6589
YaYi Llama 2 7B	0.5671	0.6210	0.5413	0.5060	0.5842	0.6498	0.5491	0.4904	0.7450
TowerBase Llama 2 7B	0.5633	0.6250	0.5290	0.5220	0.5778	0.6435	0.5367	0.4957	0.7429
TowerInstruct Llama 2 7B	0.5705	0.6290	0.5400	0.5200	0.5924	0.6683	0.5453	0.4970	0.7400
Occiglot Mistral 7B v0.1	0.5667	0.6280	0.5337	0.5200	0.5810	0.6518	0.5328	0.5003	0.7461
Occiglot Mistral 7B v0.1 Instruct	0.5655	0.6285	0.5297	0.5280	0.5939	0.6694	0.5433	0.5063	0.7293
BLOOM 7B	0.5689	0.5995	0.5587	0.5080	0.5930	0.6199	0.5905	0.5308	0.7013
BLOOMZ 7B	0.5487	0.5635	0.5460	0.5060	0.5712	0.6114	0.5582	0.4967	0.6795
mGPT	0.5504	0.5710	0.5440	0.5060	0.5443	0.5496	0.5453	0.5291	0.5969
mGPT 13B	0.5618	0.5975	0.5513	0.4820	0.5644	0.5776	0.5635	0.5331	0.6359
XGLM 7.5B	0.6064	0.6400	0.6037	0.4880	0.6075	0.6242	0.6036	0.5738	0.6884
YaYi 7B	0.5664	0.5955	0.5550	0.5180	0.6067	0.6490	0.5940	0.5265	0.6979
Llama 3 8B	0.6171	0.6835	0.5903	0.5120	0.6341	0.6850	0.6203	0.5344	0.7684
Llama 3.1 8B	0.6171	0.6930	0.5880	0.4880	0.6358	0.6866	0.6209	0.5387	0.7552
Gemma 7B	0.6364	0.7035	0.6143	0.5000	0.6501	0.6946	0.6449	0.5493	0.7741
Gemma 2 9B	0.6633	0.7340	0.6427	0.5040	0.6767	0.7247	0.6669	0.5764	0.8007
Qwen 2 7B	0.6031	0.6865	0.5640	0.5040	0.6146	0.6945	0.5697	0.5050	0.7644
Qwen 1.5 7B	0.5944	0.6685	0.5590	0.5100	0.5985	0.6662	0.5604	0.5056	0.7259
EMMA-500 Llama 2 7B	0.6311	0.6660	0.6257	0.5240	0.6638	0.6892	0.6573	0.6132	0.7280

(Shi et al., 2022). The results show that Llama 2 7B is a weak model in the MGSM math task. The base model and its variants failed in most settings. Multilingual LLMs such as BLOOM, XGLM, and YaYi are also subpar at this task. Recent LLMs like Llama 3, Qwen, and Gemma obtain reasonable performance, and Qwen series models are the best. Our model improves the Llama 2 7B model remarkably in both prompt strategies. For direct prompting, our model has an average accuracy of 0.1702, 7% higher than the Llama 2 7B chat model. For CoT-based prompting, our model increases the score of the Llama 2 7B chat model from 0.1 to 0.3 (20% higher) and slightly outperforms Llama 3 and 3.1 8B models.

E.7 MACHINE READING COMPREHENSION

BELEBELE (Bandarkar et al., 2023) is a machine reading comprehension dataset covering 122 languages including high- and low-resource languages. Each question offers four multiple-choice answers based on a short passage from the FLORES-200 dataset. This benchmark is very challenging, even the English version of it presents remarkable challenges for advanced models. Table 22 shows the zero-shot results in different resource groups. Our continual pre-training improves the Llama 2 7B base model. But Llama 2 7B-based models mostly fail in this task and get quasi-random results. Recent advanced models like Llama 3.1 and Qwen 2 get reasonable results.

We then move to a more challenging task, the ARC multilingual test, which is a machine-translated benchmark (Lai et al., 2023) from the ARC dataset (Clark et al., 2018) that contains English science exam questions for multiple grade levels. We test the five-shot performance, with results shown in Table 23. The evaluation results show a similar pattern to BELEBELE that EMMA-500 improves Llama 2 but the 7B model is not capable of this challenging task, all Llama2 7B-based models obtain close to random results, and recent advances like Llama 3, Gemma, and Qwen get much better results.

E.8 CODE GENERATION

We conduct code generation evaluations on the Multipl-E (Cassano et al., 2022) benchmark in the interest of measuring the effects of massively multilingual continual pre-training on a model’s code generation utility and detecting if any catastrophic forgetting (Luo et al., 2023) has occurred on this front. Importantly, this also has implications for a model’s reasoning (Yang et al., 2024b) and entity-tracking abilities (Kim et al., 2024).

Table 13: 0-shot results (ACC) on XCOPA in all languages

Model	Avg	et-acc	stderr	ht-acc	stderr	id-acc	stderr	lt-acc	stderr	qp-acc	stderr	sw-acc	stderr	ta-acc	stderr	tb-acc	stderr	tr-acc	stderr	vi-acc	stderr	zh-acc	stderr
Llama 2 7B	0.5667	0.4860	0.0234	0.5060	0.0234	0.6240	0.0217	0.6580	0.0212	0.5160	0.0224	0.5220	0.0224	0.5340	0.0223	0.5630	0.0223	0.5480	0.0223	0.6280	0.0216	0.6550	0.0214
Llama 2 7B Chat	0.5585	0.4780	0.0234	0.5080	0.0234	0.6217	0.0217	0.6720	0.0206	0.5060	0.0224	0.5220	0.0224	0.5060	0.0224	0.5590	0.0223	0.5520	0.0223	0.6120	0.0218	0.6140	0.0218
CodeLlama 2 7B	0.5469	0.4680	0.0223	0.5180	0.0224	0.5740	0.0221	0.6300	0.0210	0.5160	0.0224	0.4880	0.0224	0.5500	0.0223	0.5540	0.0223	0.5380	0.0223	0.5580	0.0222	0.6220	0.0217
LLaMAX Llama 2 7B Alpaca	0.5438	0.4920	0.0224	0.5260	0.0224	0.5380	0.0223	0.5260	0.0224	0.5140	0.0224	0.5400	0.0223	0.5800	0.0221	0.5720	0.0221	0.5300	0.0223	0.6340	0.0216		
LLaMAX Llama 2 7B Alpaca v2	0.5399	0.4710	0.0234	0.5140	0.0234	0.5720	0.0223	0.5400	0.0219	0.5200	0.0224	0.5280	0.0223	0.5760	0.0222	0.5420	0.0223	0.5230	0.0224	0.5720	0.0219	0.6280	0.0224
MaLA-500 Llama 2 10B v1	0.5309	0.4860	0.0224	0.5140	0.0223	0.5340	0.0223	0.5300	0.0223	0.5940	0.0220	0.5060	0.0224	0.5220	0.0223	0.5760	0.0223	0.5420	0.0223	0.5540	0.0223	0.6230	0.0216
MaLA-500 Llama 2 10B v2	0.5309	0.4860	0.0224	0.5140	0.0223	0.5340	0.0223	0.5300	0.0223	0.5940	0.0220	0.5060	0.0224	0.5220	0.0223	0.5760	0.0223	0.5420	0.0223	0.5540	0.0223	0.6230	0.0216
YiYi Llama 2 10B	0.5671	0.4880	0.0224	0.5080	0.0224	0.6260	0.0217	0.6700	0.0210	0.5060	0.0224	0.5320	0.0223	0.5520	0.0223	0.5420	0.0223	0.5540	0.0223	0.6320	0.0216	0.6280	0.0216
TowerBase Llama 2 7B	0.5633	0.4600	0.0223	0.5020	0.0224	0.6020	0.0219	0.7080	0.0209	0.5200	0.0224	0.5660	0.0224	0.5380	0.0222	0.5920	0.0222	0.6620	0.0221	0.6740	0.0220		
TowerInstruct Llama 2 7B	0.5705	0.4880	0.0224	0.5160	0.0224	0.6200	0.0217	0.7100	0.0209	0.5200	0.0224	0.5660	0.0224	0.5460	0.0223	0.5860	0.0223	0.6740	0.0220	0.6740	0.0220		
Occiglot Mistral 7B v0.1	0.5676	0.4720	0.0224	0.5040	0.0224	0.5700	0.0222	0.5700	0.0219	0.5200	0.0224	0.5700	0.0224	0.5200	0.0223	0.5700	0.0223	0.5700	0.0223	0.6220	0.0219	0.6550	0.0219
Occiglot Mistral 7B v0.1 Instruct	0.5654	0.4820	0.0223	0.5100	0.0224	0.5840	0.0220	0.5780	0.0197	0.5200	0.0223	0.5800	0.0223	0.5660	0.0222	0.5580	0.0222	0.5650	0.0222	0.6250	0.0219	0.6560	0.0219
BLOOM 7B	0.5689	0.4820	0.0224	0.5020	0.0224	0.6280	0.0226	0.5280	0.0223	0.5080	0.0224	0.5180	0.0224	0.5540	0.0223	0.5420	0.0223	0.5100	0.0224	0.7080	0.0204	0.6520	0.0213
BLOOMZ 7B	0.5487	0.4920	0.0224	0.5400	0.0223	0.6060	0.0219	0.5140	0.0224	0.5060	0.0224	0.5340	0.0223	0.5740	0.0221	0.5300	0.0223	0.5220	0.0224	0.5980	0.0219	0.6200	0.0217
mGPT	0.5504	0.5300	0.0223	0.4980	0.0224	0.5880	0.0223	0.5220	0.0210	0.5060	0.0224	0.5620	0.0223	0.5600	0.0222	0.5320	0.0223	0.5520	0.0223	0.6020	0.0219	0.5400	0.0223
mGPT 1.0B	0.5416	0.5300	0.0223	0.4980	0.0224	0.5880	0.0223	0.5220	0.0210	0.5060	0.0224	0.5620	0.0223	0.5600	0.0222	0.5320	0.0223	0.5520	0.0223	0.6020	0.0219	0.5400	0.0223
mGPT 1.5B	0.5564	0.5140	0.0218	0.5740	0.0221	0.6940	0.0206	0.6380	0.0215	0.5600	0.0224	0.5460	0.0224	0.5940	0.0220	0.5840	0.0221	0.7020	0.0205	0.6380	0.0215		
YiYi 7B	0.5664	0.5040	0.0223	0.5300	0.0223	0.6340	0.0216	0.5180	0.0224	0.5540	0.0223	0.5620	0.0222	0.5460	0.0223	0.5200	0.0224	0.6640	0.0211	0.6800	0.0209		
EMMA-500 Llama 2 7B	0.6311	0.6140	0.0218	0.5800	0.0221	0.7420	0.0196	0.6940	0.0206	0.5240	0.0224	0.6620	0.0212	0.6000	0.0219	0.5560	0.0222	0.6200	0.0217	0.7020	0.0205	0.6480	0.0214

Table 14: 0-shot results (ACC) on XStoryCloze in all languages

Model	Avg	ar-acc	stderr	en-acc	stderr	es-acc	stderr	eu-acc	stderr	fr-acc	stderr	hi-acc	stderr	id-acc	stderr	my-acc	stderr	ru-acc	stderr	sw-acc	stderr	te-acc	stderr	zh-acc	stderr
Llama 2 7B	0.5755	0.4990	0.0129	0.7704	0.0108	0.6737	0.0121	0.5036	0.0129	0.5374	0.0128	0.5923	0.0126	0.4805	0.0129	0.6300	0.0124	0.5050	0.0129	0.5433	0.0128	0.5956	0.0126		
Llama 2 7B Chat	0.5841	0.5010	0.0129	0.7860	0.0105	0.6711	0.0121	0.5080	0.0129	0.5407	0.0128	0.5963	0.0126	0.4864	0.0129	0.6552	0.0122	0.5200	0.0129	0.5334	0.0128	0.6222	0.0125		
CodeLlama 2 7B	0.5561	0.5010	0.0129	0.7414	0.0116	0.6340	0.0124	0.5044	0.0129	0.4970	0.0129	0.5586	0.0126	0.4937	0.0129	0.5923	0.0126	0.5003	0.0129	0.5374	0.0128	0.5917	0.0126		
LLaMAX Llama 2 7B	0.6036	0.5390	0.0127	0.7511	0.0111	0.6625	0.0123	0.5447	0.0120	0.5817	0.0127	0.6062	0.0126	0.5248	0.0127	0.6122	0.0125	0.5718	0.0127	0.5930	0.0126	0.6082	0.0126		
LLaMAX Llama 2 7B Alpaca	0.6083	0.6036	0.0127	0.8147	0.0104	0.7068	0.0124	0.5486	0.0125	0.6214	0.0125	0.6465	0.0125	0.5748	0.0124	0.6744	0.0125	0.6166	0.0126	0.5950	0.0124	0.6545	0.0122		
MaLA-500 Llama 2 10B v1	0.5827	0.5010	0.0129	0.6900	0.0105	0.6710	0.0124	0.5600	0.0129	0.5220	0.0126	0.5765	0.0125	0.5676	0.0124	0.6215	0.0125	0.5740	0.0126	0.6360	0.0124	0.6769	0.0124		
MaLA-500 Llama 2 10B v2	0.5827	0.4818	0.0129	0.7233	0.0114	0.6341	0.0125	0.4990	0.0129	0.4729	0.0129	0.5820	0.0128	0.5427	0.0128	0.6462	0.0128	0.5403	0.0128	0.5236	0.0128	0.6169	0.0124		
YiYi 2 7B	0.5842	0.4997	0.0129	0.7099	0.0105	0.6670	0.0119	0.5093	0.0129	0.5427	0.0128	0.6142	0.0125	0.4745	0.0129	0.6479	0.0128	0.5003	0.0129	0.5394	0.0128	0.6234	0.0125		
TowerBase Llama 2 7B	0.5778	0.4917	0.0129	0.7223	0.0108	0.6982	0.0118	0.5076	0.0129	0.5288	0.0128	0.5813	0.0127	0.4838	0.0129	0.6704	0.0121	0.5324	0.0128	0.5850	0.0127				
TowerInstruct Llama 2 7B	0.5924	0.4931	0.0129	0.8087	0.0101	0.7161	0.0116	0.5060	0.0129	0.5295	0.0126	0.4871	0.0129	0.6936	0.0119	0.5149	0.0129	0.5414	0.0128	0.6300	0.0124				
Occiglot Mistral 7B v0.1	0.5810	0.5129	0.0129	0.7737	0.0108	0.7340	0.0114	0.5208	0.0129	0.5848	0.0128	0.5684	0.0127	0.4798	0.0129	0.6294	0.0124	0.5314	0.0128	0.6089	0.0126				
Occiglot Mistral 7B v0.1 Instruct	0.5899	0.5268	0.0128	0.7942	0.0104	0.7113	0.0113	0.5010	0.0128	0.5275	0.0128	0.6036	0.0126	0.4825	0.0129	0.6506	0.0123	0.5043	0.0129	0.5381	0.0128	0.6334	0.0124		
BLOOM 7B	0.5930	0.5857	0.0127	0.7055	0.0117	0.6618	0.0122	0.5725	0.0127	0.6042	0.0126	0.6453	0.0123	0.4891	0.0127	0.5725	0.0128	0.5394	0.0128	0.6188	0.0125				
BLOOMZ 7B	0.5712	0.6552	0.0128	0.7303	0.0114	0.6459	0.0125	0.5109	0.0129	0.5764	0.0127	0.5533	0.0128	0.4825	0.0129	0.5215	0.0129	0.5817	0.0127	0.5943	0.0126				
mGPT	0.5443	0.4911	0.0129	0.5994	0.0126	0.5546	0.0128	0.5275	0.0128	0.5318	0.0126	0.5222	0.0125	0.5665	0.0128	0.5260	0.0128	0.5725	0.0127	0.5341	0.0128				
mGPT 1.0B	0.5644	0.5162	0.0129	0.6424	0.0128	0.5864	0.0127	0.5539	0.0128	0.5056	0.0128	0.5811	0.0127	0.5122	0.0129	0.5943	0.0128	0.5461	0.0128	0.5764	0.0127	0.5493	0.0128		
mGPT 1.5B	0.6075	0.5612	0.0128	0.6982	0.0118	0.6386	0.0124	0.5771	0.0127	0.5884	0.0127	0.5300	0.0124	0.5705	0.0127	0.6340	0.0128	0.5904	0.0126	0.6023	0.0126	0.5894	0.0127		
YiYi 7B	0.6067	0.6181	0.0125	0.7432	0.0112	0.6942	0.0119	0.5609	0.0126	0.5637	0.0124	0.5421	0.0125	0.4924	0.0129	0.5215	0.0128	0.5361	0.0128	0.6767	0.0121				
LLama 3 7B	0.6341	0.5864	0.0127	0.7869	0.0105	0.7062	0.0117	0.5579	0.0128	0.6281	0.0124	0.5902	0.01												

1782 **E.4 MACHINE TRANSLATION**

1784 **Table 16:** 3-shot results on FLORES-200 (Eng-X, BLEU/chrF++). EMMA-500 Llama 2 7B has
1785 better average performance than all baselines.

Model	Avg	High	Medium-High	Medium	Medium-Low	Low
Llama 2 7B	4.62/ 15.13	10.77/ 24.38	8.56/ 21.4	2.55/ 13.72	0.74/ 8.72	0.7/ 7.92
Llama 2 7B Chat	4.95/ 16.95	10.87/ 24.51	8.54/ 22.69	3.25/ 15.5	1.52/ 12.08	0.94/ 10.03
CodeLlama 2 7B	4.27/ 14.94	10.04/ 23.48	7.79/ 20.79	2.57/ 14.2	0.71/ 9.27	0.58/ 7.49
LLaMAX Llama 2 7B	0.8/ 7.42	1.85/ 12.06	1.2/ 9.74	0.54/ 6.55	0.22/ 4.52	0.38/ 4.81
LLaMAX Llama 2 7B Alpaca	12.51/ 28.35	24.8/ 41.76	18.69/ 38.42	10.1/ 27.27	3.79/ 16.53	6.68/ 18.15
MaLA-500 Llama 2 10B v1 [‡]	0.6/ 6.08	1.51/ 9.0	1.13/ 8.19	0.35/ 5.99	0.07/ 4.5	0.02/ 2.9
MaLA-500 Llama 2 10B v2 [‡]	0.54/ 6.38	1.4/ 9.19	1.02/ 8.42	0.24/ 5.99	0.07/ 5.14	0.02/ 3.27
Yayi Llama 2 7B	4.41/ 14.87	10.49/ 24.0	8.21/ 21.27	2.52/ 13.57	0.6/ 8.49	0.53/ 7.42
TowerBase Llama 2 7B	4.83/ 16.03	11.89/ 24.15	8.33/ 21.46	2.57/ 14.49	1.38/ 11.6	0.74/ 8.9
TowerInstruct Llama 2 7B	3.23/ 15.64	7.22/ 22.65	4.99/ 20.0	2.2/ 14.9	1.62/ 12.31	0.73/ 8.97
Occiglot Mistral 7B v0.1	4.32/ 16.1	10.5/ 23.74	6.95/ 20.91	2.87/ 15.44	1.47/ 12.0	0.79/ 9.09
Occiglot Mistral 7B v0.1 Instruct	3.99/ 15.8	9.46/ 23.17	6.46/ 20.73	2.68/ 15.29	1.31/ 11.31	0.84/ 9.04
BLOOM 7B	2.81/ 11.8	7.53/ 19.0	3.12/ 13.36	2.05/ 11.48	0.85/ 8.0	2.09/ 9.22
BLOOMZ 7B [†]	7.44/ 16.1	23.64/ 32.22	7.46/ 16.62	6.98/ 16.05	1.28/ 9.99	4.17/ 11.77
mGPT	2.59/ 12.56	5.24/ 17.04	4.75/ 16.92	1.14/ 9.75	0.78/ 9.24	0.84/ 9.08
mGPT-13B	3.88/ 14.57	8.32/ 21.7	6.84/ 20.55	2.06/ 12.23	0.9/ 8.4	1.33/ 9.58
Yayi 7B	4.37/ 13.5	13.72/ 26.28	4.68/ 14.31	3.35/ 12.89	0.91/ 8.51	2.55/ 10.08
Llama 3 8B	9.93/ 24.08	20.38/ 36.87	14.95/ 32.05	8.89/ 24.28	2.83/ 14.26	4.2/ 14.29
Llama 3.1 8B	10.11/ 24.69	20.82/ 37.39	15.3/ 32.82	8.85/ 24.85	2.9/ 14.83	4.23/ 14.81
Gemma 2 9B	12.09/ 26.48	24.62/ 40.69	17.82/ 35.51	10.68/ 26.58	3.38/ 15.02	5.94/ 15.98
Gemma 7B	9.05/ 23.05	17.58/ 34.5	13.62/ 30.16	7.96/ 22.85	2.64/ 14.11	4.47/ 14.82
Qwen 1.5 7B	5.87/ 17.77	14.05/ 28.6	8.88/ 23.57	3.85/ 17.07	1.7/ 10.85	2.35/ 10.21
Qwen 2 7B	5.56/ 17.17	13.22/ 27.65	8.21/ 22.36	4.15/ 16.93	1.56/ 10.47	2.19/ 10.11
EMMA-500 Llama 2 7B	15.58/ 33.25	26.37/ 42.4	21.96/ 41.98	13.4/ 32.06	9.15/ 27.99	7.92/ 21.25

1806 **Table 17:** 0-shot results on XNLI (ACC).

Model	Avg	High	Medium	Low
Llama 2 7B	0.4019	0.4526	0.3772	0.3497
Llama 2 7B Chat	0.3858	0.4277	0.3675	0.3387
CodeLlama 2 7B	0.4019	0.4627	0.3729	0.3386
LLaMAX Llama 2 7B	0.4427	0.4653	0.4264	0.4303
LLaMAX Llama 2 7B Alpaca	0.4509	0.4847	0.4280	0.4289
MaLA-500 Llama 2 10B v1	0.3811	0.4210	0.3585	0.3465
MaLA-500 Llama 2 10B v2	0.3811	0.4210	0.3585	0.3465
YaYi Llama 2 7B	0.4128	0.4732	0.3841	0.3494
TowerBase Llama 2 7B	0.3984	0.4608	0.3633	0.3439
TowerInstruct Llama 2 7B	0.4036	0.4707	0.3692	0.3379
Occiglot Mistral 7B v0.1	0.4235	0.4990	0.3839	0.3519
Occiglot Mistral 7B v0.1 Instruct	0.4081	0.4758	0.3718	0.3452
BLOOM 7B	0.4160	0.4513	0.3969	0.3838
BLOOMZ 7B	0.3713	0.4002	0.3556	0.3451
mGPT	0.4051	0.4297	0.3965	0.3730
XGLM 7.5B	0.4375	0.4572	0.4216	0.4300
YaYi 7B	0.3987	0.4385	0.3824	0.3515
Llama 3 8B	0.4497	0.4882	0.4384	0.3956
Llama 3.1 8B	0.4562	0.4961	0.4404	0.4083
Gemma 7B	0.4258	0.4644	0.4100	0.3801
Gemma 2 9B	0.4674	0.4850	0.4511	0.4649
Qwen 2 7B	0.4277	0.4731	0.4135	0.3653
Qwen 1.5 7B	0.3947	0.4095	0.3880	0.3783
EMMA-500 Llama 2 7B	0.4514	0.4609	0.4440	0.4471

1829 **F RELATED WORK**

1831 **Multilingual LLMs** Multilingual large language models (LLMs) have made significant progress
1832 in processing and understanding multiple languages within a unified framework. Models like mT5
1833 (Xue et al., 2021) and XGLM (Lin et al., 2022) leverage both monolingual and multilingual datasets
1834 to perform tasks such as translation and text summarization across a wide spectrum of languages.
1835 However, the predominant focus on English has led to disparities in performance, particularly for
low-resource languages. Recent work on multilingual LLMs, such as BLOOM (Scao et al., 2022), has

Table 18: 0-shot results (ACC) on XNLI in all languages.

Model	Avg	ar-acc	sider	bg-acc	sider	de-acc	sider	el-acc	sider	en-acc	sider	es-acc	sider	fr-acc	sider	hi-acc	sider	ro-acc	sider	sv-acc	sider	th-acc	sider	tr-acc	sider	zh-acc	sider	vocab	sider	zh-acc	sider
Llama 2 7B	0.4019	0.3542	0.0096	0.2465	0.0099	0.4711	0.0100	0.3667	0.0097	0.5530	0.0100	0.4052	0.0098	0.5024	0.0100	0.4822	0.0098	0.4237	0.0099	0.3494	0.0096	0.3635	0.0096	0.3727	0.0097	0.3361	0.0095	0.3663	0.0097	0.3616	0.0096
Llama 2 7B Chat	0.3858	0.3442	0.0095	0.3707	0.0097	0.4309	0.0100	0.3815	0.0097	0.5024	0.0100	0.4842	0.0098	0.3578	0.0096	0.3422	0.0095	0.3349	0.0095	0.3695	0.0097	0.3390	0.0095	0.3811	0.0097	0.3695	0.0097	0.3695	0.0097		
CodeLlama 2 7B	0.4019	0.3541	0.0095	0.3775	0.0097	0.4723	0.0100	0.3763	0.0097	0.5475	0.0100	0.4438	0.0100	0.4920	0.0100	0.3856	0.0096	0.4669	0.0100	0.3329	0.0094	0.3502	0.0096	0.3659	0.0098	0.3325	0.0094	0.4440	0.0098	0.3594	0.0096
LLaMAX Llama 2 7B	0.4421	0.3642	0.0095	0.4251	0.0097	0.5101	0.0100	0.4550	0.0097	0.5613	0.0100	0.4550	0.0098	0.4526	0.0100	0.4526	0.0097	0.4526	0.0100	0.4526	0.0098	0.4526	0.0100	0.4526	0.0098	0.4526	0.0100	0.4526	0.0098		
LLaMAX Llama 2 7B Alpaca	0.4509	0.3442	0.0095	0.4639	0.0100	0.4876	0.0100	0.3431	0.0099	0.5811	0.0099	0.4896	0.0100	0.3197	0.0100	0.4562	0.0100	0.4627	0.0100	0.4357	0.0099	0.4880	0.0099	0.4386	0.0099	0.4430	0.0100	0.4389	0.0099	0.3578	0.0096
MaLa-500 Llama 2 10B v1	0.3811	0.3594	0.0096	0.4120	0.0099	0.4751	0.0100	0.3446	0.0099	0.5618	0.0099	0.3410	0.0098	0.4759	0.0100	0.3365	0.0095	0.3394	0.0095	0.3222	0.0096	0.3369	0.0095	0.3383	0.0095	0.3502	0.0096	0.3602	0.0096	0.3325	0.0094
MaLa-500 Llama 2 10B v2	0.3812	0.3594	0.0095	0.4161	0.0099	0.4884	0.0100	0.3735	0.0097	0.5618	0.0099	0.3414	0.0098	0.4668	0.0100	0.3370	0.0096	0.3467	0.0100	0.3570	0.0096	0.3396	0.0095	0.3468	0.0100	0.3502	0.0096	0.3354	0.0096		
Yayi Llama 2 7B	0.4126	0.3441	0.0095	0.4161	0.0099	0.4884	0.0100	0.3735	0.0097	0.5618	0.0099	0.3414	0.0098	0.4668	0.0100	0.3370	0.0096	0.3467	0.0100	0.3570	0.0096	0.3396	0.0095	0.3468	0.0100	0.3502	0.0096	0.3354	0.0096		
TowerBase Llama 2 7B	0.3984	0.3390	0.0095	0.4137	0.0099	0.4787	0.0100	0.3526	0.0098	0.5653	0.0099	0.4169	0.0099	0.4944	0.0100	0.3454	0.0095	0.4594	0.0100	0.3502	0.0096	0.3478	0.0095	0.3571	0.0096	0.3337	0.0095	0.3719	0.0097	0.3522	0.0096
TowerInstruct Llama 2 7B	0.4036	0.3365	0.0095	0.4293	0.0099	0.4847	0.0100	0.3498	0.0098	0.5695	0.0099	0.4651	0.0100	0.3474	0.0095	0.4627	0.0100	0.3394	0.0096	0.3390	0.0095	0.3787	0.0097	0.3335	0.0095	0.3735	0.0097	0.3747	0.0097		
Occiglot Mistral 7B v0.1	0.4019	0.3600	0.0095	0.4543	0.0097	0.5177	0.0100	0.3815	0.0098	0.5920	0.0100	0.4193	0.0098	0.4765	0.0100	0.3476	0.0095	0.4765	0.0100	0.3449	0.0095	0.4765	0.0100	0.3467	0.0095	0.4765	0.0100	0.3467	0.0095		
Occiglot Mistral 7B v0.1 Instruct	0.4019	0.3600	0.0095	0.4543	0.0097	0.5177	0.0100	0.3815	0.0098	0.5920	0.0100	0.4193	0.0098	0.4765	0.0100	0.3476	0.0095	0.4765	0.0100	0.3449	0.0095	0.4765	0.0100	0.3467	0.0095	0.4765	0.0100	0.3467	0.0095		
BLOOM 7B	0.4160	0.3385	0.0097	0.3884	0.0098	0.4884	0.0100	0.3663	0.0098	0.5568	0.0100	0.4009	0.0098	0.4248	0.0100	0.3428	0.0096	0.3790	0.0095	0.3409	0.0095	0.3409	0.0095	0.3409	0.0095	0.3409	0.0095	0.3409	0.0095		
BLOOMZ 7B	0.3711	0.3269	0.0094	0.3402	0.0095	0.3978	0.0099	0.3537	0.0096	0.5397	0.0097	0.4850	0.0097	0.4980	0.0097	0.4651	0.0097	0.4303	0.0097	0.3788	0.0099	0.3505	0.0097	0.3509	0.0097	0.4220	0.0097	0.4279	0.0097	0.3535	0.0098
mGPT	0.4071	0.3384	0.0095	0.4169	0.0094	0.4699	0.0100	0.3582	0.0098	0.4687	0.0100	0.3409	0.0098	0.4040	0.0098	0.3747	0.0097	0.3371	0.0097	0.3636	0.0097	0.3667	0.0097	0.3502	0.0097	0.3667	0.0097	0.3502	0.0097		
XGLM 7.5B	0.4373	0.3349	0.0095	0.4365	0.0099	0.4755	0.0100	0.4040	0.0098	0.5309	0.0100	0.4707	0.0100	0.4542	0.0100	0.4598	0.0100	0.4570	0.0100	0.4137	0.0099	0.4679	0.0100	0.4193	0.0099	0.4289	0.0099	0.3522	0.0096		
Yayi 7B	0.3987	0.3980	0.0099	0.3578	0.0096	0.4099	0.0100	0.4261	0.0099	0.4819	0.0100	0.4004	0.0098	0.3912	0.0099	0.3406	0.0095	0.3434	0.0095	0.3317	0.0094	0.3707	0.0097	0.4418	0.0100	0.3494	0.0096	0.3494	0.0096		
EMMA-500 Llama 2 7B	0.4514	0.3478	0.0095	0.4627	0.0100	0.4707	0.0100	0.4586	0.0100	0.5378	0.0100	0.4707	0.0100	0.4887	0.0100	0.4759	0.0100	0.4655	0.0100	0.4618	0.0100	0.4480	0.0100	0.4598	0.0100	0.3522	0.0096				

Table 19: 3-shot results (ACC) on MGSM obtained by direct and CoT prompting.

Model	Direct Prompting				CoT Prompting			
	Avg	High	Medium	Low	Avg	High	Medium	Low
Llama 2 7B	0.0669	0.0807	0.0213	0.0120	0.0636	0.0760	0.0213	0.0080
Llama 2 7B Chat	0.1022	0.1373	0.0213	0.0080	0.1091	0.1353	0.0280	0.0160
CodeLlama 2 7B	0.0593	0.0707	0.0293	0.0120	0.0664	0.0873	0.0267	0.0200
LLaMAX Llama 2 7B	0.0335	0.0400	0.0200	0.0080	0.0362	0.0433	0.0227	0.0240
LLaMAX Llama 2 7B Alpaca	0.0505	0.0520	0.0400	0.0160	0.0635	0.0807	0.0413	0.0080
MaLA-500 Llama 2 10B v1	0.0091	0.0133	0.0027	0.0000	0.0073	0.0127	0.0000	0.0000
MaLA-500 Llama 2 10B v2	0.0091	0.0133	0.0027	0.0000	0.0073	0.0127	0.0000	0.0000
TowerBase Llama 2 7B	0.0615	0.0833	0.0173	0.0080	0.0616	0.0860	0.0240	0.0080
TowerInstruct Llama 2 7B	0.0724	0.0953	0.0173	0.0200	0.0824	0.1047	0.0187	0.0120
Occiglot Mistral 7B v0.1	0.1331	0.1687	0.0453	0.0160	0.1407	0.1880	0.0360	0.0120
Occiglot Mistral 7B v0.1 Instruct	0.2276	0.2980	0.0747	0.0280	0.2216	0.3040	0.0787	0.0280
BLOOM 7B	0.0287	0.0260	0.0280	0.0360	0.0229	0.0220	0.0147	0.0200
BLOOMZ 7B	0.0255	0.0267	0.0240	0.0120	0.0215	0.0167	0.0307	0.0200
mGPT	0.0135	0.0167	0.0053	0.0000	0.0142	0.0193	0.0067	0.0000
mGPT 13B	0.0131	0.0180	0.0067	0.0000	0.0153	0.0167	0.0120	0.0000
XGLM 7.5B	0.0102	0.0120	0.0067	0.0200	0.0116	0.0107	0.0120	0.0280
Yayi 7B	0.0276	0.0293	0.0147	0.0240	0.0302	0.0293	0.0240	0.0200
Llama 3 8B	0.2745	0.2787	0.2613	0.0560	0.2813	0.2853	0.2667	0.0520
Llama 3.1 8B	0.2836	0.2900	0.2613	0.0440	0.2731	0.2727	0.2547	0.0840
Gemma 7B	0.3822	0.3660	0.3827	0.2720	0.3578	0.3467	0.3707	0.2680
Gemma 2 9B	0.3295	0.2800	0.3573	0.3080	0.4469	0.3607	0.5200	0.4720
Qwen 2 7B	0.4895	0.5440	0.3880	0.1480	0.4469	0.3607	0.5200	0.4720
Qwen 1.5 7B	0.3156	0.4000	0.1600	0.0400	0.5147	0.5893	0.3907	0.1440
EMMA-500 Llama 2 7B	0.1702	0.1920	0.1187	0.0240	0.3036	0.4060	0.1480	0.0240

shown that adapting these English-centric models through vocabulary extension based on multilingual corpora and continual pre-training (CPT) can improve performance across languages, especially low-resource ones. These models highlight the importance of efficient tokenization and adaptation, which can bridge the performance gap between high-resource and low-resource languages.

Multilingual Corpora The availability and use of multilingual corpora play a crucial role in training multilingual LLMs. CC100 Corpus (Conneau et al., 2020), launched in 2020, encompasses hundreds of billions of tokens and over 100 languages. Further, CC100-XL Corpus (Lin et al., 2022), created for the training of XGLM, extends across 68 Common Crawl Snapshots and 134 languages, aiming to balance language presentation and improve performance in few-shot and zero-shot tasks. The ROOTS Corpus (Laurençon et al., 2022), released in July 2022, supports BLOOM with approximately 341 billion tokens across 46 natural languages. It emphasizes underrepresented languages such as Swahili and Catalan, drawing from diverse sources including web crawls, books and academic publications. Besides, Occiglot Fineweb⁴⁷, which began to be released in early 2024, consists of around 230 million documents in 10 European languages. It combines curated and cleaned web data to support efficient training for both high- and low-resource European languages. Additionally, recent efforts parallel corpus construction from web crawls, such as ParaCrawl (Bañón et al., 2020) and CCMATRIX (Schwenk et al., 2021b), have contributed to large-scale multilingual training too.

⁴⁷<https://occiglot.eu/posts/occiglot-fineweb/>

Tables 20 and 21 show 3-shot results on MGSM by direct and Chain-of-Thought prompting respectively. All the scores are obtained by flexible matching.

Table 20: 3-shot results (ACC) on MGSM by direct prompting and flexible matching.

Model	Avg	bn	bn-stder	de	de-stder	en	en-stder	es	es-stder	fr	fr-stder	ja	ja-stder	ru	ru-stder	sw	sw-stder	te	te-stder	th	th-stder	zh	zh-stder		
Llama 2 7B	0.069	0.290	0.0105	0.0900	0.0172	0.1600	0.0241	0.1120	0.0300	0.1200	0.0200	0.0200	0.0540	0.0097	0.0800	0.0172	0.0280	0.0105	0.0120	0.069	0.0080	0.0056	0.080	0.0160	
Llama 2 7B Chat	0.102	0.280	0.0105	0.1680	0.0237	0.2280	0.0266	0.1960	0.0252	0.0200	0.0250	0.0240	0.0097	0.0140	0.0222	0.0080	0.0056	0.0280	0.0105	0.0100	0.0190				
CodeLlama 2 7B	0.093	0.160	0.0088	0.0880	0.0180	0.1280	0.0212	0.0880	0.0180	0.0400	0.0193	0.0500	0.0146	0.0600	0.0151	0.0200	0.0089	0.0120	0.0069	0.0520	0.0141	0.0280	0.0105		
LLaMAX Llama 2 7B	0.035	0.280	0.0105	0.0360	0.0118	0.0600	0.0151	0.0200	0.0088	0.0720	0.0164	0.0320	0.0112	0.0240	0.0097	0.0160	0.0080	0.0056	0.0160	0.0080	0.0056	0.0560	0.0146		
LLaMAX Llama 2 7B Alpaca	0.050	0.080	0.0124	0.0369	0.0118	0.0600	0.0197	0.0200	0.0081	0.0640	0.0155	0.0320	0.0112	0.0240	0.0135	0.0400	0.0135	0.0369	0.0080	0.0320	0.0112	0.0700	0.0168		
MaLA-500 Llama 2 10B v1	0.099	0.080	0.0090	0.0090	0.0090	0.0120	0.0089	0.0200	0.0089	0.0240	0.0097	0.0120	0.0069	0.0200	0.0089	0.0400	0.0090	0.0000	0.0340	0.0040	0.0040	0.0040	0.0040		
MaLA-500 Llama 2 10B v2	0.0991	0.0800	0.0000	0.0000	0.0000	0.0120	0.0089	0.0200	0.0089	0.0240	0.0097	0.0120	0.0141	0.0560	0.0141	0.0200	0.0089	0.0400	0.0000	0.0340	0.0040	0.0040	0.0040	0.0040	
YaYi Llama 2 7B	0.0709	0.0320	0.0112	0.0840	0.0230	0.1600	0.0230	0.0160	0.0232	0.0160	0.0193	0.0150	0.0141	0.0200	0.0180	0.0480	0.0135	0.0100	0.0190	0.0120	0.0069	0.0040	0.0040	0.0120	
TowerBase Llama 2 7B	0.0615	0.0240	0.0097	0.0840	0.0176	0.1160	0.0203	0.0920	0.0185	0.0880	0.0180	0.0480	0.0135	0.0100	0.0190	0.0120	0.0069	0.0056	0.0160	0.0080	0.0080	0.0180			
TowerInstruct Llama 2 7B	0.0724	0.0180	0.0088	0.0840	0.0190	0.1520	0.0228	0.1560	0.0200	0.1380	0.0212	0.0460	0.0155	0.0150	0.0240	0.0180	0.0089	0.0200	0.0089	0.0200	0.0089	0.0720	0.0164		
Occiglot Mistral 7B v0.1	0.0682	0.0240	0.0120	0.1200	0.0230	0.1600	0.0200	0.0800	0.0200	0.0200	0.0200	0.0200	0.0155	0.0200	0.0155	0.0200	0.0155	0.0200	0.0105	0.1120	0.0200	0.0320	0.0294		
Occiglot Mistral 7B v0.1 Instruct	0.2276	0.0480	0.0135	0.3400	0.0300	0.4640	0.0316	0.4000	0.0310	0.3160	0.0295	0.1840	0.0246	0.2360	0.0269	0.0640	0.0155	0.0280	0.0105	0.1120	0.0200	0.0320	0.0294		
BLOOM 7B	0.087	0.0240	0.0097	0.1600	0.0080	0.0400	0.0124	0.0360	0.0118	0.0200	0.0069	0.0200	0.0124	0.0340	0.0089	0.0360	0.0112	0.0280	0.0105	0.0120	0.0069	0.0240	0.0097		
BLOOMZ 7B	0.0255	0.0320	0.0112	0.0160	0.0080	0.0360	0.0112	0.0160	0.0080	0.0320	0.0112	0.0280	0.0105	0.0360	0.0112	0.0280	0.0105	0.0120	0.0069	0.0240	0.0097	0.0118			
mGPT	0.0171	0.0200	0.0080	0.0080	0.0080	0.0120	0.0080	0.0080	0.0080	0.0120	0.0080	0.0120	0.0080	0.0240	0.0080	0.0080	0.0080	0.0000	0.0000	0.0000	0.0000	0.0056			
ngGPT 1.3B	0.0131	0.0200	0.0080	0.0080	0.0080	0.0120	0.0080	0.0080	0.0080	0.0120	0.0080	0.0120	0.0080	0.0240	0.0080	0.0080	0.0080	0.0000	0.0000	0.0000	0.0000	0.0059			
XGLM 7.5B	0.0102	0.0000	0.0000	0.0080	0.0056	0.0000	0.0000	0.0120	0.0099	0.0000	0.0000	0.0040	0.0040	0.0089	0.0200	0.0089	0.0200	0.0089	0.0000	0.0000	0.0000	0.0000	0.0135		
YaYi 7B	0.0726	0.0240	0.0097	0.1600	0.0089	0.0600	0.0151	0.0240	0.0097	0.0560	0.0146	0.0160	0.0089	0.0120	0.0069	0.0240	0.0097	0.0040	0.0000	0.0000	0.0000	0.0000	0.0056		
Llama 3 8B	0.2745	0.1760	0.0240	0.3960	0.0310	0.5080	0.0317	0.4720	0.0316	0.3640	0.0305	0.0360	0.0118	0.3760	0.0307	0.4000	0.0271	0.0560	0.0146	0.3680	0.0300	0.0280	0.0105		
Llama 3.1 8B	0.288	0.2360	0.0240	0.3960	0.0310	0.5080	0.0317	0.4720	0.0316	0.3640	0.0305	0.0360	0.0118	0.3760	0.0307	0.4000	0.0271	0.0560	0.0146	0.3680	0.0300	0.0280	0.0105		
Gemma 2 8B	0.3232	0.2440	0.0301	0.4320	0.0312	0.5580	0.0318	0.4890	0.0317	0.3660	0.0310	0.3160	0.0237	0.4120	0.0311	0.3760	0.0307	0.2720	0.0282	0.2800	0.0314	0.2920	0.0288		
Gemma 2 9B	0.3295	0.2960	0.0289	0.0400	0.0311	0.5640	0.0318	0.5080	0.0317	0.3600	0.0304	0.2120	0.0269	0.3520	0.0303	0.3760	0.0307	0.3800	0.0400	0.0310	0.0440	0.0310	0.0430		
Qwen 2 7B	0.4895	0.4440	0.0315	0.6560	0.0301	0.8090	0.0250	0.7600	0.0271	0.6960	0.0292	0.0160	0.0088	0.6720	0.0289	0.1600	0.0232	0.1480	0.0225	0.5600	0.0315	0.6460	0.0316		
Qwen 1.5 7B	0.3156	0.1240	0.0200	0.0400	0.0315	0.5520	0.0315	0.4840	0.0317	0.4520	0.0315	0.1680	0.0237	0.3420	0.0314	0.2720	0.0164	0.0400	0.0124	0.2840	0.0286	0.0420	0.0313		
EMMA-500 Llama 2 7B	0.1702	0.0880	0.018	0.2320	0.0268	0.3400	0.0300	0.2800	0.0285	0.2560	0.0277	0.0920	0.0183	0.2280	0.0261	0.0800	0.0172	0.2280	0.0266	0.2120	0.0259	0.0240	0.0190		

Table 21: 3-shot results (ACC) on MGSM by CoT prompting and flexible matching.

Model	Avg	bn	bn-stder	de	de-stder	en	en-stder	es	es-stder	fr	fr-stder	ja	ja-stder	ru	ru-stder	sw	sw-stder	te	te-stder	th	th-stder	zh	zh-stder
Llama 2 7B	0.0636	0.0200	0.0089	0.0760	0.0168	0.1600	0.0232	0.1240	0.0209	0.0920	0.0183	0.0360	0.0118	0.0800	0.0212	0.0200	0.0089	0.0088	0.0056	0.0160	0.0080	0.0172	
Llama 2 7B Chat	0.1091	0.240	0.0097	0.0760	0.0241	0.2720	0.0282	0.1920	0.0250	0.1880	0.0248	0.0460	0.0124	0.1520	0.0228	0.2000	0.0089	0.0088	0.0056	0.0280	0.0105	0.1160	0.0203
CodeLlama 2 7B	0.0664	0.160	0.0088	0.0720	0.0183	0.1320	0.0225	0.1080	0.0197	0.1080	0.0197	0.0460	0.0130	0.0760	0.0168	0.120	0.0069	0.0160	0.0060	0.0140	0.0400	0.0124	
LLaMAX Llama 2 7B	0.0635	0.280	0.0112	0.0840	0.0135	0.1520	0.0228	0.0960	0.0187	0.0520	0.0141	0.0440	0.0130	0.0360	0.0181	0.1480	0.0135	0.0480	0.0089	0.0080	0.0056	0.0160	
LLaMAX Llama 2 7B Alpaca	0.0673	0.0300	0.0000	0.0040	0.0040	0.0040	0.0040	0.0040	0.0040	0.0160	0.0080	0.0160	0.0088	0.0240	0.0097	0.0040	0.0040	0.0040	0.0000	0.0000	0.0080	0.0056	
MaLA-500 Llama 2 10B v1	0.0673	0.0300	0.0000	0.0040	0.0040	0.0040	0.0040	0.0040	0.0040	0.0160	0.0080	0.0160	0.0089	0.0240	0.0097	0.0040	0.0040	0.0040	0.0000	0.0000	0.0080	0.0056	
MaLA-500 Llama 2 10B v2	0.0703	0.0300	0.0000	0.0040	0.0040	0.0040	0.0040	0.0040	0.0040	0.0160	0.0080	0.0160	0.0089	0.0240	0.0097	0.0040	0.0040	0.0040	0.0000	0.0000	0.0080	0.0056	
YaYi Llama 2 7B	0.0722	0.0240	0.0088	0.0840	0.0118	0.1600	0.0220	0.0760	0.0197	0.0800	0.0169	0.0460	0.0135	0.0720	0.0160	0.0240	0.0089	0.0240	0.0089	0.0320	0.0112	0.0205	
TowerInstruct Llama 2 7B	0.0824	0.120	0.0089	0.0840	0.0215	0.1880	0.0248	0.1520	0.0228	0.1200	0.0206	0.0460	0.0097	0.1360	0.0217	0.160	0.0089	0.0200	0.0069	0.0320	0.0112	0.0197	
Occiglot Mistral 7B v0.1	0.1487	0.0240	0.0097	0.1600	0.0261	0.3320	0.0298	0.4200	0.0264	0.3640	0.0264	0.0460	0.0155	0.1720	0.0239	0.0320	0.0112	0.0210	0.0069	0.0600	0.0151	0.1120	0.0200
Occiglot Mistral 7B v0.1 Instruct	0.2216	0.0400	0.0124	0.3320	0.0298	0.4320	0.0314	0.4200	0.0313	0.3160	0.0295	0.1680	0.0237	0.2440	0.0272	0.560	0.0146	0.0160	0.0080	0.0840	0.0176	0.2480	0.0274
BLOOM 7B	0.2441	0.2425	0.0245	0.4543	0.0314	0.4200	0.0320	0.3540	0.0315	0.3600	0.0320	0.0400	0.0234	0.2440	0.0274	0.560	0.0146	0.0160	0.0080	0.0840	0.0176	0.2480	0.0274
BLOOMZ 7B	0.3932	0.4543	0.4367	0.4151	0.2315	0.2315	0.2315	0.															

Table 23: 5-shot results (ACC) on ARC multilingual.

Model	Avg	High	Medium	Low
Llama 2 7B	0.2756	0.3312	0.2731	0.2102
Llama 2 7B Chat	0.2802	0.3369	0.2779	0.2129
CodeLlama 2 7B	0.2523	0.2886	0.2464	0.2165
LLaMAX Llama 2 7B	0.2609	0.3000	0.2592	0.2148
LLaMAX Llama 2 7B Alpaca	0.3106	0.3689	0.3185	0.2249
MaLA-500 Llama 2 10B v1	0.2116	0.2192	0.2048	0.2132
MaLA-500 Llama 2 10B v2	0.2116	0.2192	0.2048	0.2132
YaYi Llama 2 7B	0.2840	0.3430	0.2835	0.2111
TowerBase Llama 2 7B	0.2794	0.3532	0.2682	0.2051
TowerInstruct Llama 2 7B	0.3010	0.3888	0.2885	0.2116
Occiglot Mistral 7B v0.1	0.2977	0.3839	0.2851	0.2103
Occiglot Mistral 7B v0.1 Instruct	0.3088	0.4029	0.2965	0.2113
BLOOM 7B	0.2365	0.2627	0.2272	0.2189
BLOOMZ 7B	0.2395	0.2694	0.2274	0.2218
mGPT	0.2024	0.2011	0.1965	0.2137
mGPT 13B	0.2176	0.2299	0.2114	0.2123
XGLM 7.5B	0.2221	0.2461	0.2105	0.2110
YaYi 7B	0.2444	0.2796	0.2329	0.2191
Llama 3 8B	0.3480	0.4243	0.3553	0.2406
Llama 3.1 8B	0.3493	0.4243	0.3589	0.2400
Gemma 7B	0.3868	0.4646	0.4047	0.2606
Gemma 2 9B	0.4415	0.5459	0.4618	0.2782
Qwen 2 7B	0.3382	0.4388	0.3264	0.2317
Qwen 1.5 7B	0.2893	0.3555	0.2814	0.2192
EMMA-500 Llama 2 7B	0.2953	0.3410	0.2982	0.2334

Table 24: Results on Multipl-E. For language-level breakdowns, refer to Tables 25 to 27 in the appendix.

Model	Avg Pass@1	Avg Pass@10	Avg Pass@25
Llama 2 7B	8.92%	19.45%	25.68%
CodeLlama 2 7B	28.43%	50.83%	63.92%
LLaMAX 2 7B	0.35%	1.61%	2.67%
MaLA-500 Llama 2 10B V2	0.0%	0.0%	0.0%
TowerBase Llama 2 7B	3.61%	6.65%	8.97%
Occiglot Mistral 7B v0.1	21.26%	31.37%	45.86%
Bloom 7B	5.34%	10.49%	14.65%
BloomZ 2 7B	5.85%	11.40%	15.76%
Aya23 8B	9.19%	23.52%	32.09%
Mistral 7B v0.3	26.10%	48.68%	59.05%
Llama 3 8B	30.09%	53.82%	64.01%
LLaMAX 3 8B	3.00%	7.23%	10.67%
Gemma 7B	28.55%	54.27%	64.75%
CodeGemma 7B	31.51%	63.13%	72.65%
Qwen 1.5 7B	21.05%	37.19%	47.63%
Qwen 2 7B	38.68%	62.63%	71.55%
EMMA-500 Llama 2 7B	11.38%	19.02%	26.16%

languages. Despite its benefits, CPT can lead to catastrophic forgetting, where models lose previously learned information. To address the potential degraded performance issue, Ibrahim et al. (2024) presents a simple yet effective approach to continual pre-train models, demonstrating that with a combination of learning rate re-warming, re-decaying, and replay of previous data, it is possible to match the performance of fully re-trained models. Also, recent studies have delved into the effectiveness of continual pre-training with parallel data. As highlighted in the study by Gilabert et al. (2024), the use of a Catalan-centric parallel dataset has enabled the training of models good at translating in various directions. Besides, the research by Kondo et al. (2024), proposed a two-phase continual training approach with parallel data. In the first phase, a pre-trained LLM is continually pre-trained on parallel data, followed by a second phase of supervised fine-tuning with a small amount of high-quality parallel data. Their experiments with a 3.8B-parameter model across various data formats revealed that alternating between source and target sentences during continual pre-training is crucial for enhancing translation accuracy in the corresponding direction.

1998
1999
2000
2001
2002

2003
2004
2005
2006
2007

2008
2009
2010

2011

2012

2013

2014

2015

2016

2017

2018

2019

2020

2021

2022

2023

2024

2025

2026

2027

2028

2029

2030

2031

2032

2033

2034

2035

2036

2037

2038

2039

2040

2041

2042

2043

2044

2045

2046

2047

2048

2049

2050

2051

Table 25: Pass@1 on Multipl-E.

Model	C++	C#	Java	JavaScript	Python	Rust	TypeScript
Llama 2 7B	6.74%	5.65%	8.54%	11.34%	11.98%	6.12%	12.04%
CodeLlama 2 7B	26.70%	20.44%	30.56%	32.89%	28.76%	26.23%	33.45%
LLaMAX 2 7B	0.0%	0.0%	0.02%	0.68%	0.0%	0.0%	1.78%
MaLA-500 Llama 2 10B V2	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
TowerBase Llama 2 7B	1.01%	6.85%	7.86%	0.87%	8.68%	0.0%	0.0%
Occiglot Mistral 7B v0.1	23.16%	16.14%	19.83%	26.92%	21.45%	16.34%	24.97%
Bloom 7B	5.25%	2.97%	6.37%	6.76%	7.65%	1.01%	7.34%
BloomZ 2 7B	6.87%	3.59%	6.02%	6.96%	7.33%	2.13%	8.04%
Aya23 8B	16.03%	7.91%	14.29%	6.23%	3.56%	11.43%	4.90%
Mistral 7B v0.3	26.12%	22.87%	25.54%	35.24%	24.91%	19.24%	28.76%
Llama 3 8B	34.12%	21.06%	26.56%	36.12%	30.22%	25.19%	37.34%
LLaMAX 3 8B	0.0%	0.0%	0.0%	9.73%	0.91%	0.21%	10.12%
Gemma 7B	29.66%	21.40%	27.35%	35.29%	30.11%	25.48%	30.57%
CodeGemma 7B	33.91%	21.05%	29.43%	37.78%	32.56%	28.70%	37.14%
Qwen 1.5 7B	22.04%	12.42%	17.68%	27.58%	32.11%	8.32%	27.21%
Qwen 2 7B	43.51%	21.47%	38.95%	46.31%	37.49%	34.61%	48.41%
EMMA-500 Llama 2 7B	11.34%	8.94%	11.93%	11.67%	18.97%	6.86%	9.94%

Table 26: Pass@10 on Multipl-E.

Model	C++	C#	Java	JavaScript	Python	Rust	TypeScript
Llama 2 7B	17.92%	14.22%	20.78%	23.12%	24.86%	13.08%	22.19%
CodeLlama 2 7B	51.39%	37.13%	50.38%	59.08%	56.58%	47.68%	53.60%
LLaMAX 2 7B	0.96%	0.0%	0.53%	3.99%	0.0%	0.12%	5.64%
MaLA-500 Llama 2 10B V2	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
TowerBase Llama 2 7B	1.01%	6.85%	7.86%	0.87%	8.68%	0.0%	0.0%
Occiglot Mistral 7B v0.1	36.08%	27.49%	35.87%	47.45%	41.76%	29.80%	43.13%
Bloom 7B	11.35%	6.86%	12.79%	12.44%	14.24%	3.68%	12.08%
BloomZ 2 7B	12.55%	8.21%	12.94%	14.10%	14.84%	3.91%	13.27%
Aya23 8B	28.49%	17.34%	27.12%	23.19%	26.72%	26.14%	15.67%
Mistral 7B v0.3	50.23%	35.18%	46.83%	57.23%	54.29%	42.77%	54.22%
Llama 3 8B	55.13%	36.67%	54.34%	62.65%	59.24%	45.68%	63.02%
LLaMAX 3 8B	0.72%	0.0%	0.97%	22.91%	1.58%	1.38%	23.04%
Gemma 7B	55.21%	39.09%	52.02%	61.88%	60.09%	50.34%	61.23%
CodeGemma 7B	62.29%	45.11%	59.74%	70.96%	69.76%	61.27%	72.76%
Qwen 1.5 7B	40.11%	28.95%	37.56%	48.35%	40.34%	20.19%	44.85%
Qwen 2 7B	64.58%	43.17%	63.28%	73.21%	54.34%	65.48%	74.32%
EMMA-500 Llama 2 7B	22.84%	14.29%	21.28%	18.22%	24.94%	15.68%	15.91%

Table 27: Pass@25 on Multipl-E.

Model	C++	C#	Java	JavaScript	Python	Rust	TypeScript
Llama 2 7B	24.89%	18.77%	26.45%	30.31%	31.09%	18.70%	29.55%
CodeLlama 2 7B	64.12%	45.96%	61.98%	72.76%	71.77%	59.98%	70.89%
LLaMAX 2 7B	1.47%	0.0%	1.04%	6.78%	0.0%	0.34%	9.04%
MaLA-500 Llama 2 10B V2	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
TowerBase Llama 2 7B	1.86%	15.74%	22.45%	2.53%	20.18%	0.0%	0.0%
Occiglot Mistral 7B v0.1	44.17%	33.89%	42.77%	56.63%	49.87%	37.94%	55.76%
Bloom 7B	17.27%	10.07%	17.23%	17.55%	17.93%	5.70%	16.81%
BloomZ 2 7B	17.41%	11.98%	18.66%	18.37%	18.78%	6.83%	18.31%
Aya23 8B	37.56%	21.87%	36.45%	33.57%	36.09%	35.34%	23.75%
Mistral 7B v0.3	60.93%	44.80%	56.12%	66.34%	65.33%	55.31%	64.51%
Llama 3 8B	64.42%	43.34%	62.14%	73.51%	71.82%	57.78%	75.09%
LLaMAX 3 8B	1.71%	0.0%	2.32%	35.61%	2.77%	2.08%	30.19%
Gemma 7B	66.14%	47.08%	63.24%	70.22%	70.47%	63.52%	72.56%
CodeGemma 7B	73.66%	50.79%	69.96%	79.48%	78.68%	74.87%	81.12%
Qwen 1.5 7B	51.87%	37.69%	48.71%	59.71%	49.69%	29.78%	55.98%
Qwen 2 7B	74.33%	51.87%	72.03%	81.67%	65.38%	74.97%	80.59%
EMMA-500 Llama 2 7B	31.93%	18.69%	29.85%	26.45%	32.55%	21.32%	22.34%