

# MAPLE: Multilingual Evaluation of Parameter Efficient Finetuning of Large Language Models

Anonymous ACL submission

## Abstract

Parameter Efficient Finetuning (PEFT) has emerged as a viable solution for improving the performance of Large Language Models (LLMs) without requiring massive resources and compute. Prior work on multilingual evaluation has shown that there is a large gap between the performance of LLMs on English and other languages. Further, there is also a large gap between the performance of smaller open-source models and larger LLMs. Finetuning can be an effective way to bridge this gap and make language models more equitable. In this work, we finetune the LLAMA-2-7B and MISTRAL-7B models on two synthetic multilingual instruction tuning datasets to determine its effect on model performance on six downstream tasks covering forty languages in all. Additionally, we experiment with various parameters, such as rank for low-rank adaptation and values of quantisation to determine their effects on downstream performance and find that higher rank and higher quantisation values benefit low-resource languages. We find that PEFT of smaller open-source models sometimes bridges the gap between the performance of these models and the larger ones, however, English performance can take a hit. We also find that finetuning sometimes improves performance on low-resource languages, while degrading performance on high-resource languages.

## 1 Introduction

Large Language Models (LLMs) show impressive performance on several tasks, sometimes even surpassing human performance. This has been attributed to the vast amounts of training data used during the pretraining phase, as well as various techniques used to align the models during the finetuning phase. Several variants of finetuning exist, including supervised finetuning (SFT) and instruction tuning (OpenAI, 2023), where the model is finetuned with task specific data, or instructions on

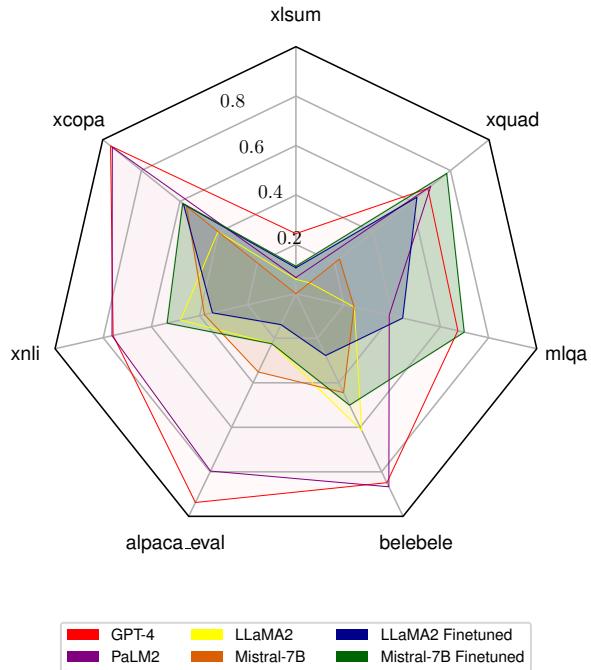


Figure 1: Comparison of best parameter efficient instruction finetuned models with other off the shelf LLMs. Notably, the best Mistral instruction finetuned model is able to outperform even GPT-4 and PaLM2 on “MLQA” and “XQUAD” tasks.

how to perform tasks. However, finetuning all the parameters of the model can be expensive and time consuming, due to the increasing size of language models in recent years. Parameter Efficient Fine-tuning (PEFT) has emerged as a viable alternative to full finetuning (Chen et al., 2023).

Most studies on LLMs focus on training, finetuning and evaluating models in performing tasks in English. Recent work on comprehensive evaluation of LLM capabilities in non-English settings (Ahuja et al., 2023a) have shown that LLMs perform far worse on languages other than English. Studies that compare multilingual performance across different models (Ahuja et al., 2023b), show that there is a large performance gap between large, proprietary

and closed models such as GPT-4 and PaLM2 and smaller open-source models like LLAMA-2-7B and MISTRAL-7B.

PEFT techniques like LoRA (Hu et al., 2022) have been shown to strengthen the multilingual capabilities of these open-source LLMs (Zhao et al., 2024). Moreover, Adapters (Pfeiffer et al., 2020a) have been proposed to boost the capabilities of language models to newer languages. Since full finetuning of models is not always feasible due to resource and compute constraints, exploring how far PEFT techniques can take us in boosting performance on non-English languages is a promising direction. Adoption of model quantisation techniques (Dettmers et al., 2022; yang Liu et al., 2023) has also made PEFT LLMs more accessible.

Not much work has been done on analyzing the impact of different choices, configurations and settings for PEFT on multilingual downstream tasks. In this work we aim to analyse how LoRA rank and quantisation affects the performance of finetuned models across 6 downstream tasks, covering 40 languages in all. We are interested in knowing whether multilingual PEFT can lead to reasonable gains in performance, or whether full finetuning of models is required. Furthermore, we check if it is better to generate multilingual instructions from larger base models or translate existing english instruction to more languages. Lastly, we also study the effect of multilingual finetuning on English performance.

Prior works have demonstrated that LoRA finetuning is the most effective PEFT out of existing techniques so far (Zhuo et al., 2024). Hence, we keep our study limited to analysing different LoRA and QLoRA (Dettmers et al., 2023) configurations to answer our research questions. Our contributions are as follows:

- We benchmark effects of various ranks and quantisation with LLAMA-2-7B and MISTRAL-7B models finetuned on MULTIALPACA and BACTRIAN-X-22 dataset. We analyse the effects of % of trainable parameters and quantisation on 6 various tasks and 40 languages.
- We study efficacy of finetuning by comparing results with non-finetuned models of similar or larger sizes.
- We analyse the effects of multilingual PEFT on English performance to check for degradations due to forgetting.

- We experiment with the choice of instruction finetuning dataset to study any variations in model performance on our downstream tasks.
- We present results and an analysis of trends across these models and instruction finetuning datasets with directions for future research.

## 2 Related Work

**Parameter Efficient Finetuning:** Recently, Parameter Efficient Finetuning has gained significant attention in the NLP research community since full finetuning of LLMs is prohibitively expensive for most organizations. Following early works on adapters (Houlsby et al., 2019; Pfeiffer et al., 2020a), several finetuning techniques like LoRA (Hu et al., 2022), (IA)<sup>3</sup> (Liu et al., 2022a), P-Tuning (Liu et al., 2022b) and Prefix Tuning (Li and Liang, 2021) have been proposed. These techniques make the compute costs manageable by significantly reduce the number of trainable parameters during finetuning. Several works have used these techniques for efficient cross lingual transfer (Ansell et al., 2022), to tackle catastrophic forgetting (Vu et al., 2022) or compose multiple adapters (Pfeiffer et al., 2021) for multi-task performance.

**Multilingual Instruction Finetuning:** Recently, there has been a lot of interest in creating multilingual instruction finetuning datasets to enhance the reasoning capabilities of LLMs on languages other than English (Li et al., 2023a; Wei et al., 2023). Such datasets, are being used to create LLMs that can serve to larger demographic and can also be more efficient during inference time (Jiang et al., 2023). Li et al. (2023a) gained significant performance on multilingual tasks by LoRA finetuning LLaMA and BLOOM on the BACTRIAN-X dataset. These instruction datasets are generally derived by generation or translation. Wei et al. (2023) translated seed tasks of ALPACA dataset (Taori et al., 2023) to 11 languages and then prompted GPT-3.5-Turbo to generated more instructions in those languages, while Li et al. (2023a) translated ALPACA instructions to 50 other languages using google translate API and generated responses using GPT-3.5-Turbo. Moreover, efforts like (Singh et al., 2024) attempts to create crowdsourced multilingual instruction datasets to capture better linguistic and cultural nuances.

156 **Quantisation for Model Compression:** Model  
157 quantisation is another way of reducing the overall  
158 memory footprint of the LLM. While many popular  
159 LLMs (notably LLAMA-2-7B (Touvron et al.,  
160 2023) and MISTRAL-7B (Jiang et al., 2023)) are  
161 pre-trained with weights represented in 16 bit float-  
162 ing point numbers (Wu et al., 2020), it is shown  
163 that finetuning with lower quantisation yields sim-  
164 ilar performance. The most popular quantisation  
165 techniques – LLM:Int8() (Dettmers et al., 2022)  
166 and 4 bit (yang Liu et al., 2023) are usually com-  
167 bined with LoRA (Dettmers et al., 2023) to further  
168 reduce the memory footprint of LLM finetuning.

169 **LLM Evaluation:** Principled LLM evaluation  
170 has gained significant interest with demonstrations  
171 of increasingly complex abilities of LLMs (Brown  
172 et al., 2020; Cobbe et al., 2021; Wei et al., 2022; Shi  
173 et al., 2023) on various tasks. However, many eval-  
174 uations are monolingual or English-only and mul-  
175 tilingual evaluation of LLMs (Ahuja et al., 2023a;  
176 Asai et al., 2023; Ahuja et al., 2023b) remains a  
177 challenging problem. Past work by Ramesh et al.  
178 (2023) has evaluated the effects of model compres-  
179 sion techniques such as quantisation, distillation  
180 and pruning on LLMs performance on downstream  
181 tasks in multilingual setting.

### 182 3 Experiments

#### 183 3.1 Setup

184 **Finetuning Models:** We finetune open-source,  
185 multilingual LLMs on multilingual instruction fine-  
186 tuning datasets. We pick models that are pre-  
187 trained on multilingual data as it would be un-  
188 fair to compare English-only LLMs when finetuning  
189 on multilingual data. Specifically, we explore  
190 PEFT on LLAMA-2-7B (Touvron et al., 2023)  
191 and MISTRAL-7B (Jiang et al., 2023) models.

192 **Finetuning Dataset:** We finetuned our  
193 models on MULTIALPACA (Wei et al., 2023) and  
194 BACTRIAN-X (Li et al., 2023a) datasets for all  
195 our experiments.

196 **MULTIALPACA** is a self instruct dataset which  
197 follows the same approach as (English-only) AL-  
198 PACA dataset (Taori et al., 2023) by translating  
199 seed tasks to 11 languages and then using GPT-3.5-  
200 Turbo for response collection. The languages in-  
201 cluded in the dataset are Arabic, German, Spanish,  
202 French, Indonesian, Japanese, Korean, Portuguese,  
203 Russian, Thai and Vietnamese.

204 **BACTRIAN-X** is a machine translated dataset of  
205 the original alpaca-52k and dolly-15k (Conover  
206 et al., 2023) datasets. In this dataset, the in-  
207 structions were translated using google translate  
208 API to 52 diverse languages and responses were  
209 generated using GPT-3.5-Turbo. In our exper-  
210 iments we finetune our models on a subset of  
211 22 languages namely, Afrikaans, Arabic, Ben-  
212 gali, Chinese-Simplified, Dutch, French, German,  
213 Gujarati, Hindi, Indonesian, Japanese, Korean,  
214 Marathi, Portuguese, Russian, Spanish, Swahili,  
215 Tamil, Telugu, Thai, Urdu and Vietnamese. We  
216 rename this dataset to BACTRIAN-X-22. Each  
217 language in BACTRIAN-X-22 consists of 67k in-  
218 structions parallel and responses whereas MULTI-  
219 ALPACA consists of nearly 100k instructions. To  
220 get a dataset of comparable size, we take 10% in-  
221 structions (6700) from each language giving us  
222 close to 150k instructions in 22 languages. We first  
223 shuffle and partition indices 0-67k and divide them  
224 into 10 partitions. Each partition now consists of  
225 random indices from 0-67k. Then we iterate over  
226 all 22 languages assigning language  $i$  to partition  
227  $i \bmod 10$ . This gives us a partition of 6700 indices  
228 from each languages which we use to form the in-  
229 struction tuning dataset from each language. This  
230 means that every instruction at index from 0-67k  
231 is included in at least two languages. We study  
232 whether this enhanced sampling ensuring at least  
233 two languages per instruction helps in cross-lingual  
234 transfer.

235 **Finetuning Techniques:** We follow the LoRA  
236 (Dettmers et al., 2023; Hu et al., 2022) finetuning  
237 recipe for each finetuning run. We finetune mod-  
238 els on various ranks and quantisations, specifically  
239 LoRA Ranks 8, 16, 32, 64 and 128 and 4bit, 8bit  
240 and 16bit quantisation.

#### 241 3.2 Evaluation

242 We evaluate multilingual capabilities of our fine-  
243 tuned models on three Classification tasks, two  
244 Question-Answering tasks and one Summarisation  
245 task. We use prompts that are similar to those  
246 used in the MEGA benchmarking study (Ahuja  
247 et al., 2023a) but adapted to the Alpaca-style (Taori  
248 et al., 2023) instruction format. We study the im-  
249 pact of multilingual finetuning on English capabili-  
250 ties using Alpaca Eval (Li et al., 2023b). We use  
251 lm-eval-harness (Gao et al., 2021) for the evalua-  
252 tions. LM-Evaluation-Harness is a unified frame-  
253 work for few shot evaluation of language models.

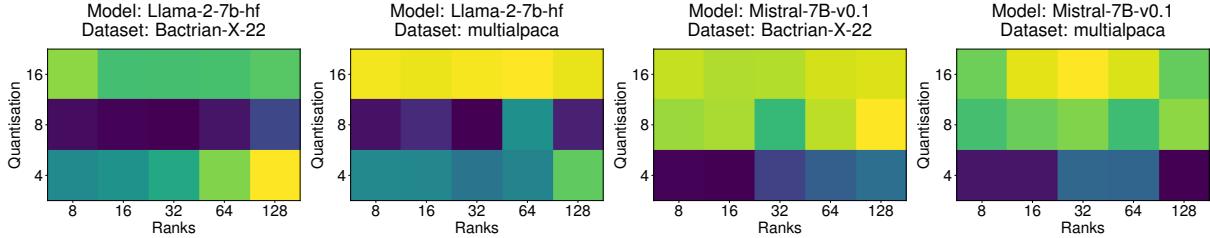


Figure 2: Average model performance of LLAMA-2-7B and MISTRAL-7B finetuned on BACTRIAN-X-22 and MULTIALPACA across tasks on all rank-quantisation configurations.

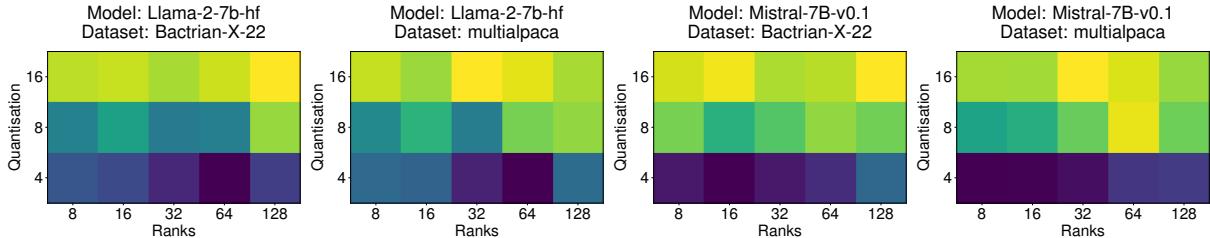


Figure 3: Belebele evaluation results of LLAMA-2-7B and MISTRAL-7B finetuned on BACTRIAN-X-22 and MULTIALPACA across tasks on all rank-quantisation configurations.

This framework standardises the inference and few shot example selection pipeline across tasks and models. We created the task configurations from MEGAVERSE (Ahuja et al., 2023b) with the ALPACA-style prompt template. We mention these prompts in detail in the Appendix Section C.

**Classification Datasets:** As part of our evaluation process, we benchmark our finetuned models on several datasets. The **Belebele** dataset (Bopardkar et al., 2023), which is parallel across 122 languages, is evaluated on a subset of 33 languages. We report our results on 30% of the test split in the zero-shot setting due to resource constraints. We report the results in Table 4, 5 6, 7, 8 and 9. The **XNLI** dataset (Conneau et al., 2018) consists of 122k training, 2490 validation, and 5010 test examples in 15 languages. We have evaluated our models on 1000 examples from test split with 4 in-context examples sampled from the validation split and report our results in Table 28, 29, 30, 31, 32 and 33. The **XCOPA** dataset (Ponti et al., 2020) covers 11 languages, and we evaluate our models on Estonian, Indonesian, Italian, Quechua, Thai and Vietnamese in the 4-shot setting similar to XNLI. We report our results in Table 22, 23, 24, 25, 26 and 27.

**Question Answering Datasets:** The **MLQA** dataset (Lewis et al., 2020) contains 5K extractive question-answering instances in 7 languages. For the interest of time, we evaluate our models for

1000 examples of the test split in a 4-shot setting and report our results in Table 16, 17, 18, 19, 20 and 21. The **XQuAD** dataset (Artetxe et al., 2020) consists of a subset of 240 paragraphs and 1190 question-answer pairs across 11 languages. We use a 4-shot setting similar to MLQA and evaluate 1000 examples of the test split. We report our results in Table 34, 35, 36, 37, 38 and 39.

**Summarisation Dataset:** The **XSUM** dataset (Hasan et al., 2021) spans 45 languages, and we evaluate our models in Arabic, Chinese-Simplified, English, French, Hindi, Japanese and Spanish. We evaluate our models on 100 text-summarization pairs from the test split in a zero-shot setting and report our results in Table 40, 42, 41, 43, 44 and 45.

**English Instruction Following Dataset:** We also use **AlpacaEval** (Li et al., 2023b) to benchmark English proficiency. We evaluate our models against text-davinci-003 responses on 800 instructions and use GPT4 (gpt-4-32k) as the evaluator. We report our results in Tables 10, 11, 12, 13, 14 and 15.

We discuss more about these benchmark datasets in detail in Appendix Section B.

## 4 Analysis of Results

**Analysis of Rank and Quantisation** In this study we aim to analyse the trade-offs between

254  
255  
256  
257  
258  
259

260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279

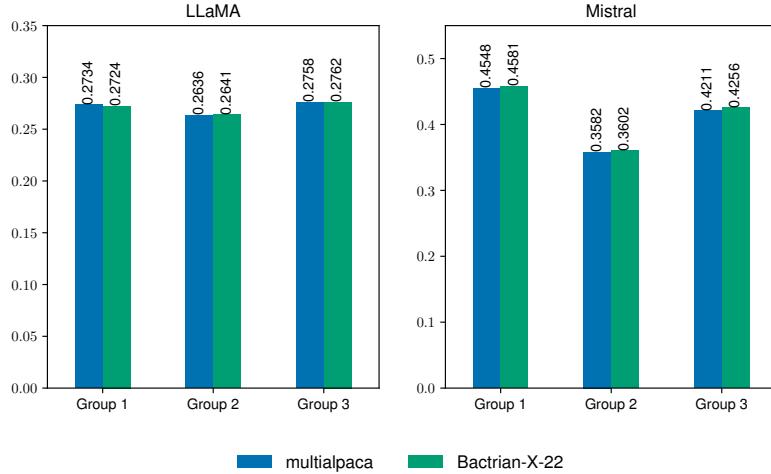
280  
281  
282  
283

284  
285  
286  
287  
288  
289  
290  
291

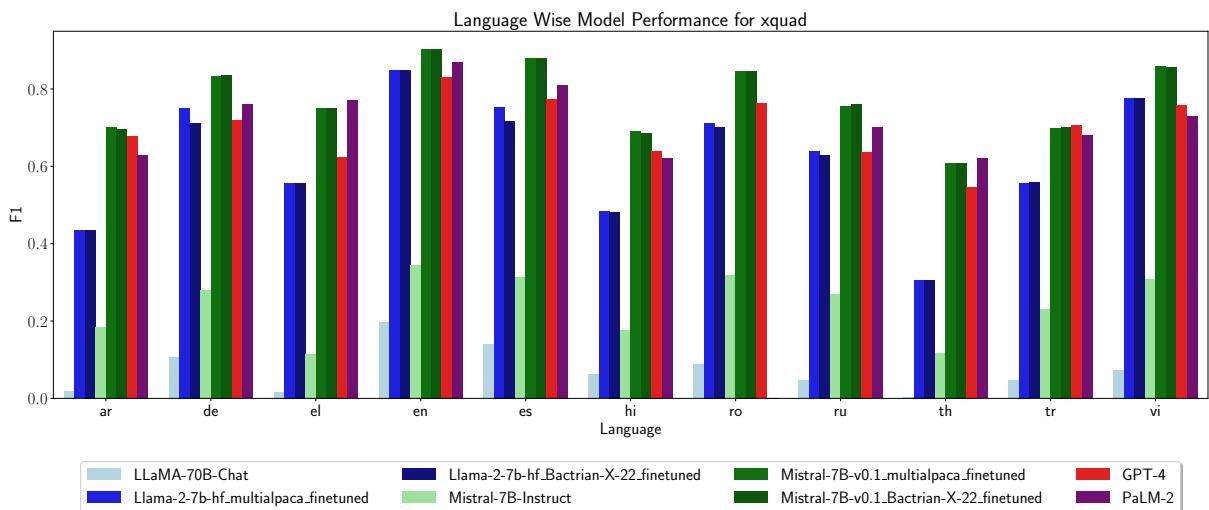
292  
293  
294  
295  
296  
297  
298  
299

300  
301  
302  
303  
304  
305  
306  
307  
308

309  
310  
311



**Figure 4: Effect of diversity of languages in fine-tuning on downstream task (belebele).** Here Group 1 is the set of 11 languages from MULTIALPACA, Group 2 is the set of 11 languages in BACTRIAN-X-22 but not in MULTIALPACA and Group 3 contains 13 languages present in neither. We find that both models trained on either datasets perform very similar to each other across all 3 groups. Additional details in Tables 4 to 9.



**Figure 5: Detailed language-wise comparison of our finetuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Arabic, German, Greek, English, Spanish, Hindi, Romanian, Russian, Thai, Turkish and Vietnamese for XQUAD (Artetxe et al., 2020).**

cost of compute and model performance. Both the LLAMA-2-7B and MISTRAL-7B models were finetuned on all rank-quantisation configurations using the MULTIALPACA and BACTRIAN-X-22 datasets resulting into 60 models.

We evaluate our finetuned models across the six benchmarking datasets mentioned in Section 3.2. We present the averaged results across these datasets in Fig 2. Lighter colours (yellow) indicate higher performance and it decreases with darker shades (blue). For MISTRAL-7B we can see a clear trend for both the finetuning datasets. Decreasing the quantisation can lead to a hit in model performance. For LLAMA-2-7B the trend is not very clear but the highest quantisation gives the best results. Additionally, higher ranks seem to give slightly better performance. According to our studies, using Rank 32 or 64 with 16bit Quantisa-

tion works the best on average. This can be inferred very clearly from our results on Belebele in Fig 3. To delve deeper, we provide a detailed task-wise performance in Fig 14.

**MULTIALPACA v/s BACTRIAN-X-22 as Instruction Finetuning Dataset** Here, we aim to study the model performance finetuned on multilingual instruction dataset created in 2 different settings i.e. LLM generated (MULTIALPACA) and machine-translated (BACTRIAN-X-22). In MULTIALPACA, both multilingual instructions and their responses are generated using GPT-3.5-Turbo from translated ALPACA seed instructions. While in BACTRIAN-X-22, the final set of ALPACA instructions were translated and then responses were collected using GPT-3.5-Turbo.

In our findings, we observe that models trained

312  
313  
314  
315  
316

317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329

330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345

347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
on both datasets give similar performance on average across tasks and languages, implying that method of instruction finetuning data creation has little to no effect on the model performance. Rather, we observe that multilingual capabilities of base model is a good indicator of the finetuned model performance across tasks. From Ahuja et al. (2023b) we know that MISTRAL-7B is a better base multilingual model than LLAMA-2-7B and we observe that multilingual capabilities of MISTRAL-7B also reaps greater benefits of multilingual instruction finetuning.

359  
360  
361  
362  
363  
Hence, in our experiments, we observe that for creating multilingual instruction datasets both approaches are equivalent - generating multilingual data from seed tasks or translating an existing English instruction dataset to more languages.

364  
365  
366  
367  
368  
369  
**Effect of Number of Languages in Training Data** Our finetuning datasets MULTIALPACA and BACTRIAN-X-22 have 11 and 22 languages respectively. We want to study if the additional 11 languages in BACTRIAN-X-22 help the model perform better at multilingual tasks.

370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
We do a case study on the “belebele” task which consists of 50+ languages. We sample 35 of these and divide them into 3 groups. The first group (Group 1) consists of 11 languages from MULTIALPACA, the next group (Group 2) consists of 11 languages present in BACTRIAN-X-22 but absent from MULTIALPACA. The final group (Group 3) contains 13 languages that are not present in any of our finetuning datasets. We compute average accuracy of finetuned models across all ranks on the 3 groups and present them in Figure 4. We find that the larger number of languages in BACTRIAN-X-22 do not necessarily help. This behavior is consistent with the findings of Shaham et al. (2024). Moreover, we observe that for other tasks as well BACTRIAN-X-22 dataset has no edge over MULTIALPACA as we can see from Table 1 and Figure 6.

388  
389  
390  
391  
392  
393  
394  
395  
396  
**Effect of Multilingual Finetuning on performance on downstream tasks for High Resource v/s Low Resource Languages** For XNLI and MLQA, we observe that finetuning improves performance on low-resource languages but worsens performance on high-resource languages. In Belebele, we find that finetuning worsens the performance for all languages for both models. For XCOPA, we get the same or better results for all lan-

397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
guages with finetuning. For XQUAD, multilingual finetuning boosts performance for all languages for both LLAMA-2-7B and MISTRAL-7B and even surpasses GPT-4 as seen in Fig 5. Moreover, for XLSUM parameter efficient finetuned models does not perform at par with fully finetuned or larger LLMs due to the generative nature of the task. This shows that overall, PEFT on smaller LLMs using multilingual instruction data can prove to be beneficial and can bridge the gap between smaller open-source models and large proprietary models.

408  
409  
410  
411  
412  
413  
414  
415  
Moreover, we observe that language-wise performance does not get affected by the choice of the training dataset, i.e. MULTIALPACA and BACTRIAN-X-22 have similar performances across languages and tasks. BACTRIAN-X-22 sometimes, leads to better performance on some low-resources languages as it includes more languages in the training data.

416  
417  
418  
419  
420  
421  
422  
423  
424  
**Analysis of Performance on English** In Table 2 we compare models on the English AlpacaEval benchmark. Overall, the capability of model to follow English instruction reduces drastically after multilingual finetuning. In general, finetuning on English-only ALPACA dataset or using higher capacity adapters (higher rank, better quantisation) seem to help in preserving the performance on English instruction finetuning.

425  
426  
427  
428  
429  
430  
More concisely, we observe that MISTRAL-7B is able to preserve more English capabilities than LLAMA-2-7B on AlpacaEval. While, BACTRIAN-X-22 and MULTIALPACA have difference in performance in English for the respective finetuned model.

431  
432  
433  
434  
435  
436  
437  
Furthermore, in a task-wise analysis we observe that multilingual finetuning leads to deterioration in English performance in Belebele, while in XNLI it deteriorates for LLAMA-2-7B and improves for MISTRAL-7B. For question-answering tasks MLQA and XQUAD, multilingual finetuning leads to improvement in performance. In XLSUM the performance improves for LLAMA-2-7B when finetuned on multilingual data, while for MISTRAL-7B the performance decreases or remains the same.

441  
442  
443  
444  
445  
446  
**Task Wise Performance Analysis** We analyse the average performance across languages on each task for GPT-4, PaLM-2, LLaMA-2-70B-Chat, Mistral-7B-Instruct and LLAMA-2-7B and MISTRAL-7B models finetuned using MULTIALPACA, ALPACA and BACTRIAN-X-22. In Figure 6

model	finetuning dataset	xnli	xcopa	xquad	belebele	mlqa	xlsum	Model Average
GPT-4	NA	0.75	0.90	0.69	0.85	0.67	<b>0.25</b>	0.69
Mistral-7B-Instruct	NA	0.38	0.53	0.23	0.44	0.24	NA	0.37
Llama-2-70b-chat	NA	0.48	0.39	0.07	0.61	0.24	0.08	0.31
PaLM2	NA	<b>0.76</b>	<b>0.96</b>	0.70	<b>0.87</b>	0.39	0.07	0.62
Llama-2-7b	MULTIALPACA	0.35	0.58	0.64	0.28	0.41	0.10	0.39
	Bactrian-X-22	0.35	0.58	0.63	0.28	0.44	0.08	0.39
	alpaca	0.35	0.58	0.63	0.28	0.35	0.07	0.38
Mistral-7b	MULTIALPACA	0.53	0.59	<b>0.79</b>	0.43	<b>0.70</b>	0.14	0.53
	Bactrian-X-22	0.52	0.59	<b>0.79</b>	0.42	<b>0.70</b>	0.14	0.53
	alpaca	0.53	0.59	0.78	0.45	<b>0.70</b>	0.10	0.52

Table 1: Detailed Task Wise Performance Comparison between GPT-4, PaLM-2, LLaMA-70B-chat, Mistral-7B-Instruct and finetuned models with best rank quantisation. Baseline numbers are referred from Ahuja et al. (2023b).

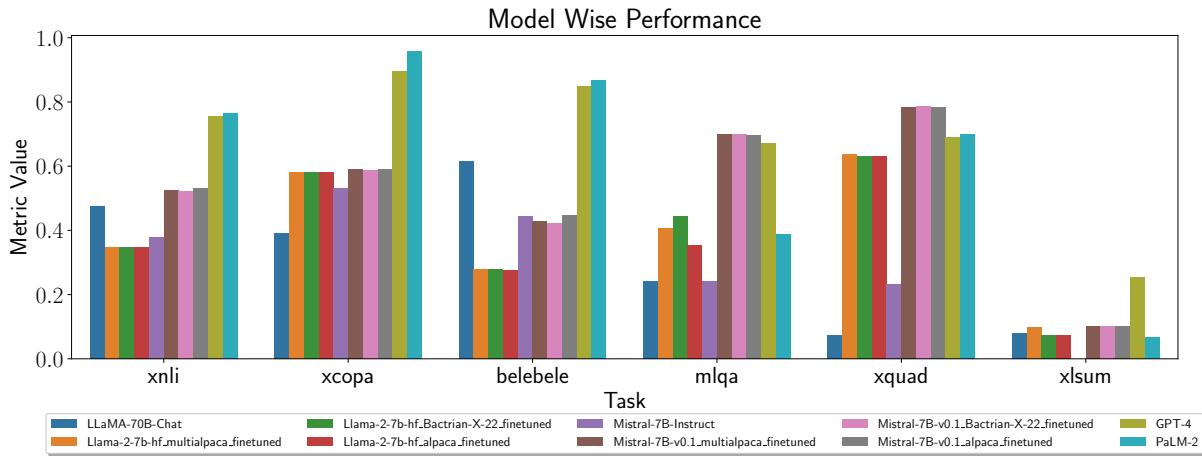


Figure 6: Task Wise Performance Comparison of Llama-2-70B-Chat, GPT-4, PaLM-2, Mistral-7B-Chat and our finetuned models averaged across languages.

model	winrate			
GPT-4	<b>93.78</b>			
PaLM2	79.66			
LLaMA-70B-Chat	22.36			
Mistral-7B-Instruct	35.12			
model	dataset	rank	quantisation	winrate
Llama-2-7B	Alpaca	128	16	13.28
	Bactrian-X-22	64	16	<b>13.73</b>
	MULTIALPACA	128	16	<b>13.73</b>
Mistral-7B	Alpaca	64	8	<b>24.47</b>
	Bactrian-X-22	16	8	22.07
	MULTIALPACA	32	8	22.45

Table 2: Best AlpacaEval Scores for each model, dataset, rank and quantisation configuration and GPT-4, PaLM-2, LLaMA-70B-Chat and Mistral-7B-Instruct baselines.

we can see that finetuning is usually better or at par with LLaMA-70B-Chat and Mistral-7B-Instruct which are not finetuned on multilingual data. We also compare the performance of finetuned models with LLMs like GPT-4 and PaLM2 and English instruction tuned versions of these models provided in (Ahuja et al., 2023b). Table 1 shows the best

model score averaged across languages per task. We observe that MISTRAL-7B beats GPT-4 and fairs as the best model on XQUAD and MLQA. While in XNLI, XCOPA and Belebele, it bridges the gap between GPT-4 and PaLM-2 by 20% on an average.

While we see some gap being bridged (20% on average) on classification tasks (XNLI, XCOPA and Belebele) using PEFT, it tends to beat larger models on question answering tasks (MLQA and XQUAD) like GPT-4. While on XLSUM there is no significant difference in performance after PEFT.

Interestingly, we observe that finetuning LLAMA-2-7B and MISTRAL-7B on ALPACA leads to comparable results with finetuning on MULTIALPACA and BACTRIAN-X-22. This can be due to the parameter efficient nature of the fine-tuning which prevents catastrophic forgetting and helps the model learn the instruction following ability from the English instruction data. Second

475 reason can be the difference in the token fertility  
476 of LLAMA-2-7B and MISTRAL-7B as shown  
477 by Ahuja et al. (2023b). We can deduce that  
478 MISTRAL-7B having higher token fertility and be-  
479 ing a better base multilingual model can benefit  
480 greatly from English instruction dataset and show  
481 excellent cross-lingual transfer. While LLAMA-2-  
482 7B having a lower token fertility does not benefit  
483 greatly from even multilingual instruction finetun-  
484 ing, as most multilingual LLaMAs resort to vocabu-  
485 lary expansion during pre-finetuning phase (Zhao  
486 et al., 2024). We illustrate the language-wise analy-  
487 sis of our model results in Figure 5, 15, 16, 17, 18,  
488 20, 21, 22, 22 and 23.

489 While we compare our multilingual finetuned  
490 models with models finetuned on English, we  
491 should note that we do not have complete informa-  
492 tion about instruction datasets used for Llama-70B-  
493 chat and Mistral-7B-Instruct. Hence, there may be  
494 chances of data contamination for some datasets in  
495 these models (Ahuja et al., 2023a) or the presence  
496 of multilingual instruction data in them.

## 497 5 Conclusion

498 In this paper we perform an extensive analysis of  
499 how rank, quantisation, finetuning dataset and base  
500 LLM effects the performance of the finetuned mod-  
501 els on 6 multilingual tasks and AlpacaEval when  
502 finetuned in a parameter efficient manner.

- 503 • Crosslingual transfer DOES happen even in  
504 parameter efficient finetuning.
- 505 • ALPACA (English-only instruction finetuning  
506 dataset) is comparable to MULTIALPACAand  
507 Bactrian-X-22 in multilingual downstream  
508 task performance. We hypothesize that this  
509 is due to Mistral being a superior model due  
510 to its better tokenizer, and that PEFT prevents  
511 catastrophic forgetting compared to full fine-  
512 tuning.
- 513 • Having more languages in the finetuning  
514 datasets does not necessarily mean signifi-  
515 cantly better multilingual performance (Sec-  
516 tion 4) if the dataset sizes are comparable.
- 517 • Quality and abilities of the base model far out-  
518 weigh the dataset or training method for pa-  
519 rameter efficient multilingual instruction fine-  
520 tuning.
- 521 • Higher capacity adapters (i.e. higher ranks or  
522 better quantisations) are better at maintaining

523 English performance along with multilingual  
524 downstream task performance.

525 We also beat GPT-4 on question-answering tasks  
526 (MLQA and XQUAD) using just multilingual  
527 PEFT on MISTRAL-7B showing that multilingual  
528 finetuning of 7B parameter LLMs is a promising  
529 direction for the future to bridge the gap of perfor-  
530 mance on multilingual downstream tasks.

## 531 6 Future Work

532 **More PEFT Techniques** This study explores the  
533 effects of PEFT using LoRA, while newer tech-  
534 niques by Ansell et al. (2024) can also be promising  
535 to study the parameter efficient techniques for mul-  
536 tilingual instruction tuning for these LLMs. We can  
537 also work towards building better PEFT techniques  
538 for specifically multilingual settings or crosslingual  
539 transfer for LLMs like MAD-X for encoder-like  
540 models (Pfeiffer et al., 2020b).

541 **Better Use of Adapters** With LLMs becoming  
542 more popular in the NLP research, it is prohibitive  
543 to re-pretrain models on newer languages. Hence,  
544 modular techniques like Pfeiffer et al. (2020a) were  
545 introduced to finetune these models in a more pa-  
546 rameter efficient manner. Furthermore, the works  
547 of Chen et al. (2024) introduced a novel technique  
548 of adapting language models to newer languages  
549 without further pretraining. However, such tech-  
550 niques can be compute intensive in LLMs and a  
551 direction of future can be pursued where similar  
552 approach can be taken in a parameter efficient man-  
553 ner.

554 **Better Multilingual Instruction Datasets**  
555 While the datasets used in this study are derived  
556 from ALPACA, newer instruction datasets using  
557 Chain Of Thought prompting (Mukherjee et al.,  
558 2023; Mitra et al., 2023) leads to better reasoning  
559 capabilities on smaller LLMs (7 billion param-  
560 eters), which can be explored as future work.  
561 More controlled crowd-sourcing efforts like Singh  
562 et al. (2024) can also lead to better multilingual  
563 instruction datasets.

## 564 7 Limitations

565 Our evaluation is performed using standard bench-  
566 marks, which has known limitations. Datasets used  
567 to create benchmarks may have been seen by mod-  
568 els during pretraining or finetuning, and due to lack  
569 of transparency about the datasets used for training

we cannot rule out test data contamination. Second, we use synthetic datasets that are created by prompting LLMs to finetune our models, this can lead to bias, which is also a limitation of the work. Finally, we compare the results obtained by our models to results from the MEGAVERSE benchmarking study while comparing the differences between finetuned models and models that are not finetuned for multilingual performance, which may have some differences in prompting and setup.

## References

- |     |   |     |
|-----|---|-----|
| 570 | Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. <i>Language models are few-shot learners</i> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc. | 623 |
| 571 |   | 624 |
| 572 |   | 625 |
| 573 |   | 626 |
| 574 |   | 627 |
| 575 |   | 628 |
| 576 |   | 629 |
| 577 |   | 630 |
| 578 |   | 631 |
| 579 |   | 632 |
| 580 |   | 633 |
| 581 | Jiaao Chen, Aston Zhang, Xingjian Shi, Mu Li, Alex Smola, and Diyi Yang. 2023. <i>Parameter-efficient fine-tuning design spaces</i> .   | 634 |
| 582 |   | 635 |
| 583 |   | 636 |
| 584 | Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Ifeoluwa Adelani, Pontus Stenetorp, Sebastian Riedel, and Mikel Artetxe. 2024. <i>Improving language plasticity via pretraining with active forgetting</i> .  | 637 |
| 585 |   | 638 |
| 586 |   | 639 |
| 587 |   | 640 |
| 588 |   | 641 |
| 589 |   | 642 |
| 590 |   | 643 |
| 591 |   | 644 |
| 592 | Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichiro Nakano, Christopher Hesse, and John Schulman. 2021. <i>Training verifiers to solve math word problems</i> .   | 645 |
| 593 |   | 646 |
| 594 |   | 647 |
| 595 |   | 648 |
| 596 | Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. <i>XNLI: Evaluating cross-lingual sentence representations</i> . In <i>Proceedings of EMNLP 2018</i> , pages 2475–2485.   | 649 |
| 597 |   | 650 |
| 598 |   | 651 |
| 599 |   | 652 |
| 600 | Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. <i>Free dolly: Introducing the world’s first truly open instruction-tuned llm</i> .  | 653 |
| 601 |   | 654 |
| 602 |   | 655 |
| 603 | Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 2022. <i>8-bit optimizers via block-wise quantization</i> .   | 656 |
| 604 |   | 661 |
| 605 |   | 662 |
| 606 | Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. <i>QLoRA: Efficient finetuning of quantized LLMs</i> . In <i>Neural Information Processing Systems (NeurIPS)</i> .  | 663 |
| 607 |   | 664 |
| 608 |   | 665 |
| 609 |   | 666 |
| 610 |   | 667 |
| 611 | Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. 2021. <i>A framework for few-shot language model evaluation. Version v0. 0.1. Sept.</i>   | 668 |
| 612 |   | 669 |
| 613 |   | 670 |
| 614 |   | 671 |
| 615 |   | 672 |
| 616 | Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. <i>XL-Sum: Large-scale multilingual abstractive summarization for 44 languages</i> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> ,   | 673 |
| 617 |   | 674 |
| 618 |   | 675 |
| 619 |   | 676 |
| 620 |   | 677 |
| 621 |   | 678 |
| 622 |   |     |
| 623 |   |     |
| 624 |   |     |
| 625 |   |     |
| 626 |   |     |
| 627 |   |     |
| 628 |   |     |
| 629 |   |     |
| 630 |   |     |
| 631 |   |     |
| 632 |   |     |
| 633 |   |     |
| 634 |   |     |
| 635 |   |     |
| 636 |   |     |
| 637 |   |     |
| 638 |   |     |
| 639 |   |     |
| 640 |   |     |
| 641 |   |     |
| 642 |   |     |
| 643 |   |     |
| 644 |   |     |
| 645 |   |     |
| 646 |   |     |
| 647 |   |     |
| 648 |   |     |
| 649 |   |     |
| 650 |   |     |
| 651 |   |     |
| 652 |   |     |
| 653 |   |     |
| 654 |   |     |
| 655 |   |     |
| 656 |   |     |
| 657 |   |     |
| 658 |   |     |
| 659 |   |     |
| 660 |   |     |
| 661 |   |     |
| 662 |   |     |
| 663 |   |     |
| 664 |   |     |
| 665 |   |     |
| 666 |   |     |
| 667 |   |     |
| 668 |   |     |
| 669 |   |     |
| 670 |   |     |
| 671 |   |     |
| 672 |   |     |
| 673 |   |     |
| 674 |   |     |
| 675 |   |     |
| 676 |   |     |
| 677 |   |     |
| 678 |   |     |

679	pages 4693–4703, Online. Association for Computational Linguistics.	732
680		733
681	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	734
682	Bruna Morrone, Quentin De Laroussilhe, Andrea	735
683	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	736
684	<b>Parameter-efficient transfer learning for NLP.</b> In	737
685	<i>Proceedings of the 36th International Conference on Machine Learning</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 2790–2799.	738
686	PMLR.	
687		
688		
689	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-	739
690	Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu	740
691	<b>LoRA: Low-rank adaptation of large</b>	741
692	<b>language models.</b> In <i>International Conference on Learning Representations</i> .	742
693		
694	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	744
695	sch, Chris Bamford, Devendra Singh Chaplot, Diego	745
696	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	746
697	laume Lample, Lucile Saulnier, Lélio Renard Lavaud,	747
698	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	748
699	Thibaut Lavril, Thomas Wang, Timothée Lacroix,	749
700	and William El Sayed. 2023. <b>Mistral 7b.</b>	750
701		751
702	Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian	752
703	Riedel, and Holger Schwenk. 2020. <b>MIQA: Evalu-</b>	753
704	<b>uating cross-lingual extractive question answering.</b>	754
705	In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7315–	755
706	7330.	756
707		757
708	Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji,	758
709	and Timothy Baldwin. 2023a. <b>Bactrian-x : A multi-</b>	759
710	<b>lingual replicable instruction-following model with low-rank adaptation.</b>	
711		
712	Xiang Lisa Li and Percy Liang. 2021. <b>Prefix-tuning: Optimizing continuous prompts for generation.</b>	
713		
714	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori,	760
715	Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and	761
716	Tatsunori B. Hashimoto. 2023b. <b>AlpacaEval: An Automatic Evaluator of Instruction-following Models.</b>	762
717		763
718	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mo-	764
719	hta, Tenghao Huang, Mohit Bansal, and Colin Raffel.	765
720	2022a. <b>Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning.</b>	766
721		
722	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengx-	767
723	iao Du, Zhilin Yang, and Jie Tang. 2022b. <b>P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks.</b> In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 61–68,	768
724	Dublin, Ireland. Association for Computational Lin-	769
725	guistics.	770
726		771
727		
728		
729	Ilya Loshchilov and Frank Hutter. 2019. <b>Decoupled weight decay regularization.</b> In <i>International Conference on Learning Representations</i> .	781
730		782
731		783
		784
		785
	Arindam Mitra, Luciano Del Corro, Shweta Mahajan,	786
	Andres Codas, Clarisse Simoes, Sahaj Agarwal, Xuxi	787
	Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Ag-	788
	garwal, Hamid Palangi, Guoqing Zheng, Corby Ros-	789
	set, Hamed Khanpour, and Ahmed Awadallah. 2023.	
	<b>Orca 2: Teaching small language models how to reason.</b>	
	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawa-	
	har, Sahaj Agarwal, Hamid Palangi, and Ahmed	
	Awadallah. 2023. <b>Orca: Progressive learning from</b>	
	<b>complex explanation traces of gpt-4.</b>	
	OpenAI. 2023. <b>Gpt4 technical report.</b>	743
	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé,	744
	Kyunghyun Cho, and Iryna Gurevych. 2021. <b>AdapterFusion: Non-destructive task composition for transfer learning.</b> In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 487–503, Online. Association for Computational Lin-	745
	guistics.	746
		747
	Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya	748
	Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun	749
	Cho, and Iryna Gurevych. 2020a. <b>AdapterHub: A framework for adapting transformers.</b> In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 46–54, Online. Association for Computational Linguistics.	750
		751
	Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Se-	752
	bastian Ruder. 2020b. <b>MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer.</b> In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7654–7673, Online. Association for Computational Linguistics.	753
		754
		755
		756
		757
		758
		759
	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska,	760
	Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020.	761
	<b>Xcopa: A multilingual dataset for causal common-</b>	762
	<b>sense reasoning.</b> In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2362–2376.	763
		764
		765
		766
	Krithika Ramesh, Arnav Chavan, Shrey Pandit, and	767
	Sunayana Sitaram. 2023. <b>A comparative study on the impact of model compression techniques on fairness in language models.</b> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15762–	774
	15782, Toronto, Canada. Association for Computational Linguistics.	775
		776
		777
		778
		779
		780
	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. <b>Choice of plausible alternatives: An evaluation of commonsense causal reasoning.</b> In <i>AAAI spring symposium: logical formalizations of commonsense reasoning</i> , pages 90–95.	781
		782
		783
		784
		785
	Uri Shaham, Jonathan Herzig, Roei Aharoni, Idan	786
	Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. <b>Mul-</b>	787
	<b>tilingual instruction tuning with just a pinch of multi-</b>	788
	<b>linguality.</b>	789

790	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. <a href="#">Language models are multilingual chain-of-thought reasoners</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	850
791		851
792		852
793		853
794		854
795		855
796		
797	Shivalika Singh, Freddie Vargas, Daniel Dsouza, Börje F. Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Souza Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergün, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Minh Chien, Sebastian Ruder, Surya Guthikonda, Emad A. Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. <a href="#">Aya dataset: An open-access collection for multilingual instruction tuning</a> .	856
798		857
799		858
800		859
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .	860
812		861
813		862
814		
815		
816		
817	NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. <a href="#">No language left behind: Scaling human-centered machine translation</a> .	866
818		867
819		868
820		869
821		870
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	871
833		872
834		
835		
836		
837		
838	Tu Vu, Aditya Barua, Brian Lester, Daniel Cer, Mohit Iyyer, and Noah Constant. 2022. <a href="#">Overcoming catastrophic forgetting in zero-shot cross-lingual generation</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9279–9300, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	880
839		881
840		882
841		
842		
843		
844		
845	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. <a href="#">Chain of thought prompting elicits reasoning in large language models</a> . In <i>Advances in Neural Information Processing Systems</i> .	883
846		884
847		885
848		886
849		
850	Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li, Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong Yang, Fei Huang, and Jun Xie. 2023. <a href="#">Polym: An open source polyglot large language model</a> .	850
851		851
852		852
853		853
854		854
855		855
856		
857		
858		
859		
860		
861		
862		
863		
864		
865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878		
879		

## A Hyperparameters for Finetuning and Training Setup

Our code for finetuning is based on the open source `axolotl`<sup>1</sup> framework. We plan to release our configuration files for better reproducibility. Each finetuning experiment took ~16-24 hours to complete on a single NVIDIA A100 GPU with 80 GB RAM. Exact hyperparameters for finetuning are mentioned below:

Hyperparameter	Value
Learning rate	$1 \times 10^{-6}$
Epochs	5
Global batch size	16
Scheduler	Cosine
Warmup	Linear
Warmup steps	10
Optimizer	AdamW (Loshchilov and Hutter, 2019)
Weight decay	0

Table 3: Hyperparameters for finetuning.

## B Evaluation Dataset Details

The detailed description of the datasets that we use for evaluation are as follows:

**XNLI:** The XNLI (Cross-lingual Natural Language Inference) dataset (Conneau et al., 2018) is an extension of the Multi-Genre NLI (MultiNLI) corpus to 15 languages. The dataset was created

<sup>1</sup><https://github.com/OpenAccess-AI-Collective/axolotl>

887 by manually translating the validation and test sets  
888 of MultiNLI into each of those 15 languages. The  
889 English training set was machine translated for all  
890 languages. The dataset is composed of 122k train,  
891 2490 validation, and 5010 test examples. XNLI  
892 provides a robust platform for evaluating cross-  
893 lingual sentence understanding methods. We eval-  
894 uated our models on the test split with 4 in-context  
895 examples sampled from the validation split. We  
896 report our results in Table 28, 29, 30, 31, 32 and  
897 33.

898 **XCOPA:** The XCOPA (Cross-lingual Choice of  
899 Plausible Alternatives) dataset (Ponti et al., 2020)  
900 is a benchmark for evaluating the ability of ma-  
901 chine learning models to transfer commonsense  
902 reasoning across languages. It is a translation and  
903 re-annotation of the English COPA dataset (Roem-  
904 mele et al., 2011) and covers 11 languages from 11  
905 families and several areas around the globe. The  
906 dataset is challenging as it requires both the com-  
907 mand of world knowledge and the ability to gener-  
908 alize to new languages. We evaluated our models  
909 on Estonian, Thai, Italian, Indonesian, Vietnamese  
910 and Southern Quechua. We evaluated our models  
911 in the 4-shot setting similar to XNLI. We report our  
912 results in Table 22, 23, 24, 25, 26 and 27.

913 **Belebele:** Belebele (Bandarkar et al., 2023) is  
914 a multiple choice machine reading compre-  
915 hension (MRC) dataset parallel across 122 languages.  
916 Each question is linked to a short passage from the  
917 FLORES-200 dataset (Team et al., 2022). The hu-  
918 man annotation procedure was carefully curated to  
919 create questions that discriminate between different  
920 levels of language comprehension. We evaluated  
921 our models in the zero-shot setting and report re-  
922 sults in Table 4, 5 6, 7, 8 and 9.

923 **MLQA:** MLQA (Lewis et al., 2020) is a mul-  
924 tilingual question answering dataset designed for  
925 cross lingual question answering. It contains 5K  
926 extractive question answering instances. It consists  
927 of 7 languages i.e. English, Arabic, Vietnamese,  
928 German, Spanish, Hindi and Simplified Chinese.  
929 The evaluation uses a 4-shot setting similar to that  
930 of XNLI. We report our results in Table 34, 35, 36,  
931 37, 38 and 39.

932 **XQuAD:** The XQuAD (Cross-lingual Question  
933 Answering Dataset) (Artetxe et al., 2020) is a  
934 benchmark dataset for evaluating cross-lingual  
935 question answering performance. It consists of

936 a subset of 240 paragraphs and 1190 question-  
937 answer pairs from the development set of SQuAD  
938 v1.1, along with their professional translations into  
939 ten languages: Spanish, German, Greek, Russian,  
940 Turkish, Arabic, Vietnamese, Thai, Chinese, and  
941 Hindi. As a result, the dataset is entirely parallel  
942 across 11 languages. This dataset provides a robust  
943 platform for developing and evaluating models on  
944 cross-lingual question answering tasks. For evalua-  
945 tion, we use a 4-shot setting similar to MLQA. We  
946 report our results in table 34, 35, 36, 37, 38 and 39.

947 **XLSUM:** XLSUM (Hasan et al., 2021) is a com-  
948 prehensive and diverse dataset for abstractive sum-  
949 marization comprising 1 million human annotated  
950 article-summary pairs from BBC. The dataset cov-  
951 ers 44 languages ranging from low to high-resource,  
952 for many of which no public dataset is currently  
953 available. We evaluate our models on a subset of  
954 7 languages, namely, Arabic, Chinese-Simplified,  
955 English, Hindi, French, Japanese and Spanish in a  
956 zero-shot setting. We present our results in Table  
957 40, 42, 41, 43, 44 and 45.

958 **AlpacaEval:** AlpacaEval (Li et al., 2023b) is an  
959 LLM based automatic evaluator for instruction fol-  
960 lowing models. It consists of around 800 instruc-  
961 tions and corresponding responses obtained from  
962 (text-davinci-003) GPT3. The benchmark com-  
963 pares responses from GPT3 (or any other “oracle”  
964 model) with target (finetuned) model using another  
965 LLM (typically GPT4) as an evaluator. The eval-  
966 uator LLM decides which response is better and  
967 overall win rate (higher the better) is computed for  
968 the target model. For our evaluation, we use the  
969 text-davinci-003 responses from the dataset as our  
970 oracle/gold responses and use GPT4 (gpt-4-32k) as  
971 our evaluator. We report our results in Table 10, 11,  
972 12, 13, 14 and 15.

## C Evaluation Prompts

973 For XNLI, XCOPA, Belebele, MLQA, XQUAD,  
974 XLSUM we use the standard ALPACA system  
975 prompt "**Below is an instruction that describes a**  
976 **task, paired with an input that provides further**  
977 **context. Write a response that appropriately**  
978 **completes the request.**".

```

### Instruction:
The task is to solve Natural Language Inference (NLI) problems. NLI is the task of determining the inference relation between two (short, ordered) texts: entailment, contradiction, or neutral. Answer as concisely as possible in the same format as the examples below:

{{premise}}
```

Question: {{hypothesis}} True, False, or Neither?

```

### Response
```

Figure 7: XNLI Prompt

```

### Instruction:
The task is to perform reading comprehension task. Given the following passage, query, and answer choices, output the letter corresponding to the correct answer.

Passage: {{florespassage}}
Query : {{question}}
Choices :
A : {{mcansWER1}}
B : {{mcansWER2}}
C : {{mcansWER3}}
D : {{mcansWER4}}
```

```

### Response :
```

Figure 9: Belebele Prompt

```

### Instruction:
The task is to perform open-domain commonsense causal reasoning. You will be provided a premise and two alternatives, where the task is to select the alternative that more plausibly has a causal relation with the premise. Answer as concisely as possible in the same format as the examples below:

Given this premise:
{{premise}}
```

What's the best option?

```

-choice1 : {{choice1}}
-choice2 : {{choice2}}
```

We are looking for % if question == "cause" %} a cause % else %} an effect % endif %}

```

### Response:
```

Figure 8: XCOPA Prompt

```

### Instruction:
The task is to solve reading comprehension problems. You will be provided questions on a set of passages and you will need to provide the answer as it appears in the passage. The answer should be in the same language as the question and the passage.

Context:{{context}}
Question:{{question}}
```

Referring to the passage above, the correct answer to the given question is

```

### Response:
```

Figure 10: MLQA Prompt

```

### Instruction:
The task is to solve reading comprehension problems. You will be provided questions on a set of passages and you will need to provide the answer as it appears in the passage. The answer should be in the same language as the question and the passage.

Context:{{context}}
Question:{{question}}
```

Referring to the passage above, the correct answer to the given question is

```

### Response:
```

Figure 11: XQUAD Prompt

- 980 **C.1 XNLI**
- 981 **C.2 XCOPA**
- 982 **C.3 Belebele**
- 983 **C.4 MLQA**
- 984 **C.5 XQUAD**
- 985 **C.6 XLSUM**
- 986 **C.7 AlpacaEval**
- 987 **D Further of Analysis of Results**
- 988 **D.1 Analysis of Rank and Quantisation**
- 989 **D.2 Tasks Wise and Language Performance Plots**

```
### Instruction:  
The task is to summarize any given  
article. You should summarize all  
important information concisely  
in the same language in which you  
have been provided the document.  
Following the examples provided  
below:  
{ {text} }  
### Response:
```

Figure 12: XLSUM Prompt

Below is an instruction that  
describes a task, paired with an  
input that provides further context.  
Write a response that appropriately  
completes the request.

Figure 13: Alpaca Prompt

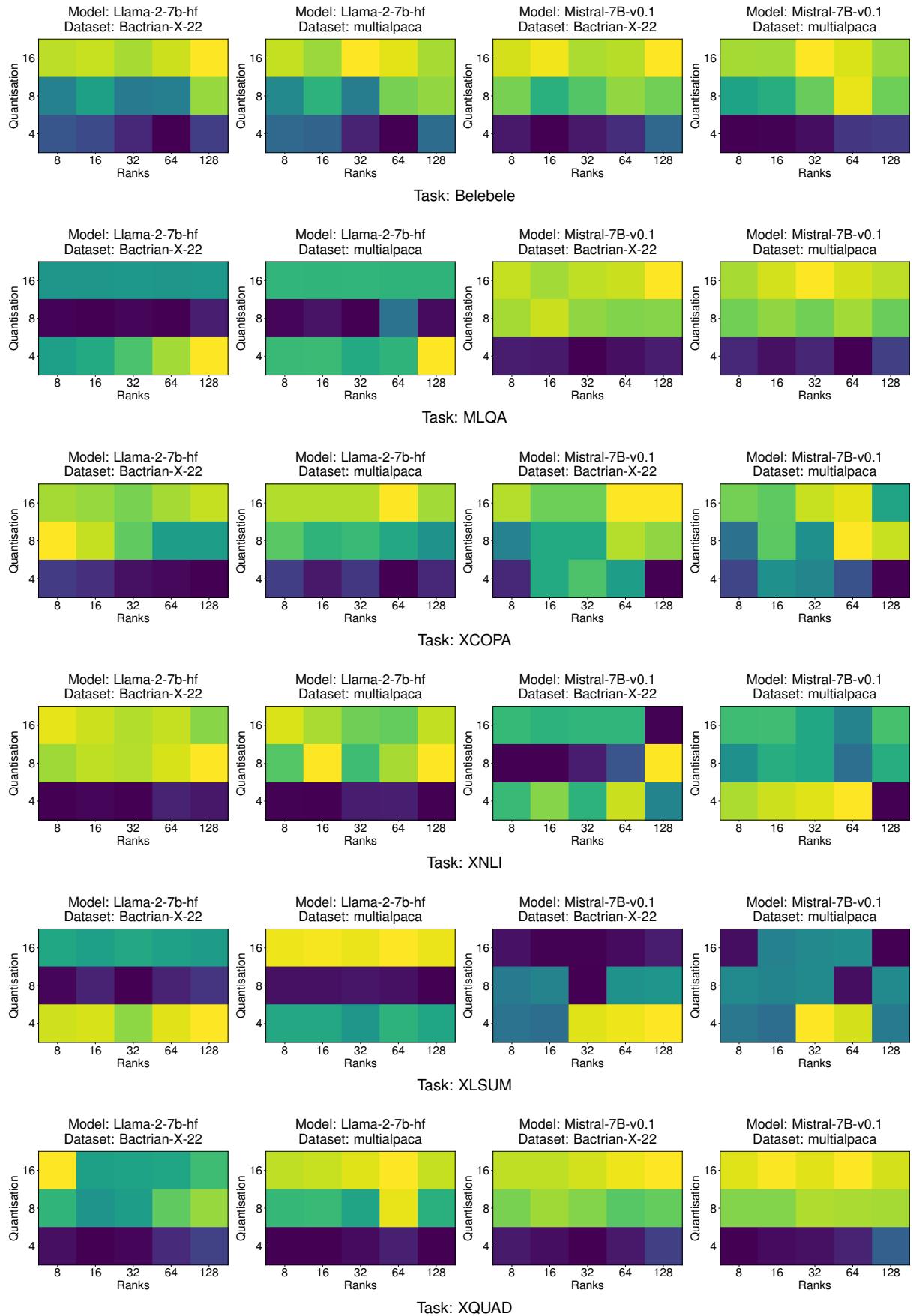


Figure 14: Task-wise performance of MISTRAL-7B and LLAMA-2-7B fine-tuned on BACTRIAN-X-22 and MULTIALPACA averaged across languages on all rank-quantisation configurations.

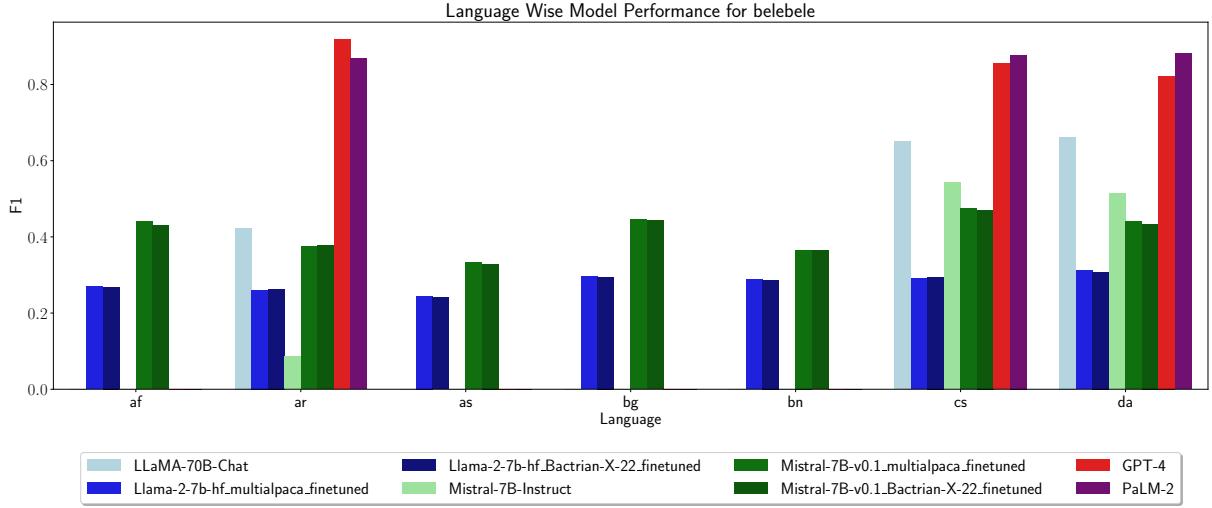


Figure 15: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Afrikaans, Arabic, Assamese, Bulgarian, Bengali, Czech and Danish for Belebele (Bandarkar et al., 2023).

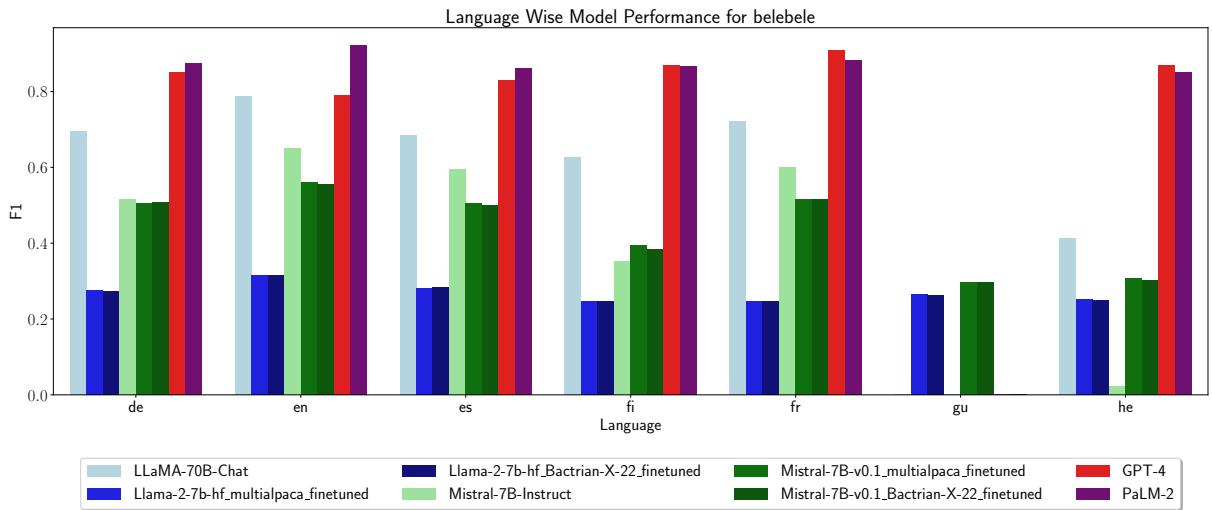


Figure 16: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on German, English, Spanish, Finnish, French, Gujarati and Hebrew for Belebele (Bandarkar et al., 2023).

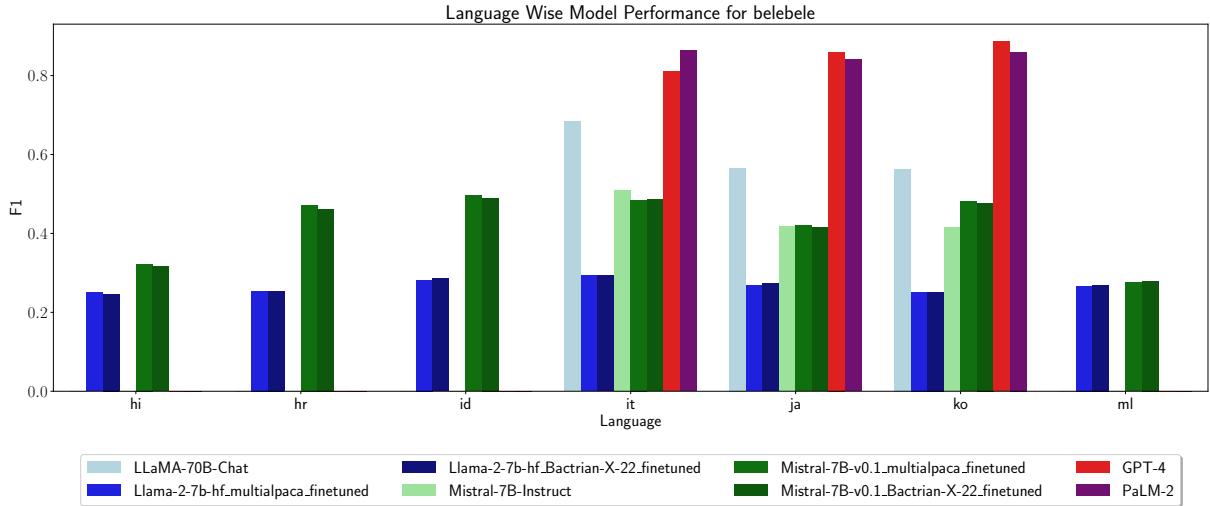


Figure 17: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Hindi, Croatian, Indonesian, Italian, Japanese, Korean and Malayalam for Belebele (Bandarkar et al., 2023).

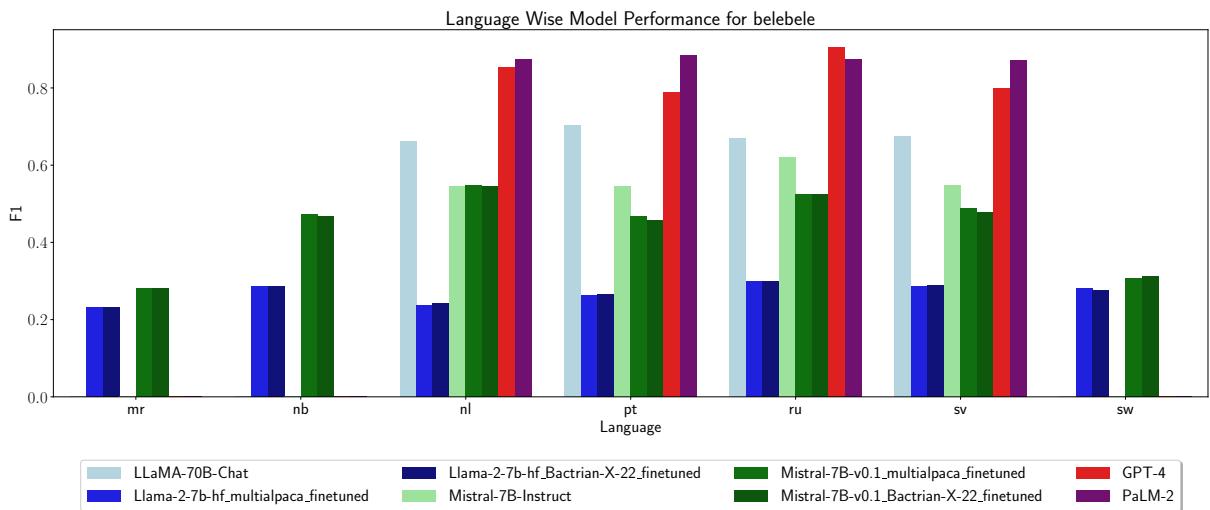


Figure 18: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Marathi, Norwegian, Dutch, Portuguese, Russian, Swedish and Swahili for Belebele (Bandarkar et al., 2023).

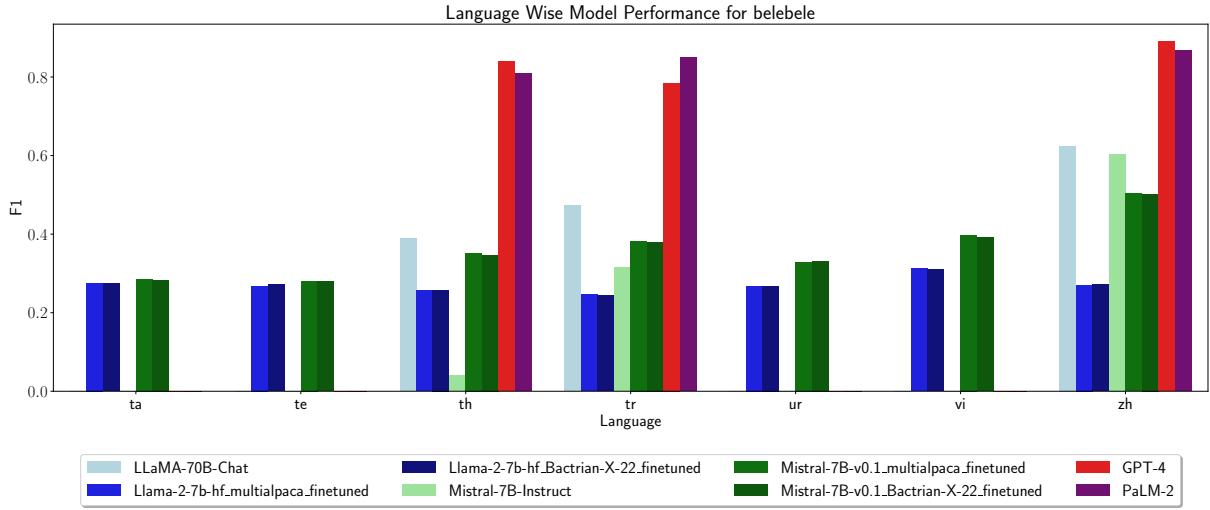


Figure 19: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Tamil, Telugu, Thai, Turkish, Urdu, Vietnamese and Chinese-Simplified for Belebele (Bandarkar et al., 2023).

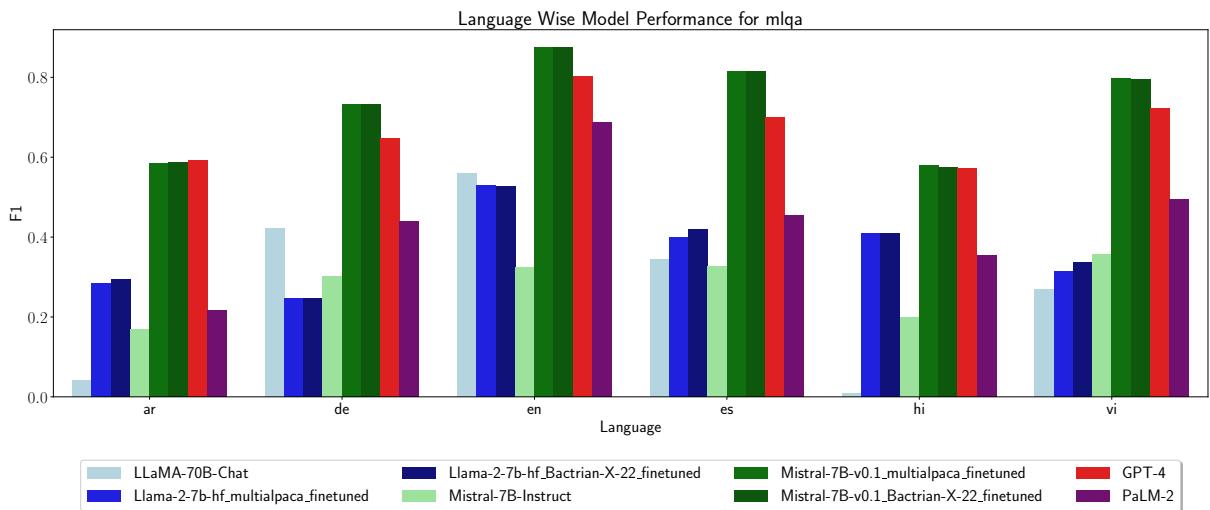


Figure 20: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Arabic, German, English, Spanish, Hindi and Vietnamese for MLQA (Lewis et al., 2020).

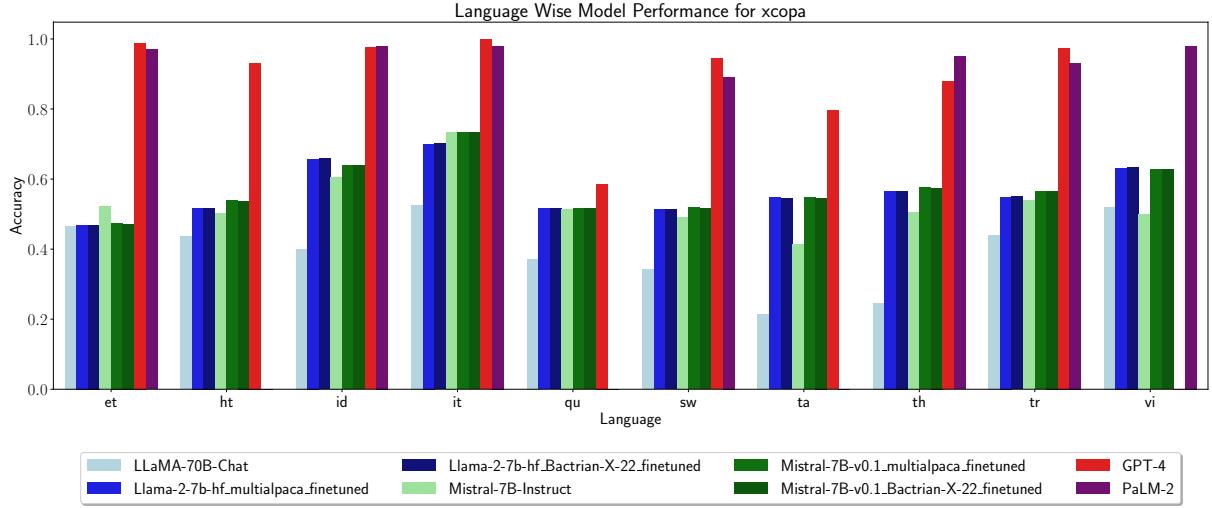


Figure 21: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Estonian, Haitian, Indonesian, Italian, Quechua, Swahili, Tamil, Thai, Turkish and Vietnamese for XCOPA (Ponti et al., 2020).

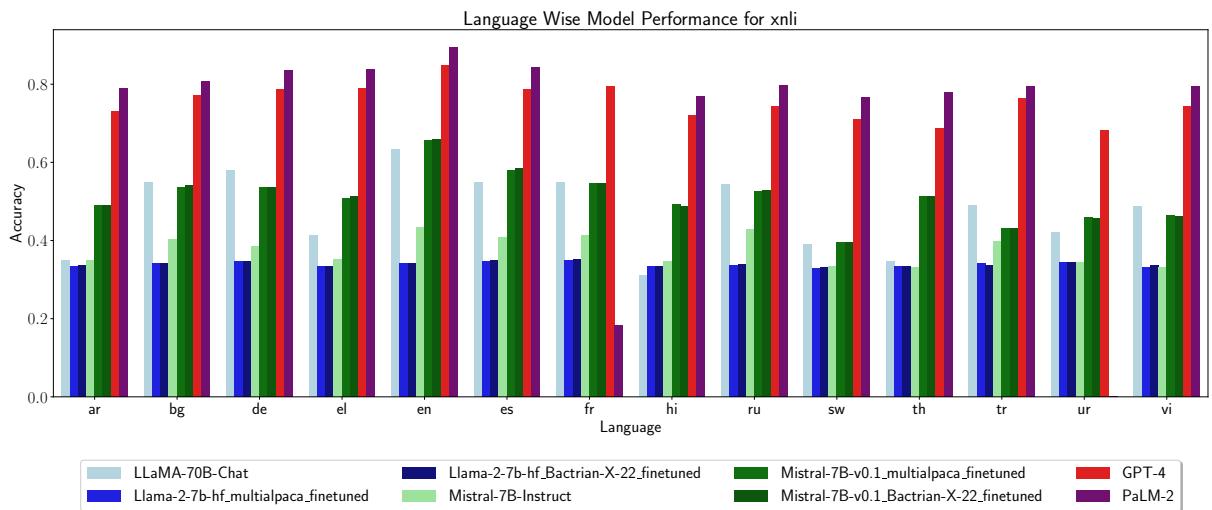


Figure 22: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Arabic, Bulgarian, German, Greek, English, Spanish, French, Hindi, Russian, Swahili, Thai, Turkish, Urdu and Vietnamese for XNLI (Conneau et al., 2018).

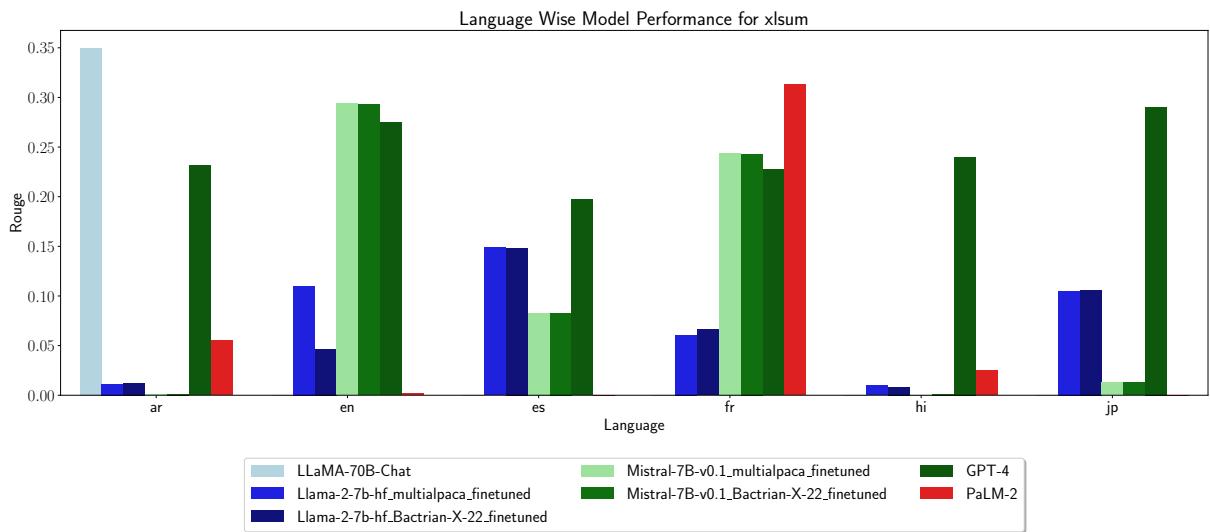


Figure 23: Detailed language-wise comparison of our fine-tuned MISTRAL-7B and LLAMA-2-7B models with other baselines (Ahuja et al., 2023b) on Arabic, English, Spanish, French, Hindi and Japanese for XLSUM(Hasan et al., 2021).

Model	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
LLaMA-70B-Chat	-	0.42	-	-	-	0.65	0.66	0.69	0.79	0.68	0.63	0.72	-	0.41	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61			
Mistral-7B-Instruct	-	0.09	-	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44			
GPT-3.5-Turbo	-	0.69	-	-	-	0.77	0.81	0.83	0.88	0.79	0.78	0.83	-	0.64	-	-	0.80	0.71	0.67	-	-	0.80	0.83	0.78	0.82	-	-	0.56	0.70	-	-	0.78	0.76			
GPT-4	-	<b>0.92</b>	-	-	-	0.85	0.82	0.85	0.79	0.83	<b>0.87</b>	<b>0.91</b>	-	<b>0.87</b>	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.78	-	-	<b>0.89</b>	0.85			
PALM2	-	0.87	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	0.87	0.88	-	0.85	-	-	<b>0.86</b>	0.84	0.86	-	-	<b>0.87</b>	<b>0.88</b>	0.87	<b>0.87</b>	-	-	0.81	<b>0.85</b>	-	-	0.87	<b>0.87</b>				
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
64 8	0.27	0.27	0.22	0.33	0.29	0.30	0.28	0.30	0.33	0.30	0.25	0.26	0.24	0.24	0.24	0.24	0.26	0.26	0.25	0.28	0.23	0.30	0.23	0.26	0.29	0.29	0.29	0.27	0.26	0.26	0.27	0.30	0.27	0.27		
64 16	<b>0.29</b>	0.25	<b>0.27</b>	0.31	0.29	0.32	0.29	0.27	0.32	0.27	0.24	0.25	<b>0.27</b>	<b>0.26</b>	<b>0.27</b>	<b>0.29</b>	<b>0.29</b>	<b>0.28</b>	<b>0.28</b>	<b>0.28</b>	<b>0.26</b>	<b>0.27</b>	<b>0.24</b>	<b>0.24</b>	<b>0.29</b>	<b>0.29</b>	<b>0.27</b>	<b>0.28</b>	<b>0.27</b>	<b>0.28</b>	<b>0.28</b>	<b>0.28</b>	<b>0.28</b>	<b>0.28</b>		
128 8	0.27	0.24	0.23	<b>0.34</b>	<b>0.30</b>	0.30	0.28	0.28	0.32	0.29	0.24	0.25	0.26	0.23	0.25	0.25	0.27	0.29	0.27	0.26	<b>0.31</b>	<b>0.31</b>	0.23	<b>0.31</b>	0.23	0.27	0.30	<b>0.30</b>	<b>0.29</b>	0.29	0.25	0.27	0.31	0.27	0.28	
128 16	<b>0.29</b>	0.24	<b>0.27</b>	0.31	0.29	0.32	0.29	0.27	0.32	0.27	0.24	0.25	<b>0.27</b>	<b>0.26</b>	<b>0.27</b>	<b>0.29</b>	<b>0.29</b>	<b>0.27</b>	<b>0.28</b>	<b>0.26</b>	<b>0.27</b>	<b>0.24</b>	<b>0.29</b>	<b>0.27</b>	<b>0.31</b>	<b>0.29</b>	<b>0.27</b>	<b>0.29</b>	<b>0.27</b>	<b>0.28</b>	<b>0.32</b>	<b>0.27</b>	<b>0.28</b>			

Table 4: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on Belebele (Bandarkar et al., 2023).

Model	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg				
LLaMA-70B-Chat	-	0.42	-	-	-	0.65	0.66	0.69	0.79	0.68	0.63	0.72	-	0.41	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61			
Mistral-7B-Instruct	-	0.09	-	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44			
GPT-3.5-Turbo	-	0.69	-	-	-	0.77	0.81	0.83	0.88	0.79	0.78	0.83	-	0.64	-	-	0.80	0.71	0.67	-	-	0.80	0.83	0.78	0.82	-	-	0.56	0.70	-	-	0.78	0.76			
GPT-4	-	<b>0.92</b>	-	-	-	0.85	0.82	0.85	0.79	0.83	<b>0.87</b>	<b>0.91</b>	-	<b>0.87</b>	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.78	-	-	<b>0.89</b>	0.85			
PALM2	-	0.87	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	0.87	0.88	-	0.85	-	-	<b>0.86</b>	0.84	0.86	-	-	<b>0.87</b>	<b>0.88</b>	0.87	<b>0.87</b>	-	-	0.81	<b>0.85</b>	-	-	0.87	<b>0.87</b>				
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
64 8	0.49	0.39	0.30	0.49	<b>0.41</b>	0.49	0.47	0.56	0.60	0.51	0.40	0.53	0.32	0.33	0.35	0.50	0.50	0.53	0.44	0.50	<b>0.31</b>	0.29	0.49	0.58	0.52	0.54	0.53	0.31	0.29	0.28	0.34	0.41	0.33	0.40	0.56	0.44
64 16	0.49	0.43	0.33	0.50	0.49	0.48	0.56	0.60	0.51	0.40	0.54	<b>0.34</b>	0.32	0.36	<b>0.51</b>	0.51	0.53	0.44	0.53	0.31	0.29	0.49	0.60	0.52	0.54	<b>0.55</b>	<b>0.33</b>	0.28	<b>0.30</b>	0.36	0.40	0.34	<b>0.40</b>	0.56	0.44	
128 8	<b>0.50</b>	0.42	0.30	0.48	0.37	0.51	0.47	0.56	0.59	0.54	0.40	0.53	0.30	0.35	0.34	0.49	<b>0.54</b>	0.53	0.44	0.51	0.31	0.27	<b>0.50</b>	0.61	0.51	0.52	0.56	0.39	0.30	0.29	0.35	0.41	0.34	0.40	0.58	0.44
128 16	0.46	0.43	0.31	<b>0.51</b>	0.39	0.51	0.48	0.56	0.61	0.54	0.43	0.53	0.31	0.33	<b>0.36</b>	0.50	<b>0.54</b>	0.54	0.45	0.53	0.29	<b>0.30</b>	0.49	0.61	0.54	0.55	0.55	0.31	0.29	0.28	0.38	0.41	<b>0.37</b>	0.39	0.55	0.45

Table 5: Detailed performance of various ALPACA finetuned MISTRAL-7B models on Belebele (Bandarkar et al., 2023).

Model	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
LLAMA-70B-Chat	-	0.42	-	-	0.65	0.66	0.69	0.79	0.68	0.63	0.72	-	0.41	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61				
Mistral-7B-Instruct	-	0.09	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44				
GPT-3.5-Turbo	-	0.69	-	-	0.77	0.81	0.83	0.88	0.79	0.83	0.87	0.91	-	0.64	-	-	0.80	0.71	0.67	-	-	0.80	0.83	0.78	0.82	-	-	0.78	0.76	-	-	0.89	0.85			
GPT-4	-	<b>0.92</b>	-	-	0.75	0.82	0.85	0.79	0.83	0.87	0.91	-	0.87	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.78	-	-	<b>0.84</b>	0.78				
PALM2	-	0.87	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	<b>0.87</b>	<b>0.91</b>	-	0.85	-	-	<b>0.86</b>	0.84	0.86	-	-	<b>0.87</b>	<b>0.88</b>	0.87	<b>0.87</b>	-	-	0.81	<b>0.85</b>	-	-	0.87	<b>0.87</b>				
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
8	4	0.26	0.29	0.24	0.27	0.29	0.27	0.37	0.29	0.29	0.25	0.23	<b>0.28</b>	0.26	0.24	0.23	0.28	0.29	0.26	0.22	0.26	0.23	0.28	0.24	0.26	0.27	0.26	0.23	0.25	0.26	0.29	0.30	0.27			
8	8	0.26	0.24	0.25	0.31	0.29	0.27	0.30	0.26	0.31	0.27	0.26	0.26	0.23	0.26	0.26	0.26	0.28	0.28	0.27	0.25	0.24	<b>0.31</b>	0.22	0.25	0.31	0.28	0.29	0.28	<b>0.29</b>	0.26	0.24	0.26	0.33	0.26	0.27
8	16	0.28	<b>0.24</b>	<b>0.26</b>	0.30	0.32	0.26	0.33	0.20	0.27	0.24	0.24	<b>0.28</b>	0.26	<b>0.27</b>	<b>0.26</b>	<b>0.29</b>	0.28	0.28	0.26	0.27	0.26	0.23	0.30	0.26	0.27	0.28	0.27	0.28	0.27	0.28	0.33	0.27			
16	4	0.26	0.29	0.24	0.27	0.29	0.27	0.37	0.29	0.29	0.25	0.23	<b>0.28</b>	0.27	0.24	0.22	0.28	0.29	0.26	0.26	0.22	0.28	0.24	0.26	0.28	0.26	0.27	0.27	0.26	0.25	0.25	0.26	0.29	0.30	0.27	
16	8	0.28	0.27	0.25	0.29	0.27	0.27	0.27	0.29	0.33	0.25	0.26	0.25	0.24	0.23	0.25	0.28	0.28	0.27	0.27	0.26	0.25	0.25	0.25	0.26	0.27	0.28	0.28	0.27	0.28	0.28	0.28	0.27			
16	16	0.28	0.24	0.26	0.30	0.30	0.31	0.30	0.27	0.27	0.23	0.20	0.26	0.27	0.23	0.26	0.27	0.26	0.28	0.28	0.27	0.23	0.20	0.25	0.27	0.23	0.26	0.28	0.27	0.28	0.28	0.27				
32	4	0.26	0.29	0.23	0.28	0.29	0.28	0.35	0.26	0.29	0.28	0.24	0.23	0.26	0.24	0.24	0.28	0.31	0.26	0.22	0.22	0.27	0.24	0.26	0.27	0.24	0.28	0.29	0.29	0.26	0.29	0.26				
32	8	0.28	0.24	0.24	0.32	0.26	0.29	0.27	0.27	0.32	0.29	0.26	0.26	0.27	0.23	0.26	0.26	0.27	0.24	0.27	0.27	0.23	0.24	0.28	0.27	0.24	0.28	0.27	0.24	0.28	0.27	0.27				
32	16	0.25	0.29	0.31	0.31	0.31	0.31	0.29	0.31	0.31	0.27	0.33	0.27	0.25	0.27	0.26	<b>0.28</b>	0.27	0.29	0.28	0.27	0.23	0.30	0.24	0.27	0.28	0.27	0.28	0.27	0.28	0.27	0.28				
64	4	0.24	0.29	0.22	0.29	0.28	0.28	0.35	0.26	0.24	0.23	0.23	0.26	0.25	0.24	0.23	0.28	0.32	0.30	0.27	0.26	0.23	0.23	0.26	0.25	0.24	0.29	0.27	0.26	0.27	0.26	0.27				
64	8	<b>0.29</b>	<b>0.25</b>	<b>0.26</b>	0.32	0.28	0.30	0.29	0.28	0.32	0.26	0.24	0.22	0.24	0.22	0.24	0.26	<b>0.28</b>	0.24	0.29	0.23	0.22	0.20	0.27	0.26	0.28	0.27	0.26	0.27	0.26	0.27	0.26				
64	16	0.28	0.29	0.24	0.24	0.30	0.29	0.26	0.29	0.30	0.26	0.24	0.24	0.24	0.24	0.24	0.28	0.30	0.28	0.26	0.26	0.25	0.27	0.27	0.26	0.27	0.27	0.26	0.27	0.26	0.27	0.26				
128	4	0.26	0.27	0.24	0.27	0.29	0.27	0.29	0.27	0.29	0.27	0.29	0.27	0.29	0.27	0.29	0.28	0.30	0.28	0.26	0.26	0.22	0.27	0.27	0.26	0.27	0.26	0.27	0.26	0.27	0.26	0.27				
128	8	0.28	0.24	<b>0.26</b>	0.30	0.30	0.31	0.28	0.27	0.32	0.27	0.23	0.25	0.26	0.26	0.27	0.26	0.28	0.29	0.28	0.27	0.23	<b>0.31</b>	0.25	0.27	0.27	0.25	0.28	0.27	0.26	0.28	0.27				
128	16	0.28	0.24	0.26	0.30	0.30	0.31	0.28	0.27	0.32	0.27	0.23	0.25	0.26	0.26	0.27	0.26	0.28	0.29	0.28	0.27	0.23	<b>0.31</b>	0.25	0.27	0.27	0.25	0.28	0.27	0.26	0.28	0.27				

Table 6: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on Belebele (Bandarkar et al., 2023).

Model	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
LLAMA-70B-Chat	-	0.42	0.34	0.31	0.40	0.33	0.44	0.43	0.48	0.52	0.50	0.36	0.53	0.30	0.30	<b>0.50</b>	0.48	0.44	0.39	0.45	0.27	0.46	0.55	0.47	0.47	0.44	0.29	0.29	0.27	0.34	0.41	0.34	0.38	0.47	0.40	
Mistral-7B-Instruct	-	0.09	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61				
GPT-3.5-Turbo	-	0.69	-	-	0.77	0.81	0.83	0.88	0.79	0.83	0.87	0.91	-	0.64	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44			
GPT-4	-	<b>0.92</b>	-	-	0.75	0.82	0.85	0.82	0.85	0.87	0.91	-	0.87	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.76	-	-	0.78	0.76				
PALM2	-	0.87	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	<b>0.87</b>	<b>0.88</b>	-	-	-	-	<b>0.86</b>	0.84	0.86	-	-	<b>0.87</b>	<b>0.88</b>	0.87	<b>0.87</b>	-	-	0.81	<b>0.85</b>	-	-	0.87	<b>0.87</b>			
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
8	4	0.44	0.34	0.46	0.38	0.50	0.44	0.43	0.48	0.52	0.50	0.36	0.53	0.30	0.30	<b>0.50</b>	0.48	0.44	0.39	0.45	0.27	0.46	0.55	0.47	0.47	0.44	0.29	0.29	0.27	0.34	0.41	0.34	0.38	0.47	0.40	
8	8	0.45	0.37	0.33	0.44	0.37	0.47	0.44	0.44	0.51	0.58	0.30	0.51	0.30	0.31	0.45	0.50	0.48	0.41	0.49	0.30	0.27	0.47	0.56	0.45	0.53	0.50	0.31	0.37	0.33	0.40	0.51	0.42			
8	16	0.45	0.39	0.34	0.43	0.40	0.45	0.45	0.47	0.51	0.58	0.35	0.52	0.27	0.29	0.46	0.50	0.49	0.43	0.50	0.28	0.28	0.48	0.55	0.45	0.54	0.52	0.33	0.39	0.47	0.50	0.42				
8	32	0.47	0.39	0.32	0.44	0.38	0.48	0.44	0.51	0.58	0.51	0.41	0.50	0.30	0.31	<b>0.55</b>	0.47	0.48	0.49	0.45	0.27	0.48	0.55	0.47	0.47	0.44	0.32	0.27	0.33	0.40	0.47	0.40				
16	4	0.45	0.34	0.31	0.41	0.34	0.45	0.44	0.47	0.52	0.50	0.36	0.53	0.30	0.30	0.46	0.50	0.49	0.44	0.40	0.27	0.48	0.55	0.45	0.45	0.42	0.33	0.27	0.33	0.40	0.47	0.40				
16	8	0.45	0.34	0.31	0.41	0.34	0.45	0.44	0.47	0.52	0.50	0.36	0.53	0.30	0.30	0.46	0.50	0.49	0.44	0.40	0.27	0.48	0.55	0.45	0.45	0.42	0.33	0.27	0.33	0.40	0.47	0.40				
16	16	0.47	0.37	0.33	0.43	0.34	0.47</td																													

Model	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
LLAMA-70B-Chat	-	0.42	-	-	0.65	0.66	0.69	0.79	0.68	0.63	0.72	-	0.41	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61				
Mistral-7B-Instruct	-	0.09	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44				
GPT-3.5-Turbo	-	0.69	-	-	0.85	0.81	0.83	0.88	0.79	0.83	0.87	<b>0.91</b>	-	0.64	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.78	-	-	0.78	0.76			
GPT-4	-	<b>0.92</b>	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	<b>0.87</b>	<b>0.91</b>	-	<b>0.87</b>	-	-	<b>0.86</b>	<b>0.84</b>	<b>0.86</b>	-	-	<b>0.87</b>	<b>0.88</b>	<b>0.87</b>	<b>0.87</b>	-	-	<b>0.81</b>	<b>0.85</b>	-	-	0.87	<b>0.87</b>			
PALM2	-	0.87	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	<b>0.87</b>	<b>0.91</b>	-	<b>0.85</b>	-	-	<b>0.86</b>	<b>0.84</b>	<b>0.86</b>	-	-	<b>0.87</b>	<b>0.88</b>	<b>0.87</b>	<b>0.87</b>	-	-	<b>0.81</b>	<b>0.85</b>	-	-	0.87	<b>0.87</b>			
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
8	4	0.26	0.29	0.24	0.27	0.29	0.27	0.37	0.29	0.29	0.25	0.23	<b>0.28</b>	0.27	0.24	0.22	0.28	0.29	0.26	0.22	0.28	0.24	0.26	0.28	0.27	0.29	0.27	0.23	0.25	0.26	0.29	0.31	0.27			
8	8	0.26	0.25	0.25	0.30	0.29	0.29	0.26	0.27	0.31	0.27	0.25	0.24	0.25	0.23	0.26	0.27	0.29	0.27	0.26	<b>0.29</b>	0.23	0.28	0.21	0.27	0.29	0.30	<b>0.29</b>	0.28	0.27	0.25	0.26	0.31	0.28	0.27	
8	16	0.28	0.24	0.27	0.31	0.29	0.32	0.29	0.26	0.32	0.27	0.24	<b>0.28</b>	0.26	<b>0.27</b>	<b>0.29</b>	0.29	0.28	0.28	0.26	0.27	0.26	0.27	0.28	0.27	0.28	0.27	0.28	0.27	0.28	0.32	0.27	0.28			
16	4	0.26	0.28	0.24	0.27	0.29	0.27	0.36	0.29	0.29	0.25	0.23	0.23	0.27	0.27	0.25	0.23	0.28	0.29	0.26	0.22	0.26	0.23	0.27	0.24	0.26	0.27	0.28	0.26	0.27	0.30	0.27				
16	8	0.28	0.25	0.24	0.29	0.26	0.29	0.28	0.27	0.33	0.29	0.23	0.26	0.26	0.23	0.22	0.25	0.29	0.29	0.26	0.28	<b>0.26</b>	0.29	0.23	0.26	0.27	0.28	0.28	0.28	0.28	0.28	0.27	0.27			
16	16	0.28	0.24	0.27	0.31	0.29	0.32	0.29	0.27	0.32	0.27	0.24	0.25	0.27	0.26	0.27	0.26	0.27	0.29	0.27	0.26	0.27	0.27	0.28	0.27	0.29	0.27	0.28	0.28	0.27	0.28	0.27	0.27			
16	32	0.26	0.28	0.24	0.28	0.28	0.30	0.28	0.27	0.35	0.29	0.29	0.25	0.23	0.27	0.24	0.23	0.28	0.31	0.26	0.22	0.26	0.23	0.27	0.24	0.26	0.25	0.29	0.27	0.28	0.29	0.27				
32	8	<b>0.29</b>	0.27	0.23	0.31	0.28	0.30	0.28	0.27	0.34	0.30	0.26	0.27	0.24	0.23	0.23	0.23	0.24	0.28	0.29	0.26	0.24	0.27	0.23	0.29	0.28	0.28	0.27	0.29	0.27	0.27					
32	16	0.28	0.24	0.26	0.30	<b>0.32</b>	0.29	0.26	0.32	0.27	0.26	0.27	0.24	0.25	0.27	0.26	0.27	0.27	0.29	0.28	0.26	0.27	0.28	0.27	0.29	0.28	0.27	0.28	0.27	0.28	0.27	0.28				
64	4	0.24	0.29	0.23	0.28	0.28	0.35	0.28	0.30	0.28	0.25	0.23	0.27	0.26	0.24	0.23	0.28	0.30	0.29	0.26	0.22	0.26	0.24	0.23	0.28	0.27	0.26	0.27	0.26	0.27	0.26					
64	8	0.26	0.27	0.22	0.30	0.29	0.29	0.27	0.26	0.32	0.29	0.27	0.24	0.21	0.23	0.23	0.26	0.28	0.30	0.28	0.27	0.26	0.24	0.30	0.23	0.27	0.27	0.26	0.27	0.26	0.27	0.26				
64	16	<b>0.29</b>	0.24	0.26	0.30	0.31	0.29	0.27	0.32	0.27	0.23	0.24	0.26	<b>0.26</b>	<b>0.27</b>	<b>0.29</b>	<b>0.29</b>	<b>0.28</b>	<b>0.27</b>	<b>0.27</b>	<b>0.29</b>	<b>0.28</b>	<b>0.27</b>	<b>0.27</b>	<b>0.29</b>	<b>0.28</b>										
64	32	0.24	0.27	0.22	0.28	0.28	0.34	0.28	0.29	0.28	0.25	0.24	0.26	0.26	0.23	0.24	0.25	0.26	0.27	0.26	0.23	0.24	0.25	0.26	0.27	0.26	0.27	0.26	0.27	0.26						
128	8	0.27	0.23	0.22	0.29	0.29	0.30	0.26	0.27	0.24	0.24	0.25	0.26	0.30	0.33	0.30	0.29	0.33	0.30	0.23	0.26	0.23	0.27	0.26	0.27	0.26	0.27	0.26	0.27	0.26	0.27					
128	16	0.28	0.29	0.24	<b>0.33</b>	0.31	0.31	0.29	0.27	0.33	0.29	0.24	0.24	0.26	0.26	0.26	0.27	<b>0.30</b>	0.30	0.27	0.28	0.27	0.23	0.29	0.25	0.27	0.26	0.27	0.28	0.27	0.26	0.27				

Table 8: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on Belebele (Bandarkar et al., 2023).

Model	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
LLAMA-70B-Chat	-	0.42	-	-	0.65	0.66	0.69	0.79	0.68	0.63	0.72	-	0.41	-	-	0.68	0.57	0.56	-	-	0.66	0.70	0.67	0.67	-	-	0.39	0.47	-	-	0.62	0.61				
Mistral-7B-Instruct	-	0.09	-	-	0.54	0.51	0.52	0.65	0.59	0.35	0.60	-	0.02	-	-	0.51	0.42	0.41	-	-	0.55	0.54	0.62	0.55	-	-	0.04	0.31	-	-	0.60	0.44				
GPT-3.5-Turbo	-	0.69	-	-	0.85	0.82	0.85	0.79	0.83	0.87	<b>0.91</b>	-	0.64	-	-	0.81	<b>0.86</b>	<b>0.89</b>	-	-	0.85	0.79	<b>0.91</b>	0.80	-	-	<b>0.84</b>	0.78	-	-	0.78	0.76				
GPT-4	-	<b>0.92</b>	-	-	-	<b>0.88</b>	<b>0.88</b>	<b>0.88</b>	<b>0.92</b>	<b>0.86</b>	<b>0.87</b>	<b>0.91</b>	-	-	-	<b>0.86</b>	<b>0.84</b>	<b>0.86</b>	-	-	<b>0.87</b>	<b>0.88</b>	<b>0.87</b>	<b>0.87</b>	-	-	<b>0.81</b>	<b>0.85</b>	-	-	0.87	<b>0.87</b>				
PALM2	-	0.87	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-				
rank quantisation	af	ar	as	bg	bn	cs	da	de	en	es	fi	fr	gu	he	hi	hr	id	it	ja	ko	ml	mr	nb	nl	pt	ru	sv	sw	ta	te	th	tr	ur	vi	zh	avg
8	4	0.42	0.33	0.33	0.41	0.33	0.44	0.43	0.47	0.52	0.49	0.35	<b>0.54</b>	<b>0.31</b>	0.29	<b>0.50</b>	0.50	0.44	0.40	0.45	0.27	0.46	0.55	0.46	0.48	0.44	0.47	0.30	0.27	0.34	0.38	0.47	0.40			
8	8	<b>0.47</b>	0.38	0.31	0.44	0.37	0.49	0.44	0.52	0.57	0.50	0.41	0.49	0.27	0.31	0.49	0.45	0.47	0.49	0.41	0.51	0.29	0.30	0.47	0.54	0.49	0.55	0.49	0.34	0.33	0.39	0.39	0.52	0.42		
8	16	0.45	0.39	0.34	0.46	0.37	0.49	0.43	0.45	0.51	0.50	0.41	0.51	0.27	0.30	0.49	0.43	0.47	0.49	0.43	0.50	0.27	0.29	0.47	0.54	0.49	0.55	0.49	0.34	0.33	0.39	0.39	0.52	0.42		
8	32	0.42	0.39	0.31	0.43	<b>0.40</b>	0.47	0.43	0.52	0.58	0.51	0.41	0.52	0.27	0.31	0.49	0.45	0.47	0.49	0.41	0.49	0.27	0.29	0.47	0.54	0.49	0.55	0.49	0.34	0.33	0.39	0.39	0.52	0.42		
16	4	0.41	0.35	0.34	0.41	0.33	0.46	0.43	0.47	0.53	0.50	0.36	0.52	0.31	0.30	0.49	0.45	0.47	0.49	0.43	0.49	0.27	0.29	0.47	0.54	0.49	0.55	0.49	0.34	0.33	0.39	0.39	0.52	0.42		
16	8	0.44	0.38	0.31	0.47	0.34	0.46	0.43	0.49	0.53	0.50	0.42	0.49	0.27	0.30	0.49	0.45	0.47	0.49	0.43	0.49	0.27	0.29	0.47	0.54	0.49	0.55	0.49	0.34	0.33	0.39	0.39	0.52	0.42		
16	16	0.44	0.39	0.33	0.45																															

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
64	8	13.05
64	16	13.11
128	8	13.28
128	16	13.23

Table 10: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on AlpacaEval (Li et al., 2023b).

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
8	4	11.86
8	8	12.81
8	16	13.23
16	4	11.97
16	8	13.06
16	16	13.36
32	4	11.72
32	8	13.15
32	16	13.36
64	4	11.47
64	8	13.67
64	16	13.73
128	4	11.10
128	8	13.05
128	16	0.00

Table 11: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on AlpacaEval (Li et al., 2023b).

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
8	4	11.61
8	8	13.56
8	16	13.36
16	4	11.99
16	8	13.38
16	16	13.73
32	4	11.86
32	8	12.47
32	16	13.73
64	4	11.61
64	8	13.47
64	16	13.86
128	4	11.35
128	8	13.56
128	16	13.73

Table 12: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on AlpacaEval (Li et al., 2023b).

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
64	8	20.47
64	16	18.89
128	8	19.55
128	16	19.33

Table 13: Detailed performance of various ALPACA finetuned MISTRAL-7B models on AlpacaEval (Li et al., 2023b).

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
8	4	15.04
8	8	20.67
8	16	21.27
16	4	15.77
16	8	22.07
16	16	21.08
32	4	15.77
32	8	21.07
32	16	21.14
64	4	16.90
64	8	21.30
64	16	21.27
128	4	17.77
128	8	22.04
128	16	21.83

Table 14: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on AlpacaEval (Li et al., 2023b).

Model	win_rate	
LLaMA-70B-Chat	22.36	
Mistral-7B-Instruct	35.13	
GPT-4	<b>93.78</b>	
PALM2	79.66	
rank	quantisation	win_rate
8	4	15.98
8	8	21.42
8	16	20.71
16	4	15.90
16	8	20.55
16	16	21.27
32	4	16.98
32	8	22.45
32	16	22.08
64	4	16.83
64	8	21.07
64	16	18.70
128	4	16.52
128	8	20.92
128	16	21.33

Table 15: Detailed performance of various MULTIAL-PACA finetuned MISTRAL-7B models on AlpacaEval (Li et al., 2023b).

Models	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRv6	<b>0.76</b>	<b>0.80</b>	<b>0.87</b>	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
64	8	0.31	0.20	0.28	0.24	0.43	0.19	0.08	0.25
64	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35
128	8	0.31	0.20	0.29	0.25	0.43	0.19	0.08	0.25
128	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35

Table 16: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on MLQA (Lewis et al., 2020).

Models	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRV6	<b>0.76</b>	<b>0.80</b>	<b>0.87</b>	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
8	4	0.24	0.18	0.75	0.49	0.37	0.41	0.08	0.36
8	8	0.31	0.20	0.31	0.24	0.43	0.19	0.08	0.25
8	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35
16	4	0.26	0.18	0.76	0.51	0.37	0.42	0.08	0.37
16	8	0.31	0.20	0.30	0.24	0.43	0.19	0.08	0.25
16	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35
32	4	0.29	0.18	0.77	0.56	0.37	0.48	0.08	0.39
32	8	0.32	0.20	0.31	0.25	0.43	0.19	0.08	0.25
32	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35
64	4	0.31	0.18	0.79	0.65	0.37	0.54	0.08	0.42
64	8	0.32	0.19	0.31	0.25	0.42	0.19	0.08	0.25
64	16	0.27	0.36	0.49	0.43	0.43	0.33	0.15	0.35
128	4	0.33	0.18	0.80	0.72	0.38	0.61	0.08	0.44
128	8	0.37	0.20	0.34	0.26	0.42	0.20	0.08	0.27
128	16	0.28	0.36	0.49	0.43	0.43	0.33	0.15	0.35

Table 17: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on MLQA (Lewis et al., 2020).

Models	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRV6	<b>0.76</b>	<b>0.80</b>	<b>0.87</b>	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
8	4	0.24	0.18	0.74	0.48	0.37	0.40	0.08	0.36
8	8	0.30	0.20	0.31	0.25	0.43	0.20	0.08	0.25
8	16	0.27	0.36	0.50	0.43	0.43	0.33	0.15	0.35
16	4	0.24	0.18	0.75	0.48	0.37	0.41	0.07	0.36
16	8	0.34	0.20	0.32	0.24	0.42	0.19	0.08	0.26
16	16	0.27	0.36	0.49	0.43	0.43	0.33	0.15	0.35
32	4	0.22	0.18	0.74	0.46	0.38	0.37	0.08	0.35
32	8	0.30	0.20	0.31	0.25	0.42	0.19	0.08	0.25
32	16	0.28	0.36	0.49	0.43	0.43	0.33	0.15	0.35
64	4	0.23	0.18	0.75	0.48	0.38	0.38	0.08	0.35
64	8	0.40	0.20	0.47	0.34	0.43	0.23	0.08	0.31
64	16	0.28	0.36	0.48	0.43	0.43	0.33	0.15	0.35
128	4	0.30	0.18	0.78	0.62	0.37	0.52	0.08	0.41
128	8	0.32	0.20	0.32	0.25	0.41	0.19	0.08	0.25
128	16	0.27	0.36	0.49	0.43	0.43	0.33	0.15	0.35

Table 18: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on MLQA (Lewis et al., 2020).

Models	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRv6	<b>0.76</b>	<b>0.80</b>	0.87	<b>0.82</b>	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
64	8	0.61	0.73	0.89	0.81	0.60	0.80	0.41	0.69
64	16	0.61	0.74	<b>0.89</b>	0.81	0.61	0.81	0.41	0.70
128	8	0.61	0.74	<b>0.89</b>	0.81	0.60	0.80	0.41	0.69
128	16	0.61	0.74	0.88	0.81	0.60	0.81	0.41	0.70

Table 19: Detailed performance of various ALPACA finetuned MISTRAL-7B models on MLQA (Lewis et al., 2020).

Model	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRv6	<b>0.76</b>	<b>0.80</b>	0.87	0.82	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
8	4	0.55	0.71	0.86	0.81	0.53	0.78	0.39	0.66
8	8	0.62	0.74	0.88	0.81	0.59	0.80	0.41	0.69
8	16	0.61	0.74	0.88	0.82	0.60	0.80	0.41	0.70
16	4	0.55	0.71	0.86	0.81	0.53	0.78	0.39	0.66
16	8	0.60	0.75	<b>0.89</b>	0.82	0.59	0.81	0.41	0.70
16	16	0.61	0.74	0.88	0.81	0.59	0.80	0.41	0.69
32	4	0.54	0.71	0.86	0.81	0.53	0.77	0.39	0.66
32	8	0.60	0.74	0.88	0.81	0.60	0.80	0.41	0.69
32	16	0.61	0.74	0.88	0.81	0.60	0.81	0.41	0.69
64	4	0.54	0.71	0.86	<b>0.82</b>	0.52	0.78	0.39	0.66
64	8	0.60	0.74	0.88	0.81	0.60	0.80	0.41	0.69
64	16	0.61	0.74	0.88	0.81	0.60	0.81	0.41	0.70
128	4	0.55	0.72	0.86	0.82	0.52	0.78	0.39	0.66
128	8	0.59	0.75	0.88	0.81	0.59	0.81	0.41	0.69
128	16	0.61	0.74	0.88	0.81	0.61	0.81	0.41	0.70

Table 20: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on MLQA (Lewis et al., 2020).

Model	ar	de	en	es	hi	vi	zh	avg	
LLaMA-7B-Chat	0.05	0.45	0.70	0.52	0.00	0.42	0.07	0.32	
LLaMA-3B-Chat	0.55	0.73	0.62	0.62	0.00	0.09	0.09	0.39	
LLaMA-70B-Chat	0.04	0.42	0.56	0.34	0.01	0.27	0.05	0.24	
Mistral-7B-Instruct	0.17	0.30	0.32	0.33	0.20	0.36	0.02	0.24	
DV003	0.38	0.58	0.75	0.63	0.25	0.48	0.32	0.48	
GPT-3.5-Turbo	0.48	0.51	0.73	0.54	0.51	0.59	0.57	0.56	
GPT-4	0.59	0.65	0.80	0.70	0.57	0.72	0.67	0.67	
TULRv6	<b>0.76</b>	<b>0.80</b>	0.87	0.82	<b>0.82</b>	<b>0.82</b>	<b>0.78</b>	<b>0.81</b>	
XLM-R	0.67	0.70	0.83	0.74	0.71	0.74	0.62	0.72	
mBERT	0.52	0.59	0.80	0.67	0.50	0.61	0.60	0.61	
mT5	0.57	0.62	0.82	0.67	0.55	0.66	0.62	0.64	
PALM2	0.22	0.44	0.69	0.45	0.36	0.49	0.06	0.39	
rank	quantisation	ar	de	en	es	hi	vi	zh	avg
8	4	0.56	0.71	0.86	0.82	0.54	0.78	0.39	0.66
8	8	0.60	0.74	0.89	0.81	0.58	0.80	0.41	0.69
8	16	0.61	0.74	0.88	0.81	0.60	0.81	0.41	0.69
16	4	0.54	0.71	0.85	0.81	0.54	0.78	0.39	0.66
16	8	0.60	0.74	<b>0.89</b>	0.81	0.59	0.81	0.41	0.69
16	16	0.61	0.75	0.88	0.82	0.60	0.81	0.41	0.70
32	4	0.55	0.71	0.86	0.81	0.54	0.78	0.39	0.66
32	8	0.58	0.74	0.88	0.81	0.60	0.80	0.41	0.69
32	16	0.60	0.75	0.88	0.82	0.62	0.81	0.41	0.70
64	4	0.55	0.71	0.86	0.81	0.53	0.77	0.39	0.66
64	8	0.59	0.74	0.88	0.82	0.61	0.81	0.41	0.69
64	16	0.60	0.74	0.88	<b>0.82</b>	0.61	0.81	0.41	0.70
128	4	0.57	0.72	0.86	0.81	0.54	0.78	0.39	0.67
128	8	0.60	0.74	0.89	0.81	0.59	0.81	0.41	0.69
128	16	0.61	0.74	0.88	0.81	0.60	0.80	0.41	0.70

Table 21: Detailed performance of various MULTIALPACA finetuned MISTRAL-7B models on MLQA (Lewis et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	–	–	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	–	–	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	–	–	0.90	
TULRv6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	–	–	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	–	–	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	–	–	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	–	–	0.50	
PALM2	0.97	–	<b>0.98</b>	0.98	–	0.89	–	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
64	8	0.47	0.52	0.66	0.71	0.52	0.51	0.54	0.57	0.55	0.63	0.70	0.58
64	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.64	0.70	0.58
128	8	0.48	0.53	0.66	0.71	0.52	0.51	0.54	0.56	0.55	0.64	0.69	0.58
128	16	0.48	0.52	0.66	0.70	0.52	0.51	0.54	0.58	0.55	0.64	0.70	0.58

Table 22: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on XCOPA (Ponti et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	—	—	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	—	—	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	—	—	0.90	
TULRV6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	—	—	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	—	—	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	—	—	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	—	—	0.50	
PALM2	0.97	—	<b>0.98</b>	0.98	—	0.89	—	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
8	4	0.46	0.51	0.66	0.69	0.52	0.51	0.56	0.56	0.55	0.63	0.69	0.58
8	8	0.47	0.53	0.66	0.71	0.52	0.52	0.54	0.57	0.56	0.63	0.69	0.58
8	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.64	0.70	0.58
16	4	0.46	0.50	0.66	0.69	0.52	0.52	0.56	0.56	0.55	0.63	0.69	0.58
16	8	0.47	0.52	0.67	0.72	0.52	0.52	0.54	0.56	0.55	0.63	0.69	0.58
16	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.64	0.70	0.58
32	4	0.46	0.50	0.66	0.69	0.52	0.51	0.55	0.56	0.55	0.63	0.69	0.57
32	8	0.47	0.53	0.65	0.71	0.51	0.52	0.54	0.58	0.55	0.64	0.69	0.58
32	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.63	0.70	0.58
64	4	0.46	0.50	0.65	0.69	0.52	0.52	0.55	0.55	0.55	0.63	0.69	0.57
64	8	0.47	0.52	0.66	0.70	0.51	0.52	0.54	0.57	0.55	0.63	0.68	0.58
64	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.64	0.70	0.58
128	4	0.46	0.50	0.65	0.69	0.52	0.51	0.56	0.56	0.55	0.63	0.69	0.57
128	8	0.47	0.52	0.66	0.70	0.52	0.52	0.54	0.57	0.55	0.63	0.70	0.58
128	16	0.47	0.52	0.66	0.70	0.51	0.51	0.54	0.57	0.55	0.64	0.70	0.58

Table 23: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on XCOPA (Ponti et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	—	—	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	—	—	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	—	—	0.90	
TULRv6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	—	—	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	—	—	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	—	—	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	—	—	0.50	
PALM2	0.97	—	<b>0.98</b>	0.98	—	0.89	—	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
8	4	0.46	0.51	0.66	0.69	0.52	0.51	0.56	0.56	0.55	0.63	0.69	0.58
8	8	0.47	0.52	0.65	0.71	0.52	0.52	0.54	0.56	0.55	0.64	0.70	0.58
8	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.63	0.70	0.58
16	4	0.46	0.51	0.66	0.69	0.52	0.51	0.56	0.55	0.55	0.63	0.69	0.57
16	8	0.47	0.53	0.65	0.70	0.52	0.51	0.54	0.57	0.55	0.63	0.69	0.58
16	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.64	0.70	0.58
32	4	0.46	0.50	0.66	0.69	0.52	0.51	0.56	0.56	0.55	0.63	0.69	0.58
32	8	0.46	0.53	0.66	0.70	0.52	0.51	0.55	0.57	0.55	0.63	0.70	0.58
32	16	0.47	0.52	0.66	0.70	0.52	0.51	0.54	0.58	0.55	0.63	0.70	0.58
64	4	0.46	0.50	0.65	0.69	0.52	0.51	0.55	0.55	0.54	0.63	0.69	0.57
64	8	0.47	0.53	0.65	0.71	0.52	0.52	0.54	0.57	0.54	0.63	0.69	0.58
64	16	0.48	0.53	0.66	0.70	0.51	0.51	0.54	0.57	0.55	0.64	0.70	0.58
128	4	0.46	0.50	0.66	0.69	0.52	0.51	0.56	0.56	0.55	0.63	0.69	0.57
128	8	0.47	0.52	0.65	0.70	0.51	0.51	0.54	0.57	0.55	0.63	0.69	0.58
128	16	0.48	0.52	0.66	0.70	0.52	0.51	0.54	0.57	0.55	0.63	0.70	0.58

Table 24: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on XCOPA (Ponti et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	—	—	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	—	—	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	—	—	0.90	
TULRv6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	—	—	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	—	—	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	—	—	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	—	—	0.50	
PALM2	0.97	—	<b>0.98</b>	0.98	—	0.89	—	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
64	8	0.48	0.55	0.64	0.74	0.51	0.52	0.55	0.57	0.57	0.63	0.73	0.59
64	16	0.48	0.53	0.64	0.74	0.51	0.51	0.55	0.57	0.57	0.63	0.72	0.59
128	8	0.48	0.54	0.65	0.74	0.52	0.51	0.55	0.56	0.57	0.64	0.73	0.59
128	16	0.48	0.53	0.64	0.73	0.52	0.51	0.55	0.57	0.57	0.63	0.72	0.59

Table 25: Detailed performance of various ALPACA finetuned MISTRAL-7B models on XCOPA (Ponti et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	—	—	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	—	—	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	—	—	0.90	
TULRV6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	—	—	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	—	—	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	—	—	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	—	—	0.50	
PALM2	0.97	—	<b>0.98</b>	0.98	—	0.89	—	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
8	4	0.47	0.54	0.63	0.72	0.52	0.52	0.54	0.59	0.57	0.61	0.74	0.59
8	8	0.48	0.54	0.64	0.74	0.51	0.51	0.54	0.56	0.57	0.62	0.73	0.59
8	16	0.48	0.53	0.64	0.74	0.51	0.51	0.55	0.57	0.57	0.64	0.73	0.59
16	4	0.46	0.54	0.63	0.72	0.52	0.52	0.55	0.58	0.57	0.62	0.74	0.59
16	8	0.47	0.54	0.64	0.73	0.51	0.51	0.55	0.57	0.57	0.63	0.73	0.59
16	16	0.48	0.53	0.64	0.74	0.52	0.51	0.55	0.57	0.57	0.63	0.72	0.59
32	4	0.47	0.54	0.64	0.73	0.52	0.52	0.55	0.58	0.56	0.62	0.74	0.59
32	8	0.47	0.53	0.64	0.74	0.51	0.52	0.54	0.57	0.57	0.63	0.73	0.59
32	16	0.47	0.53	0.64	0.74	0.52	0.51	0.55	0.58	0.56	0.64	0.72	0.59
64	4	0.47	0.54	0.64	0.73	0.52	0.53	0.55	0.58	0.56	0.62	0.73	0.59
64	8	0.48	0.54	0.64	0.74	0.51	0.51	0.54	0.57	0.57	0.63	0.73	0.59
64	16	0.48	0.53	0.64	0.73	0.51	0.51	0.55	0.58	0.56	0.64	0.73	0.59
128	4	0.46	0.53	0.64	0.73	0.52	0.51	0.54	0.58	0.56	0.62	0.73	0.58
128	8	0.47	0.55	0.64	0.74	0.51	0.51	0.55	0.56	0.56	0.63	0.72	0.59
128	16	0.48	0.54	0.64	0.73	0.51	0.51	0.55	0.57	0.57	0.63	0.74	0.59

Table 26: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on XCOPA (Ponti et al., 2020).

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.51	0.51	0.59	0.71	0.50	0.51	0.50	0.52	0.53	0.58	0.59	0.55	
LLaMA-3B-Chat	0.51	0.49	0.72	0.80	0.50	0.50	0.00	0.00	0.54	0.02	0.70	0.44	
LLaMA-70B-Chat	0.47	0.44	0.40	0.53	0.37	0.34	0.21	0.25	0.44	0.52	0.35	0.39	
Mistral-7B-Instruct	0.52	0.50	0.61	0.73	0.51	0.49	0.41	0.51	0.54	0.50	0.50	0.53	
DV003	0.88	0.75	0.91	0.96	0.55	0.64	0.54	0.67	0.88	—	—	0.75	
GPT-3.5-Turbo	0.91	0.72	0.90	0.95	0.55	0.82	0.59	0.78	0.91	—	—	0.79	
GPT-4	<b>0.99</b>	<b>0.93</b>	0.98	<b>1.00</b>	0.59	<b>0.94</b>	0.80	0.88	<b>0.97</b>	—	—	0.90	
TULRv6	0.77	0.78	0.93	0.96	<b>0.61</b>	0.69	<b>0.85</b>	0.87	0.93	—	—	0.82	
BLOOMZ	0.48	0.55	0.86	0.74	0.50	0.60	0.67	0.50	0.54	—	—	0.60	
XGLM	0.66	0.59	0.69	0.69	0.47	0.63	0.56	0.62	0.58	—	—	0.61	
mT5	0.50	0.50	0.49	0.50	0.51	0.50	0.49	0.51	0.49	—	—	0.50	
PALM2	0.97	—	<b>0.98</b>	0.98	—	0.89	—	<b>0.95</b>	0.93	<b>0.98</b>	<b>0.99</b>	<b>0.96</b>	
rank	quantisation	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	avg
8	4	0.47	0.53	0.63	0.72	0.52	0.52	0.55	0.58	0.57	0.62	0.73	0.59
8	8	0.47	0.54	0.64	0.73	0.51	0.51	0.55	0.57	0.57	0.63	0.73	0.59
8	16	0.48	0.53	0.65	0.73	0.51	0.52	0.55	0.57	0.57	0.64	0.73	0.59
16	4	0.46	0.54	0.63	0.73	0.52	0.53	0.55	0.59	0.57	0.62	0.73	0.59
16	8	0.48	0.54	0.64	0.74	0.52	0.52	0.55	0.57	0.56	0.63	0.73	0.59
16	16	0.48	0.54	0.64	0.74	0.51	0.52	0.55	0.58	0.56	0.63	0.73	0.59
32	4	0.47	0.54	0.63	0.73	0.52	0.53	0.55	0.58	0.57	0.62	0.73	0.59
32	8	0.48	0.54	0.64	0.73	0.51	0.53	0.54	0.57	0.56	0.63	0.72	0.59
32	16	0.48	0.53	0.64	0.74	0.52	0.52	0.54	0.57	0.56	0.63	0.74	0.59
64	4	0.47	0.54	0.63	0.72	0.52	0.52	0.55	0.58	0.57	0.61	0.73	0.59
64	8	0.48	0.54	0.65	0.74	0.52	0.51	0.55	0.57	0.56	0.63	0.73	0.59
64	16	0.48	0.54	0.64	0.74	0.51	0.52	0.55	0.58	0.56	0.64	0.73	0.59
128	4	0.47	0.54	0.63	0.72	0.52	0.52	0.55	0.58	0.57	0.61	0.73	0.58
128	8	0.48	0.54	0.64	0.74	0.52	0.52	0.55	0.57	0.57	0.63	0.72	0.59
128	16	0.48	0.53	0.64	0.74	0.52	0.51	0.55	0.57	0.56	0.63	0.73	0.59

Table 27: Detailed performance of various MULTIALPACA finetuned MISTRAL-7B models on XCOPA (Ponti et al., 2020).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.89	0.84	0.84	0.18	0.77	0.80	0.77	0.78	0.79	—	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
64	8	0.34	0.34	0.36	0.34	0.35	0.37	0.36	0.34	0.33	0.33	0.34	0.34	0.37	0.33	0.35	0.35
64	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.33	0.34	0.35	0.34	0.36	0.35	0.35
128	8	0.34	0.34	0.36	0.34	0.34	0.36	0.35	0.34	0.34	0.33	0.34	0.36	0.34	0.36	0.35	0.35
128	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.33	0.34	0.35	0.36	0.34	0.36	0.35

Table 28: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on XNLI (Conneau et al., 2018).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT-4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.84	0.89	0.84	0.18	0.77	0.80	0.77	0.78	0.79	–	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
8	4	0.32	0.33	0.34	0.33	0.33	0.33	0.32	0.34	0.33	0.33	0.34	0.33	0.31	0.33	0.33	0.33
8	8	0.34	0.34	0.35	0.34	0.35	0.35	0.36	0.33	0.35	0.32	0.34	0.33	0.37	0.34	0.36	0.34
8	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.34	0.34	0.34	0.35	0.34	0.36	0.35
16	4	0.32	0.33	0.34	0.33	0.33	0.33	0.32	0.34	0.33	0.33	0.33	0.33	0.32	0.34	0.33	0.33
16	8	0.35	0.33	0.36	0.34	0.35	0.36	0.35	0.33	0.34	0.33	0.34	0.34	0.37	0.34	0.36	0.35
16	16	0.34	0.35	0.35	0.34	0.34	0.34	0.37	0.33	0.35	0.33	0.34	0.35	0.35	0.34	0.36	0.35
32	4	0.32	0.34	0.34	0.33	0.33	0.33	0.32	0.34	0.32	0.33	0.33	0.32	0.32	0.34	0.33	0.33
32	8	0.34	0.34	0.35	0.33	0.34	0.34	0.37	0.36	0.34	0.33	0.34	0.34	0.37	0.33	0.36	0.35
32	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.33	0.33	0.35	0.35	0.34	0.36	0.35
64	4	0.32	0.34	0.34	0.33	0.33	0.33	0.32	0.34	0.33	0.33	0.33	0.32	0.31	0.34	0.33	0.33
64	8	0.35	0.34	0.35	0.34	0.35	0.37	0.36	0.34	0.35	0.32	0.34	0.34	0.36	0.33	0.35	0.35
64	16	0.34	0.35	0.34	0.34	0.34	0.36	0.37	0.33	0.35	0.33	0.34	0.35	0.35	0.34	0.36	0.35
128	4	0.32	0.33	0.35	0.33	0.34	0.33	0.33	0.33	0.34	0.33	0.33	0.32	0.31	0.34	0.33	0.33
128	8	0.35	0.34	0.36	0.34	0.36	0.36	0.37	0.33	0.33	0.33	0.34	0.35	0.36	0.33	0.36	0.35
128	16	0.34	0.35	0.34	0.34	0.34	0.36	0.36	0.33	0.34	0.33	0.34	0.35	0.35	0.34	0.35	0.34

Table 29: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on XNLI (Conneau et al., 2018).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT-4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.84	0.89	0.84	0.18	0.77	0.80	0.77	0.78	0.79	–	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
8	4	0.32	0.33	0.34	0.33	0.33	0.33	0.32	0.34	0.33	0.33	0.34	0.33	0.32	0.34	0.33	0.33
8	8	0.34	0.34	0.35	0.34	0.35	0.35	0.35	0.34	0.34	0.33	0.34	0.34	0.35	0.32	0.38	0.34
8	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.34	0.34	0.35	0.36	0.34	0.36	0.35
16	4	0.32	0.33	0.34	0.33	0.33	0.34	0.32	0.34	0.33	0.33	0.34	0.33	0.32	0.34	0.33	0.33
16	8	0.34	0.35	0.36	0.34	0.34	0.38	0.37	0.34	0.33	0.32	0.34	0.35	0.37	0.33	0.36	0.35
16	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.34	0.34	0.34	0.35	0.35	0.34	0.36	0.35
32	4	0.32	0.34	0.34	0.33	0.34	0.33	0.32	0.34	0.32	0.33	0.33	0.33	0.31	0.34	0.33	0.33
32	8	0.35	0.34	0.34	0.34	0.34	0.35	0.35	0.33	0.34	0.33	0.34	0.35	0.36	0.32	0.35	0.34
32	16	0.34	0.34	0.34	0.34	0.34	0.35	0.37	0.33	0.34	0.33	0.34	0.35	0.35	0.33	0.34	0.34
64	4	0.32	0.34	0.34	0.33	0.34	0.33	0.33	0.34	0.33	0.33	0.33	0.33	0.31	0.34	0.33	0.33
64	8	0.35	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.34	0.32	0.34	0.35	0.37	0.33	0.34	0.35
64	16	0.34	0.34	0.35	0.34	0.34	0.35	0.37	0.33	0.34	0.32	0.34	0.35	0.37	0.33	0.34	0.35
128	4	0.32	0.33	0.34	0.33	0.33	0.33	0.32	0.34	0.33	0.33	0.33	0.33	0.32	0.34	0.33	0.33
128	8	0.34	0.34	0.37	0.34	0.35	0.36	0.36	0.34	0.34	0.32	0.34	0.35	0.36	0.33	0.37	0.35
128	16	0.34	0.35	0.35	0.34	0.34	0.35	0.37	0.33	0.35	0.33	0.33	0.35	0.35	0.33	0.36	0.35

Table 30: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on XNLI (Conneau et al., 2018).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.84	0.89	0.84	0.18	0.77	0.80	0.77	0.78	0.79	–	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
64	8	0.47	0.55	0.56	0.50	0.69	0.59	0.55	0.50	0.54	0.37	0.52	0.45	0.47	0.47	0.59	0.52
64	16	0.47	0.57	0.57	0.53	0.68	0.62	0.57	0.51	0.54	0.38	0.53	0.45	0.46	0.47	0.60	0.53
128	8	0.46	0.53	0.54	0.51	0.68	0.58	0.54	0.49	0.53	0.36	0.52	0.43	0.46	0.46	0.57	0.51
128	16	0.48	0.56	0.55	0.51	0.66	0.61	0.56	0.50	0.53	0.37	0.53	0.43	0.46	0.45	0.59	0.52

Table 31: Detailed performance of various ALPACA finetuned MISTRAL-7B models on XNLI (Conneau et al., 2018).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.84	0.89	0.84	0.18	0.77	0.80	0.77	0.78	0.79	–	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
8	4	0.51	0.56	0.54	0.53	0.66	0.57	0.55	0.46	0.53	0.43	0.49	0.45	0.43	0.47	0.58	0.52
8	8	0.49	0.51	0.53	0.51	0.67	0.58	0.53	0.49	0.52	0.37	0.53	0.42	0.46	0.46	0.58	0.51
8	16	0.48	0.55	0.53	0.53	0.65	0.60	0.56	0.50	0.53	0.38	0.52	0.42	0.46	0.47	0.58	0.52
16	4	0.52	0.56	0.54	0.52	0.65	0.58	0.56	0.47	0.53	0.44	0.49	0.45	0.43	0.48	0.58	0.52
16	8	0.48	0.51	0.52	0.51	0.67	0.57	0.54	0.49	0.52	0.37	0.51	0.42	0.45	0.45	0.59	0.51
16	16	0.48	0.55	0.54	0.51	0.66	0.60	0.55	0.50	0.53	0.38	0.52	0.42	0.47	0.47	0.58	0.52
32	4	0.53	0.56	0.53	0.52	0.65	0.57	0.55	0.47	0.53	0.44	0.48	0.44	0.44	0.47	0.58	0.52
32	8	0.49	0.52	0.52	0.51	0.68	0.58	0.54	0.49	0.53	0.38	0.52	0.42	0.45	0.46	0.59	0.51
32	16	0.47	0.54	0.55	0.51	0.65	0.60	0.55	0.50	0.53	0.38	0.52	0.42	0.47	0.46	0.59	0.52
64	4	0.52	0.56	0.54	0.52	0.65	0.58	0.55	0.48	0.52	0.44	0.49	0.44	0.46	0.47	0.58	0.52
64	8	0.48	0.52	0.53	0.51	0.68	0.58	0.54	0.49	0.53	0.38	0.53	0.42	0.46	0.46	0.59	0.51
64	16	0.47	0.54	0.56	0.51	0.65	0.60	0.55	0.50	0.53	0.38	0.52	0.43	0.47	0.46	0.59	0.52
128	4	0.50	0.57	0.54	0.49	0.64	0.58	0.54	0.47	0.52	0.42	0.50	0.45	0.47	0.45	0.58	0.51
128	8	0.48	0.53	0.55	0.51	0.69	0.59	0.55	0.49	0.54	0.38	0.53	0.43	0.46	0.48	0.60	0.52
128	16	0.46	0.53	0.55	0.50	0.65	0.59	0.54	0.50	0.52	0.37	0.52	0.43	0.45	0.45	0.58	0.51

Table 32: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on XNLI (Conneau et al., 2018).

Model	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg	
LLaMA-7B-Chat	0.39	0.45	0.45	0.39	0.56	0.50	0.50	0.37	0.48	0.33	0.35	0.40	0.36	0.41	0.45	0.43	
LLaMA-3B-Chat	0.37	0.51	0.50	0.00	0.55	0.51	0.52	0.00	0.50	0.36	0.00	0.45	0.29	0.48	0.48	0.37	
LLaMA-70B-Chat	0.35	0.55	0.58	0.41	0.63	0.55	0.55	0.31	0.55	0.39	0.35	0.49	0.42	0.49	0.51	0.48	
Mistral-7B-Instruct	0.35	0.40	0.38	0.35	0.43	0.41	0.41	0.35	0.43	0.34	0.33	0.40	0.34	0.33	0.40	0.38	
DV003	0.52	0.62	0.66	0.60	0.80	0.71	0.66	0.48	0.62	0.50	0.51	0.58	0.50	0.56	0.58	0.59	
GPT-3.5-Turbo	0.59	0.64	0.67	0.65	0.76	0.70	0.68	0.55	0.62	0.56	0.54	0.63	0.49	0.61	0.62	0.62	
GPT-4	0.73	0.77	0.79	0.79	0.85	0.79	0.79	0.72	0.74	0.71	0.69	0.76	0.68	0.74	0.75	0.75	
TULRv6	<b>0.89</b>	<b>0.91</b>	<b>0.90</b>	<b>0.90</b>	<b>0.93</b>	<b>0.91</b>	<b>0.91</b>	<b>0.86</b>	<b>0.89</b>	<b>0.85</b>	<b>0.88</b>	<b>0.88</b>	<b>0.83</b>	<b>0.89</b>	<b>0.88</b>	<b>0.89</b>	
BLOOMZ	0.61	0.47	0.54	0.47	0.67	0.61	0.61	0.57	0.53	0.50	0.44	0.43	0.50	0.61	0.57	0.54	
XLM-R	0.77	0.83	0.82	0.81	0.89	0.84	0.82	0.76	0.79	0.71	0.77	0.78	0.72	0.79	0.78	0.79	
mBERT	0.64	0.68	0.70	0.65	0.81	0.73	0.73	0.59	0.68	0.50	0.54	0.61	0.57	0.69	0.68	0.65	
XGLM	0.46	0.49	0.46	0.49	0.53	0.46	0.49	0.47	0.49	0.45	0.47	0.45	0.43	0.48	0.49	0.47	
mT5	0.73	0.79	0.77	0.77	0.85	0.80	0.79	0.71	0.77	0.69	0.73	0.73	0.68	0.74	0.74	0.75	
PALM2	0.79	0.81	0.84	0.84	0.89	0.84	0.18	0.77	0.80	0.77	0.78	0.79	–	0.79	0.80	0.76	
rank	quantisation	ar	bg	de	el	en	es	fr	hi	ru	sw	th	tr	ur	vi	zh	avg
8	4	0.52	0.56	0.54	0.52	0.66	0.58	0.56	0.47	0.53	0.43	0.49	0.45	0.44	0.47	0.60	0.52
8	8	0.49	0.53	0.51	0.51	0.67	0.58	0.53	0.49	0.53	0.37	0.52	0.42	0.46	0.46	0.59	0.51
8	16	0.48	0.55	0.55	0.51	0.65	0.59	0.55	0.51	0.53	0.38	0.53	0.43	0.46	0.46	0.58	0.52
16	4	0.52	0.56	0.54	0.52	0.66	0.57	0.55	0.48	0.53	0.44	0.50	0.45	0.45	0.48	0.58	0.52
16	8	0.49	0.52	0.54	0.51	0.68	0.58	0.54	0.49	0.53	0.37	0.53	0.42	0.46	0.46	0.60	0.51
16	16	0.48	0.54	0.54	0.51	0.65	0.60	0.56	0.51	0.53	0.38	0.52	0.43	0.46	0.47	0.58	0.52
32	4	0.53	0.57	0.54	0.52	0.65	0.58	0.55	0.48	0.54	0.44	0.49	0.45	0.45	0.49	0.59	0.52
32	8	0.49	0.51	0.53	0.50	0.68	0.58	0.54	0.51	0.53	0.38	0.52	0.43	0.46	0.45	0.61	0.51
32	16	0.47	0.53	0.55	0.49	0.65	0.60	0.55	0.51	0.53	0.38	0.52	0.43	0.46	0.46	0.58	0.51
64	4	0.51	0.57	0.54	0.52	0.65	0.57	0.55	0.48	0.53	0.44	0.51	0.45	0.46	0.49	0.60	0.53
64	8	0.48	0.49	0.53	0.50	0.67	0.57	0.55	0.49	0.51	0.37	0.52	0.42	0.46	0.45	0.58	0.51
64	16	0.46	0.52	0.55	0.51	0.65	0.59	0.54	0.50	0.52	0.38	0.51	0.43	0.47	0.45	0.57	0.51
128	4	0.48	0.54	0.51	0.47	0.60	0.55	0.53	0.46	0.49	0.41	0.52	0.41	0.46	0.45	0.56	0.50
128	8	0.48	0.51	0.54	0.50	0.68	0.58	0.54	0.50	0.54	0.37	0.53	0.42	0.47	0.46	0.60	0.51
128	16	0.48	0.54	0.55	0.51	0.65	0.60	0.55	0.51	0.53	0.39	0.53	0.42	0.46	0.46	0.58	0.52

Table 33: Detailed performance of various MULTIALPACA finetuned MISTRAL-7B models on XNLI ([Conneau et al., 2018](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	<b>0.76</b>	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	<b>0.88</b>	<b>0.86</b>	–	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	<b>0.88</b>	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	–	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	–	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	–	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	–	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
64	8	0.44	0.76	0.58	0.86	0.66	0.50	0.68	0.61	0.33	0.58	0.80	0.56	0.61
64	16	0.44	0.75	0.58	0.85	0.79	0.50	0.72	0.66	0.34	0.58	0.78	0.55	0.63
128	8	0.44	0.76	0.58	0.86	0.66	0.49	0.68	0.62	0.33	0.59	0.79	0.54	0.61
128	16	0.44	0.75	0.58	0.85	0.79	0.50	0.72	0.66	0.34	0.58	0.78	0.55	0.63

Table 34: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	<b>0.76</b>	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	<b>0.88</b>	<b>0.86</b>	–	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	<b>0.88</b>	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	–	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	–	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	–	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	–	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
8	4	0.41	0.73	0.50	0.83	0.76	0.45	0.70	0.62	0.24	0.50	0.75	0.53	0.59
8	8	0.44	0.76	0.58	0.86	0.67	0.50	0.68	0.62	0.32	0.59	0.79	0.56	0.61
8	16	0.44	0.75	0.59	0.85	0.79	0.50	0.72	0.66	0.34	0.58	0.78	0.55	0.63
16	4	0.41	0.71	0.50	0.83	0.75	0.45	0.70	0.62	0.24	0.50	0.75	0.53	0.58
16	8	0.43	0.67	0.58	0.86	0.68	0.50	0.69	0.62	0.33	0.59	0.79	0.56	0.61
16	16	0.45	0.68	0.58	0.85	0.69	0.50	0.69	0.62	0.34	0.58	0.78	0.55	0.61
32	4	0.42	0.72	0.50	0.83	0.75	0.45	0.70	0.62	0.24	0.51	0.75	0.53	0.58
32	8	0.44	0.67	0.58	0.85	0.69	0.49	0.69	0.63	0.33	0.60	0.80	0.55	0.61
32	16	0.45	0.68	0.58	0.85	0.69	0.50	0.69	0.63	0.34	0.58	0.78	0.55	0.61
64	4	0.42	0.73	0.51	0.83	0.76	0.45	0.70	0.62	0.25	0.51	0.76	0.54	0.59
64	8	0.44	0.76	0.58	0.86	0.69	0.50	0.71	0.63	0.33	0.58	0.80	0.55	0.62
64	16	0.45	0.68	0.58	0.86	0.68	0.50	0.70	0.63	0.34	0.58	0.79	0.55	0.61
128	4	0.42	0.74	0.51	0.83	0.76	0.45	0.71	0.62	0.25	0.52	0.76	0.54	0.59
128	8	0.45	0.71	0.59	0.86	0.72	0.49	0.73	0.65	0.33	0.58	0.80	0.56	0.62
128	16	0.46	0.69	0.58	0.86	0.69	0.50	0.71	0.64	0.34	0.58	0.79	0.55	0.62

Table 35: Detailed performance of various BACTRIAN-X-22 finetuned LLAMA-2-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	<b>0.76</b>	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	<b>0.88</b>	<b>0.86</b>	—	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	<b>0.88</b>	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	—	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	—	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	—	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	—	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
8	4	0.41	0.73	0.50	0.83	0.76	0.45	0.70	0.62	0.24	0.50	0.75	0.53	0.59
8	8	0.45	0.76	0.58	0.86	0.69	0.50	0.69	0.63	0.33	0.59	0.80	0.55	0.62
8	16	0.45	0.75	0.58	0.86	0.79	0.50	0.72	0.67	0.34	0.58	0.78	0.55	0.63
16	4	0.41	0.73	0.50	0.83	0.76	0.45	0.70	0.62	0.24	0.50	0.75	0.53	0.59
16	8	0.44	0.76	0.58	0.86	0.70	0.49	0.70	0.63	0.33	0.59	0.80	0.55	0.62
16	16	0.44	0.75	0.58	0.85	0.79	0.50	0.72	0.67	0.34	0.58	0.78	0.55	0.63
32	4	0.41	0.74	0.50	0.83	0.76	0.45	0.70	0.62	0.24	0.51	0.75	0.53	0.59
32	8	0.44	0.77	0.58	0.85	0.68	0.51	0.68	0.62	0.33	0.59	0.79	0.54	0.61
32	16	0.45	0.75	0.58	0.86	0.79	0.50	0.73	0.67	0.34	0.58	0.79	0.55	0.63
64	4	0.41	0.74	0.51	0.83	0.76	0.45	0.71	0.62	0.24	0.51	0.76	0.53	0.59
64	8	0.45	0.77	0.58	0.86	0.77	0.51	0.74	0.65	0.35	0.58	0.80	0.56	0.63
64	16	0.45	0.75	0.58	0.86	0.80	0.50	0.74	0.67	0.34	0.58	0.78	0.56	0.64
128	4	0.41	0.74	0.50	0.83	0.76	0.45	0.70	0.62	0.24	0.50	0.75	0.53	0.59
128	8	0.43	0.76	0.59	0.86	0.69	0.49	0.70	0.62	0.33	0.59	0.80	0.55	0.62
128	16	0.45	0.75	0.58	0.86	0.79	0.50	0.73	0.67	0.34	0.58	0.78	0.55	0.63

Table 36: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	0.76	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	0.88	<b>0.86</b>	—	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	<b>0.88</b>	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	—	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	—	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	—	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	—	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
64	8	0.72	0.84	0.76	0.91	0.89	0.69	<b>0.86</b>	0.77	0.62	0.72	0.87	0.73	0.78
64	16	0.72	0.84	0.76	0.91	0.88	0.70	0.86	0.77	0.64	0.72	0.88	0.73	0.78
128	8	0.71	0.83	0.76	0.91	0.88	0.69	0.86	0.77	0.62	0.71	0.88	0.72	0.78
128	16	0.72	0.84	0.77	0.91	<b>0.89</b>	0.70	0.86	0.77	0.63	0.71	0.88	0.73	0.78

Table 37: Detailed performance of various ALPACA finetuned MISTRAL-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	0.76	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	0.88	<b>0.86</b>	–	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	0.88	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	–	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	–	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	–	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	–	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
8	4	0.66	0.82	0.72	0.89	0.87	0.66	0.82	0.76	0.57	0.68	0.82	0.71	0.75
8	8	0.72	0.84	0.76	0.90	0.88	0.69	0.86	0.76	0.63	0.71	0.88	0.72	0.78
8	16	0.72	0.84	0.77	0.91	<b>0.89</b>	0.70	0.86	0.77	0.62	0.71	0.88	0.73	0.78
16	4	0.66	0.82	0.71	0.89	0.87	0.66	0.82	0.75	0.58	0.68	0.82	0.70	0.75
16	8	0.72	0.84	0.76	0.91	0.88	0.70	<b>0.86</b>	0.77	0.62	0.71	0.87	0.73	0.78
16	16	0.71	0.84	0.78	0.91	<b>0.89</b>	0.71	0.85	0.77	0.62	0.71	0.88	0.73	0.78
32	4	0.65	0.82	0.71	0.89	0.87	0.66	0.82	0.75	0.57	0.67	0.81	0.70	0.74
32	8	0.71	0.84	0.76	0.91	0.88	0.69	0.86	0.76	0.62	0.71	0.86	0.73	0.78
32	16	0.73	0.84	0.78	0.91	0.89	0.70	0.85	0.77	0.62	0.71	0.88	0.73	0.78
64	4	0.65	0.82	0.71	0.89	0.87	0.66	0.83	0.76	0.58	0.68	0.82	0.71	0.75
64	8	0.71	0.84	0.76	0.91	0.88	0.69	0.85	0.76	0.62	0.71	0.87	0.71	0.78
64	16	0.73	0.84	0.78	0.91	0.88	0.71	0.86	0.76	0.63	0.72	<b>0.88</b>	0.73	0.79
128	4	0.66	0.83	0.72	0.89	0.88	0.66	0.83	0.76	0.59	0.69	0.82	0.71	0.75
128	8	0.70	0.84	0.76	0.91	0.88	0.69	0.86	0.76	0.62	0.72	0.87	0.72	0.78
128	16	0.73	0.84	0.78	0.92	0.89	0.71	<b>0.86</b>	0.77	0.63	0.72	0.88	0.73	0.79

Table 38: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg	
LLaMA-7B-Chat	0.02	0.12	0.02	0.19	0.14	0.01	0.10	0.46	0.01	0.05	0.07	0.07	0.10	
LLaMA-3B-Chat	0.02	0.11	0.02	0.17	0.13	0.01	0.09	0.04	0.01	0.01	0.07	0.06	0.06	
LLaMA-70B-Chat	0.02	0.11	0.02	0.20	0.14	0.06	0.09	0.05	0.00	0.05	0.07	0.07	0.07	
Mistral-7B-Instruct	0.18	0.28	0.11	0.34	0.31	0.18	0.32	0.27	0.12	0.23	0.31	0.14	0.23	
DV003	0.37	0.55	0.32	0.77	0.62	0.20	0.58	0.29	0.12	0.45	0.42	0.36	0.42	
GPT-3.5-Turbo	0.60	0.71	0.49	0.79	0.70	0.54	0.70	0.58	0.42	0.62	0.69	0.50	0.61	
GPT-4	0.68	0.72	0.62	0.83	0.77	0.64	0.76	0.64	0.55	0.71	0.76	0.60	0.69	
TULRv6	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	0.90	0.88	<b>0.86</b>	–	<b>0.87</b>	<b>0.87</b>	<b>0.84</b>	0.88	0.79	<b>0.86</b>	
BLOOMZ	0.83	0.76	0.50	<b>0.92</b>	0.87	0.83	0.71	0.66	0.21	0.51	0.87	<b>0.82</b>	0.71	
XLM-R	0.69	0.80	0.80	0.86	0.82	0.77	–	0.80	0.74	0.76	0.79	0.59	0.77	
mBERT	0.61	0.71	0.63	0.83	0.76	0.59	–	0.71	0.43	0.55	0.69	0.58	0.65	
PALM2	0.63	0.76	0.77	0.87	0.81	0.62	–	0.70	0.62	0.68	0.73	0.49	0.70	
mT5	0.64	0.74	0.60	0.85	0.75	0.60	–	0.58	0.58	0.68	0.71	0.66	0.67	
rank	quantisation	ar	de	el	en	es	hi	ro	ru	th	tr	vi	zh	avg
8	4	0.66	0.82	0.72	0.89	0.87	0.66	0.82	0.74	0.57	0.67	0.82	0.71	0.75
8	8	0.71	0.83	0.76	0.91	0.88	0.70	0.86	0.76	0.62	0.70	0.88	0.72	0.78
8	16	0.73	0.84	0.77	0.91	0.89	0.70	0.86	0.76	0.62	0.71	0.88	0.73	0.78
16	4	0.66	0.82	0.72	0.89	0.87	0.66	0.83	0.75	0.57	0.67	0.82	0.71	0.75
16	8	0.71	0.84	0.76	0.91	0.88	0.69	0.86	0.76	0.62	0.70	0.87	0.73	0.78
16	16	0.73	0.84	0.77	0.91	0.88	0.70	0.86	0.77	0.63	0.71	0.88	0.73	0.78
32	4	0.66	0.82	0.72	0.89	0.87	0.67	0.82	0.75	0.58	0.67	0.83	0.71	0.75
32	8	0.72	0.84	0.76	0.91	0.89	0.70	0.86	0.76	0.62	0.72	0.88	0.72	0.78
32	16	0.72	0.84	0.77	0.91	0.88	0.71	0.85	0.76	0.63	0.72	0.88	0.72	0.78
32	32	0.72	0.84	0.77	0.91	0.88	0.71	0.85	0.76	0.63	0.72	0.88	0.72	0.78
64	4	0.66	0.82	0.71	0.89	0.88	0.68	0.83	0.74	0.58	0.67	0.83	0.71	0.75
64	8	0.72	0.84	0.76	0.91	0.88	0.70	0.86	0.76	0.62	0.71	0.87	0.73	0.78
64	16	0.71	0.84	0.77	0.92	0.89	0.71	<b>0.86</b>	0.76	0.63	0.72	<b>0.88</b>	0.73	0.78
128	4	0.69	0.82	0.73	0.89	0.88	0.68	0.83	0.75	0.58	0.68	0.84	0.71	0.76
128	8	0.72	0.84	0.76	0.91	0.88	0.70	0.86	0.76	0.63	0.70	0.87	0.73	0.78
128	16	0.72	0.84	0.77	0.91	<b>0.89</b>	0.71	0.86	0.76	0.62	0.71	0.87	0.72	0.78

Table 39: Detailed performance of various MULTIALPACA finetuned MISTRAL-7B models on XQuAD ([Artetxe et al., 2020](#)).

Model	ar	en	es	fr	hi	jp	zh	avg	
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.20	0.08	
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.22	0.24	
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
64	8	0.00	0.03	0.14	0.05	0.01	0.08	0.05	0.05
64	16	0.02	0.07	0.15	0.09	0.02	0.10	0.06	0.07
128	8	0.01	0.04	0.15	0.06	0.02	0.09	0.06	0.06
128	16	0.02	0.07	0.14	0.09	0.02	0.10	0.06	0.07

Table 40: Detailed performance of various ALPACA finetuned LLAMA-2-7B models on XLSum ([Hasan et al., 2021](#)).

Model		ar	en	es	fr	hi	jp	zh	avg
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.08
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.26	0.22	0.24
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
8	4	0.00	0.03	0.15	0.05	0.00	0.13	0.05	0.06
8	8	0.02	0.04	0.15	0.06	0.02	0.09	0.06	0.06
8	16	0.02	0.26	0.14	0.08	0.02	0.10	0.06	0.10
16	4	0.00	0.03	0.15	0.05	0.00	0.13	0.05	0.06
16	8	0.02	0.04	0.15	0.06	0.02	0.09	0.05	0.06
16	16	0.02	0.26	0.14	0.08	0.02	0.10	0.06	0.10
32	4	0.00	0.03	0.15	0.05	0.00	0.11	0.05	0.06
32	8	0.01	0.04	0.15	0.06	0.01	0.10	0.06	0.06
32	16	0.02	0.26	0.14	0.08	0.02	0.10	0.06	0.10
64	4	0.00	0.04	0.15	0.05	0.00	0.12	0.05	0.06
64	8	0.01	0.04	0.15	0.06	0.01	0.09	0.06	0.06
64	16	0.02	0.26	0.14	0.07	0.02	0.11	0.06	0.10
128	4	0.00	0.03	0.15	0.05	0.00	0.13	0.05	0.06
128	8	0.01	0.04	0.15	0.06	0.02	0.08	0.06	0.06
128	16	0.02	0.26	0.14	0.08	0.02	0.10	0.06	0.10

Table 41: Detailed performance of various MULTIALPACA finetuned LLAMA-2-7B models on XLSum (Hasan et al., 2021).

Model		ar	en	es	fr	hi	jp	zh	avg
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.08
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.26	0.22	0.24
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
8	4	0.00	0.03	0.15	0.05	0.00	0.13	0.05	0.06
8	8	0.02	0.04	0.15	0.05	0.01	0.09	0.06	0.06
8	16	0.02	0.07	0.14	0.09	0.02	0.10	0.06	0.07
16	4	0.00	0.03	0.15	0.05	0.00	0.13	0.05	0.06
16	8	0.01	0.04	0.15	0.05	0.02	0.10	0.06	0.06
16	16	0.02	0.07	0.14	0.09	0.02	0.10	0.06	0.07
32	4	0.00	0.03	0.15	0.05	0.00	0.11	0.05	0.06
32	8	0.02	0.04	0.15	0.06	0.01	0.09	0.06	0.06
32	16	0.02	0.07	0.15	0.09	0.02	0.10	0.06	0.07
64	4	0.00	0.04	0.15	0.05	0.00	0.12	0.05	0.06
64	8	0.02	0.04	0.15	0.05	0.01	0.09	0.07	0.06
64	16	0.02	0.07	0.14	0.10	0.02	0.10	0.06	0.07
128	4	0.00	0.04	0.15	0.05	0.00	0.13	0.05	0.06
128	8	0.02	0.04	0.15	0.06	0.01	0.10	0.07	0.06
128	16	0.02	0.07	0.14	0.10	0.02	0.10	0.06	0.07

Table 42: Detailed performance of various BACTRIAN-X-22 finetuned LLaMA-2-7B models on XLSum (Hasan et al., 2021).

Model		ar	en	es	fr	hi	jp	zh	avg
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.08
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.26	0.22	0.24
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
64	8	0.01	0.29	0.10	0.24	0.00	0.02	0.05	0.10
64	16	0.01	0.29	0.10	0.25	0.00	0.01	0.05	0.10
128	8	0.01	0.29	0.09	0.25	0.00	0.02	0.05	0.10
128	16	0.00	0.29	0.11	0.25	0.00	0.01	0.05	0.10

Table 43: Detailed performance of various ALPACA finetuned MISTRAL-7B models on XLSum (Hasan et al., 2021).

Model		ar	en	es	fr	hi	jp	zh	avg
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.08
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.26	0.22	0.24
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
8	4	0.00	0.30	0.08	0.24	0.00	0.00	0.06	0.10
8	8	0.00	0.29	0.08	0.25	0.00	0.02	0.06	0.10
8	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
16	4	0.00	0.30	0.07	0.24	0.00	0.00	0.06	0.10
16	8	0.00	0.29	0.08	0.24	0.00	0.02	0.06	0.10
16	16	0.00	0.29	0.08	0.25	0.00	0.02	0.05	0.10
32	4	0.00	0.30	0.08	0.23	0.00	0.00	0.06	0.10
32	8	0.00	0.29	0.08	0.26	0.00	0.01	0.05	0.10
32	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
64	4	0.00	0.30	0.08	0.23	0.00	0.00	0.06	0.10
64	8	0.00	0.30	0.08	0.25	0.00	0.02	0.05	0.10
64	16	0.00	0.28	0.10	0.25	0.00	0.02	0.05	0.10
128	4	0.00	0.30	0.08	0.23	0.00	0.01	0.06	0.10
128	8	0.00	0.29	0.09	0.25	0.00	0.01	0.06	0.10
128	16	0.00	0.29	0.08	0.25	0.00	0.02	0.05	0.10

Table 44: Detailed performance of various MULTIALPACA finetuned MISTRAL-7B models on XLSum ([Hasan et al., 2021](#)).

Model		ar	en	es	fr	hi	jp	zh	avg
LLaMA-70B-Chat	<b>0.35</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.08
GPT-3.5-Turbo	0.25	0.26	0.21	0.26	0.24	0.26	0.26	0.22	0.24
GPT-4	0.23	0.27	0.20	0.23	0.24	0.29	0.31	0.25	
mT5	0.31	<b>0.35</b>	<b>0.27</b>	0.23	<b>0.34</b>	<b>0.39</b>	<b>0.40</b>	<b>0.33</b>	
PALM2	0.06	0.00	0.00	<b>0.31</b>	0.03	0.00	–	0.07	
rank	quantisation	ar	en	es	fr	hi	jp	zh	avg
8	4	0.00	0.30	0.07	0.24	0.00	0.00	0.06	0.10
8	8	0.00	0.29	0.09	0.24	0.00	0.02	0.05	0.10
8	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
16	4	0.00	0.29	0.07	0.24	0.00	0.00	0.06	0.10
16	8	0.00	0.29	0.09	0.24	0.00	0.02	0.05	0.10
16	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
32	4	0.00	0.30	0.07	0.23	0.00	0.00	0.06	0.10
32	8	0.00	0.29	0.09	0.24	0.00	0.01	0.05	0.10
32	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
64	4	0.00	0.30	0.07	0.24	0.00	0.01	0.06	0.10
64	8	0.00	0.29	0.08	0.25	0.00	0.02	0.06	0.10
64	16	0.00	0.29	0.09	0.25	0.00	0.02	0.05	0.10
128	4	0.00	0.30	0.08	0.23	0.00	0.01	0.05	0.10
128	8	0.01	0.29	0.09	0.25	0.00	0.02	0.05	0.10
128	16	0.00	0.29	0.10	0.25	0.00	0.02	0.05	0.10

Table 45: Detailed performance of various BACTRIAN-X-22 finetuned MISTRAL-7B models on XLSum (Hasan et al., 2021).