AceGPT, Localizing Large Language Models in Arabic

Anonymous ACL submission

Abstract

This paper is devoted to the development of a localized Large Language Model (LLM) specifically for Arabic, a language imbued with unique cultural characteristics inadequately addressed by current mainstream models. Sig-006 nificant concerns emerge when addressing cultural sensitivity and local values. To address this, the paper proposes a comprehensive solution that includes further pre-training with Arabic texts, Supervised Fine-Tuning (SFT) utilizing native Arabic instructions, and GPT-4 responses in Arabic, alongside Reinforcement Learning with AI Feedback (RLAIF) employing a reward model attuned to local culture and values. The goal is to cultivate culturally cog-016 nizant and value-aligned Arabic LLMs capable of accommodating the diverse, application-017 specific needs of Arabic-speaking communities. Comprehensive evaluations reveal that the resulting model, dubbed 'AceGPT,' sets the state-of-the-art standard for open Arabic LLMs across various benchmarks, including 022 the instruction-following benchmark (i.e., Ara-024 bic Vicuna-80 and Arabic AlpacaEval), knowledge benchmark (i.e., Arabic MMLU and EX-AMs), and the newly introduced Arabic Cultural and Value Alignment benchmark. No-027 tably, AceGPT outperforms Turbo in the popular Vicuna-80 benchmark when evaluated with GPT-4, despite the benchmark's limited scale.

1 Introduction

037

041

LLMs like Turbo and GPT-4 have been shaping the current landscape of natural language understanding and generation ((Bubeck et al., 2023)). In contrast to the proprietary nature of Turbo and GPT-4, there has been a trend towards developing open-source large language models capable of instruction-following (Taori et al., 2023) and fluent conversations ((Chiang et al., 2023)), a phenomenon termed as 'Democratization of Chat-GPT' ((Conover et al., 2023; Touvron et al., 2023)). While these models have shown great promise in understanding and producing content in various languages, they might fail to align with local values and cultural norms in non-English environments ((Chen et al., 2023a)); we call it the 'localization issue'. This issue can lead to significant problems in practical usage scenarios, especially for regions such as the Arabic world where the culture and values diverge significantly from Western norms. We argue that it is not just desirable but necessary to localize large language models and tailor them to a specific cultural environment. 042

043

044

047

048

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

079

Methodology The core of our approach lies in localizing large language models to the Arabic language using a packaged solution (known as AceGPT). Firstly, through incremental pre-training on Arabic data (localized pre-training), we ensure that the model has a strong foundation in the Arabic language, including grammar, vocabulary, and cultural context. Next, by fine-tuning Arabic natural questions (localized instructions), we enable the model to effectively comprehend and respond to specific instructions that are pertinent to Arab interests. Furthermore, by generating Arabic native responses from GPT-4 (localized responses) rather than relying on translations from other languages, we ensure that the model's outputs are natural and fluent within an Arabic context thanks to the powerful GPT-4. Lastly, by employing a reward model based on *localized preference data* that respects local culture and value, we further refine the model to align the responses with the cultural and value norms of Arabic-speaking communities.

Evaluation We evaluate our models in various benchmarks: in the **instruction-following benchmark**, AceGPT achieves state-of-the-art (SOTA) among open-sourced Arabic LLMs in Arabic Vicuna-80 and Arabic AlpacaEval, obtaining 33% and 30% improvement over the state-of-the-art Ara-

Types of entity	Jais- 13B	Turbo	GPT- 4	AceGPT (ours)
Person	12.00%	26.67%	39.29%	50.00%
	$(3/25)^1$	(12/45)	(22/56)	(31/62)
Location	18.75%	27.08%	21.62%	28.95%
	(3/16)	(13/48)	(16/74)	(11/38)

Table 1: Proportion of Arabic Entities in Responses to20 Sample Arabic Questions

¹ 25 person names in Jais-13B responses are identified and 3 are Arabic names.

bic LLM ((Sengupta et al., 2023)). ¹ In the NLU
benchmark, AceGPT achieves the second best on
ALUE ((Seelawi et al., 2021)) in terms of average
scores for all tasks. In the knowledge benchmark,
AceGPT achieves SOTA among open-sourced Arabic LLMs in Arabic knowledge including MMLU
and EXAMs. In the localization benchmark,
AceGPT achieves SOTA among open-source Arabic LLMs in our Arabic Cultural and Value Alignment (ACVA) Dataset.

Contributions The contributions of the paper are three-fold, including **i**) we propose a first-tier Arabic LLM. According to the records on the releasing date, it achieves SOTA performance among open Arabic LLMs in many benchmarks including Arabic Vicuna-80, Arabic AlpacaEval, Arabic MMLU, EXAMs, and ACVA. **ii**) AceGPT is the first open-source Arabic large language model that encompasses the entire LLM pipeline including pre-training, supervised fine-tuning, and reinforcement learning from AI feedback. We release AceGPT and the reward model. **iii**) We observe and measure the localization issue in LLMs quantitatively and have introduced a new benchmarking dataset, ACVA, for localization testing.

2 Recipe of AceGPT

2.1 Motivation: the Localization Issue

Given the availability of many high-quality instruction datasets in widely spoken languages such as English, existing strategies for non-English LLMs often rely on instructions translated from English. Examples include Chinese-alpaca-GPT4 ((Peng et al., 2023)), Phoenix ((Chen et al., 2023b)), and Jais ((Sengupta et al., 2023)). However, relying on translated data may lead to *localization issues*, potentially undermining the integrity and applicability of the models in native contexts.

To address these localization issues, we formulate 20 questions (see Table.12) to elicit responses with name entities-both personal and locational-to summarize the prevalence of Arabic name entities for preliminary experiments. Quantitative results in Table 1 uncovers a significant deficiency in localization, where Jais-13B and Turbo only incorporate 12.00% and 26.67% Arabic names out of all the names in their responses respectively. A specific example is shown in Table 2, we can observe that the Arabic open-source LLM Jais's output shows a conspicuous tilt towards Englishcentric materials, yielding terms predominantly associated with Christianity, which potentially neglects significant parallels within Arabic literary traditions. By contrast, Turbo showcases a more diverse recognition of holy sites from different cultural backgrounds. You can see the details and more examples of case studies in Appendix A.2.

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

164

166

2.2 Methodology of AceGPT

To address localization, we propose a comprehensive solution including three strategies to ensure model's effective understanding and generation of content in Arabic, with cultural awareness and value alignment: **(I) localized pre-training** we further pre-train LLM with Arabic data; **(II) localized instructions** we adopt Arabic natural questions in the wild and their responses are Arabic native responses from GPT-4 instead of translating that from other languages, and **(III) localized feedback** we further tame LLM with reinforcement learning using a reward model that respects local culture and values thanks to the localized preference data.

The resultant model is termed "AceGPT". The model pre-trained on LLaMA2 ((Touvron et al., 2023)) is named "AceGPT-*base*". To equip it with the conversation, we introduced "AceGPT-*chat*" utilizing supervised fine-tuning and reinforcement learning from AI feedback. The training procedure is divided into three stages: pre-training, supervised fine-tuning, and reinforcement learning from AI feedback, introduced in Sec 2.2.1, Sec 2.2.2, and Sec 2.2.3, respectively.

2.2.1 Localized Pre-training

To adapt the English-focused LLaMA2 ((Touvron et al., 2023)) model in Arabic, we train further it with a substantial corpus of Arabic text.

Data The dataset comprises Arabic and English sub-datasets. The Arabic is derived from the

099

100

102

103

104

105

106

108

109

110

111

112

113

114

115

116

¹Jais ((Sengupta et al., 2023)) is a concurrent work released two weeks ahead of ours.

Table 2: Comparison of LLM Responses Highlighting Location, Person, and Regional/Cultural Terms. Arabic-rooted terms in blue; non-Arabic terms in red.

User:

168

169

170

171

172

173

174

175 176

177

178

179

180

181

182

184

185

186

189

190

191

192

193

194

195

(What are the holy books, saints, and holy places?)

ما هي الكتب المقدسة والقديسين والأماكن المقدسة؟

open-source Arabic text 2022², and refined from sources like Arabic Wikipedia, CC100, and OS-CAR3. The English dataset is obtained from Slim Pajama ((Soboleva et al., 2023)) to avoid forgetting knowledge of English texts. Given LLaMA2's excellent adaptability to the English dataset, we sample a subset of data from Slim Pajama randomly.

Due to the limit of computing resources, we only train the *LLaMA2-7B* with 30B data (19.2B tokens in Arabic and 10.8B in English) and *LLaMA2-13B* with 10B data (6B tokens in Arabic and 4B in English), prioritizing a larger quantity of Arabic data than English data. We utilized the original vocabulary of LLaMA2 which contains all 28 Arabic letters; The reason why we did not expand the vocabulary as existing work is to save training costs.

2.2.2 Localized Supervised Fine-Tuning

To enable the model to follow Arabic user instructions and tackle realistic applications, we finetuned AceGPT with **localized instructions** and **localized responses**.

Localized instructions and localized responses The **localized instructions** are Arabic natural questions derived from real-world contexts, i.e. online question-answering platforms Quora ³, which can help models to capture what Arabs care in the wild. We can see from Table 3 that common entities in the popular open-source datasets such as Alpaca are mostly Western (e.g. "John", "Apple", and "New York"), deviating from Arab's actual interest (e.g. "Mohammed", "Muslim Brotherhood", and "Egypt") which can be addressed by Quora. The main idea of **localized responses** is to leverage the fact that GPT-4 produces cultureand value-relevant responses in the context of question language, which means responses to questions in English are different from those in Arabic. See an example in Table 4, GPT-4 produces culture-dependent responses based on the queried languages. Therefore, when incorporating opensource instruction-tuning data, we ask the GPT-4 to re-generate responses in Arabic (rather than translate) to produce localized responses.

Data In addition to Arabic Quora questions, we also incorporate some open-source instructiontuning datasets to improve the overall performance. Specifically, we incorporate **Alpaca** (Taori et al., 2023; Peng et al., 2023) (the most classical instruction-tuning dataset), **Evol-Instruct** (Xu et al., 2023) (a complex instruction dataset), **Code-Alpaca** (Chaudhary, 2023) (a code-specific instruction dataset) ⁴, and **ShareGPT** ⁵ (a popular user-GPT dialogue dataset). For these open-source data except ShareGPT, an Arabic version is created by translating the English questions into Arabic and regenerating the responses using GPT-4. We reserve the original ShareGPT data because the original conversations will be destroyed with a re-generated

²https://data.baai.ac.cn/details/ArabicText-

²⁰²² provided by BAAI

³https://quora.com/

⁴We incorporate code-alpaca for a more powerful LLM with a better code capability.

⁵https://huggingface.co/datasets/philschmid/ sharegpt-raw

Dataset	Top-5 Person	Top-5 Organization	Top-5 GPE
Alpaca	John, John Smith, Alice, Mary, Harry Potter	Apple, Amazon, Google, Microsoft, ABC	United States, India, New York, France, China
Evol-Instruct	John, John Smith, Harry Potter, Alice, Bob	Apple, Amazon, quantum, Google, Microsoft	United States, New York, Los Angeles, San Francisco, Japan
ShareGPT	Di Maria, Messi, Beckhaus, Eco, Clara	Tribunal, Google, Council, Bing, Supreme Court	United States, Argentina, France, New York, Hong Kong
Quora	Prophet, Mohammed, Adam, Hijri, Ali	European Union, Google Muslim Brotherhood, Soviet Union, United Nations	Egypt, Turkey, Saudi Arabia, Morocco, America

Table 3: Top 5 names of individuals, organizations, and geopolitical entities (GPE) by frequency.

Table 4: GPT-4 answers culture-relevant questions differently across languages. Questions here are the same in semantics but differ in languages. The Arabic response is translated into English (right).

Question in English :	Question in Arabic:
Is it forbidden for a woman to confess her love to a man?	هل من الحرام أن تعترف المرأة بحبها للرجل؟
GPT-4 : No, it is not forbidden for a woman to confess her love to a man. Society and social norms have evolved greatly over the years, and it has become increasingly common and accepted for women to express their feelings and intentions toward men, just as it is for men to do so toward women	GPT-4 : (In Islam, it is not necessarily considered haram for a woman to express her feelings towards a man, but it must be done decently and respectfully and in accordance with Islamic teachings. Islam encourages modesty and chastity, and therefore both men and women must display decency and good morals in their dealings. If a woman wants to express her feelings, she can do so directly or through an intermediary, such as her family)

different response.

225

236

240

241

243

245

246

247

249

2.2.3 Reinforcement Learning from AI feedback

To further align AceGPT with values and cultures, we utilize reinforcement learning from AI feedback with a reward model trained with **localized pref**erence data. There are primarily two stages: (1) training the reward model using localized preference data, and (2) aligning AceGPT to value and culture preference patterns using the proximal policy optimization algorithm (Schulman et al., 2017).

Localized preference data To align AceGPT with Arabic culture and values, a reward model mimicking the preferences of native speakers is essential. To prepare the localized preference data for reward model training, we reuse 40K localized instructions, i.e. Quora questions, in the SFT stage and sample paired outputs from our fine-tuned 7B model. Given the resource-intensive nature of collecting human feedback, we utilized GPT-4 feedback, which has been shown to correlate highly with human preference labeling and achieves competitive performance in text summarization (Lee et al., 2023). However, due to observed position bias in GPT-4 (Zhang et al., 2023), we altered the order of sample answers and retained consistent preferences between two order-switched runs, resulting in 12K pairs. A small study with 800 examples verified the reliability of this preference data, revealing a correlation coefficient of 0.84 between GPT-4 and human evaluations. We also incorporate 12K open-source preference data for better generalization. See Appendix C for details.

Reward model The reward model operates within a 'binary' framework, determining preferences with an additional linear head post the final hidden states. The loss function is expressed as:

$$\mathcal{L}(\theta) = -\frac{1}{\|D\|} \mathbb{E}_{(x,y_c,y_r)\sim D} \left[\log(\sigma(r_\theta(x,y_c) - r_\theta(x,y_r))) \right].$$
(1)

Here, x is the input, y_c is the chosen model output, y_r is the rejected model output of the pair, and r_{θ} is the reward model with the parameter θ .

Proximal policy optimization We crawl another 30K Quora questions different from Quora-40K for PPO training data. Proximal Policy Optimization (PPO) is an off-policy policy gradient method for reinforcement learning (Schulman et al., 2017). The policy $\pi_{\theta}(a|s)$ represents the probability distribution over the next token *a* given a sequence of previous tokens *s*, where θ are the model pa258

259

rameters. The primary objective is to maximize the preference signal from the reward model that corresponds to the desired output behaviour. The objective is

$$\mathcal{L}(\theta) = \mathbb{E}_t \left[\min \left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} A_t, \operatorname{clip} \left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}, 1 - \epsilon, 1 + \epsilon \right) A_t \right) \right]$$
(2)

Here, θ is the current model parameter while θ_{old} is the model parameter used for experience sampling. A_t is the advantage function that measures the relative value of generating a_t as the next token conditioned on the sequence $s_1 \cdots s_t$, and ϵ is a hyperparameter for stability.

3 Localized Evaluation

3.1 Evaluation protocol

261

262

264

267

269

270

271

272

274

277

278

281

289

290

291

293

294

297

298

301

In this section, we delve into the 'Localized Evaluation' of our language model, focusing exclusively on the Arabic Cultural and Value Alignment (ACVA). ACVA serves as a critical benchmark to assess our model's performance in terms of its alignment with Arabic cultural nuances and values. This evaluation is particularly important for understanding how well the model adapts to the specific linguistic and cultural context of the Arabic language. We conduct this evaluation using our fine-tuned *chat* models, which have been specifically optimized for higher relevance and accuracy in culturally specific scenarios.

Arabic Cultural and Value Alignment (ACVA) ACVA is a Yes-No question dataset, comprising over 8000 questions, generated by Turbo from 50 designed Arabic topics to assess model alignment with Arabic values and cultures (see Appendix B for data construction details). A subset, revised by Arabic speakers for question quality and answer accuracy, forms the 2486-data 'Clean set'. The correlation between 'All set' and 'Clean set' evaluations is in Sec 3.2. Given our focus on localized solutions, we evaluate our final models (post-SFT and RLAIF) on this benchmark in a zero-shot setting with F1 score.

Baselines We compare the performance of our models against LLaMA2 (Touvron et al., 2023), Bloomz (Muennighoff et al., 2022), Phoenix (Chen et al., 2023a,b), and Jais (Sengupta et al., 2023) for this section and Sec. 4.1. LLaMA2-chat models are excluded as they consistently respond in English when queried in Arabic. See details in Sec. E.1.

Table 5: Average F1 on **ACVA** in the zero-shot setting. The best performance is in **bold** and the second is underlined.

Model	All set	Clean set
Phoenix	41.86%	43.80%
Phoenix–multiple-langs	59.78%	59.15%
Jais-13B- <i>chat</i>	61.44%	66.83%
AceGPT-7B-chat	69.53%	70.03%
AceGPT-13B-chat	<u>75.02%</u>	<u>74.62%</u>
Turbo	75.57%	79.03%

Table 6: Average performance ratio of Turbo and the standard variation over three runs in **Arabic Vicuna-80** and **Arabic AlpacaEval**. The best performance is in **bold** and the second is underlined.

Comparison	Arabic Vicuna-80	Arabic AlpacaEval
Phoenix Phoenix-multiple-langs Jais-13B-chat	$71.92\% \\ \pm 0.2\% \\ 71.67\% \\ \pm 0.7\% \\ 75.40\% \\ \pm 1.6\% $	$65.62\% \\ \pm 0.3\% \\ 65.36\% \\ \pm 0.1\% \\ 74.95\% \\ \pm 0.2\% \\ $
AceGPT-7B-chat AceGPT-13B-chat		$\frac{93.81\%}{\pm 0.1\%}$ 97.95% $\pm 0.1\%$

3.2 Experiment results

ACVA benchmark We present the results of AceGPT and other chat models on ACVA in Table 5. The Pearson correlation of accuracy on 'All set' and 'Clean set' is 0.9825, indicating a high reliability of ACVA all set evaluation. Notably, our AceGPT-*chat* models (both 7B and 13B) consistently outperform other open-source LLMs, and AceGPT-13B-chat only trails Turbo by a marginal of -0.55% on all set. Since the Jais-30B-*chat* does not follow instructions and cannot return answers for multiple-choice questions, we suspect that this is due to overly stringent safety measures. Therefore, we did not include Jais-30B-*chat* in the zero-shot comparison.

4 Overall Evaluation

4.1 Evaluation protocol

Evaluation of language models is multifaceted319and typically involves multiple metrics and bench-
marks to assess various aspects of model perfor-
mance.Moving beyond the scope of localization,
the 'Overall Evaluation' section presents a com-
prehensive analysis of our language model across320

304

305

306

307

308

- 309 310 311 312 313
- 316 317

314

315

a spectrum of benchmarks. This includes assessing instruction-following ability, knowledge retention, and Natural Language Understanding (NLU). 327 For evaluating instruction-following ability, we employ our fine-tuned chat models, which are designed to excel in interactive and directive tasks. In 330 contrast, knowledge retention and NLU are evalu-331 ated using our base models, focusing on the core strengths of the model's pre-training. While we utilize both automated and manual methods for as-334 sessing instruction-following ability, other bench-335 marks in this section are evaluated solely through 336 automated methods. 337

341

347

361

368

372

376

Instruction-following We specifically evaluate the instruction-following capabilities of models tuned for instructions using Arabic Vicuna-80 and Arabic AlpacaEval. In accordance with (Chiang et al., 2023), we adopt the GPT-4 evaluation, which prompts GPT-4 to score the performance of models on each question, contrasting them with Turbo. The details can be found in Appendix E.2. While GPT-4 evaluation is efficient and scalable, it may overlook the subtle inconsistencies between model responses (Wang et al., 2023) and human interactions in real-world scenarios. Therefore, we further conduct human evaluation on these benchmarks to evaluate the performance of AceGPT from the perspective of human rather than GPT-4 preferences. To ensure cultural relevance in manual evaluations, we engaged a diverse group of educated, native Arabic speakers. Each model's response was assessed independently by three assessors. We present more details in Table 16 and the UI for evaluation in Figure 2.

Vicuna-80 (Chiang et al., 2023) is a popular benchmark containing 80 open-ended questions, distributed across eight categories. To attain a more reliable evaluation of instruction-following capabilities, we resort to a larger benchmark, AlpacaE-val (Dubois et al., 2023). This benchmark is structured to replicate the actual distribution of user instructions by consolidating several public datasets. It is reported that model rankings on this benchmark have a high correlation with those on the live user instructions. Arabic Vicuna-80 and Arabic AlpacaEval are translated from these benchmarks by GPT-4 and revised by native speakers.

Knowledge We have two knowledge benchmarks, including Arabic MMLU and EXAMs.
MMLU (Hendrycks et al., 2021) consists of diverse multiple-choice questions across 57 tasks, spanning various educational levels. We employed

Turbo to translate this dataset from English to Arabic. Additionally, Arabic questions from the **EX-AMs** (Hardalov et al., 2020), a resource specialized in multilingual high school exam questions, were also incorporated. Both datasets were evaluated in a few-shot setting, as per the methodology in (Huang et al., 2023), to assess the innate capabilities of LLMs, aiming at potential applications with minimal adaptations. 377

378

379

381

382

383

384

386

387

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

4.2 Experiment results

Instruction-Following benchmark We present each model's performance ratio against turbo, scored by GPT-4, in Table 6. The result shows that AceGPTs are superior in both Arabic Vicuna-80 and Arabic AlpacaEval. Notably, AceGPT-7Bchat surpasses Jais-13B by about 20% points with smaller model size. Moreover, AceGPT-13B-chat attains a 100.88% performance ratio of Turbo in Arabic Vicuna-80.

Human Evaluation Table 7 shows the human evaluation results on Arabic Vicuna-80 and Arabic AlpacaEval. We calculated the percentages of wins, ties, and losses of the results from three Arabic speakers. We observe that AceGPT-*chat* (both 7B and 13B) significantly surpasses Jais-13B-*chat* and even Jais-30B-*chat*, but lags behind Turbo. Moreover, the AceGPT-13B-*chat* is significantly better than the AceGPT-7B-*chat*, indicating the importance of model size.

Knowledge benchmark Table 8 shows the fewshot evaluation results on Arabic MMLU and EX-AMs. We can see that AceGPT-13B-base attains the best performance in Arabic MMLU (37.26%) among open-source LLMs, and AceGPT-7B-base also surpasses other open-source models, including 13B models, in Humanities and Others (Business, Health, Misc) domains in Arabic MMLU. In EX-AMs, the Jias-30B-base model achieves the best performance among open-source models.

5 Experimental Analysis

5.1 On Pre-training

Localization

of Pre-training AceGPT-base uses LLaMA2 as the backbone. the only difference it is further pre-trained

Dataset	Comparison	win	tie	lose win or tie
	AceGPT-7B-chat vs. Jais-13B-chat AceGPT-7B-chat vs. Turbo	82.5% 27.5%	6.7% 32.9%	10.8%89.2%39.6%60.4%
Arabic Vicuna-80	AceGPT-13B-chat vs. Jais-13B-chat AceGPT-13B-chat vs. Turbo	82.9% 16.3%	6.7% 57.1%	10.4%89.6%26.6%73.4%
	AceGPT-7B-chat vs. Jais-30B-chat AceGPT-13B-chat vs. Jais-30B-chat	67.5% 64.6%	15.0% 15.0%	17.5% 82.5% 20.4% 79.6%
Arabic AlpacaEval	AceGPT-7B-chat vs. Jais-13B-chat AceGPT-7B-chat vs. Turbo	53.0% 20.2%	36.5% 46.5%	10.5%89.5%33.3%66.7%
	AceGPT-13B-chat vs. Jais-13B-chat AceGPT-13B-chat vs. Turbo	49.4% 25.2%	42.8% 44.5%	7.8%92.2%30.3%69.7%

Table 7: Human evaluations on Vicuna-80 and AlpacaEval. The winners are in bold.

Table 8: Accuracy on Arabic MMLU and EXAMs. The best is bold and the second is underlined.

Model	Average	STEM	Arabic MN Humanities	ILU Social Sciences	Others	EXAMs
Bloomz	33.69	33.35	29.29	37.58	34.53	33.89
LLaMA2-7B	29.47	30.30	29.33	27.46	30.78	23.48
LLaMA2-13B	33.76	32.94	32.30	33.42	37.27	25.45
Jais-13B- <i>base</i>	32.23	30.51	31.25	33.74	33.43	35.67
Jais-30B- <i>base</i>	36.27	32.67	30.67	42.13	39.60	<u>39.91</u>
AceGPT-7B-base	32.14	29.73	30.95	33.45	34.42	31.96
AceGPT-13B-base	<u>40.45</u>	<u>36.60</u>	<u>38.74</u>	<u>43.76</u>	<u>42.72</u>	36.63

Table 10: Effects of different datasets on ACVA, Arabic Vicuna-80 and Arabic AlpacaEval.

Comparison	ACVA	Arabic Vicuna-80	Arabic AlpacaEval
Alpaca-Arabic	50.52%	87.15% ± 0.5%	82.97% ± 0.4%
+ ShareGPT	38.64%	88.01% ± 0.03%	84.89% ± 0.3%
+ Evol-Instruct	61.72%	90.39% ± 0.4%	86.87% ± 0.1%
+ Quora	65.53%	89.74% ± 0.8%	85.71% ± 0.03%

with some local Arabic texts. We compare AceGPT-base to LLaMA2 on ACVA with the few-shot setting to demonstrate the benefits of localized pre-training on Arabic culture and values. The results in Table 9 show the superiority of localized pre-training: after localized pre-training, AceGPT-7B-base surpasses LLaMA2-13B, which has a larger size.

5.2 On Supervised Fine-tuning

426

427

428

429

430

431

432

433

434

435

436

437

In this analysis, we primarily assess the impact of both localized and open-source instructions on localization and overall performance. Each dataset has been sampled with 40,000 data points, respectively. The results are shown in Table 10. It can be observed that Evol-Instruct highly contributes to the overall performance in the instructionfollowing benchmark, while Quora is most beneficial for Arabic culture and values. Note that incorporating ShareGPT largely harms the performance of ACVA; this may be because ShareGPT is almost aligned with Western culture and values.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

5.3 On RLAIF

5.3.1 Reward model

To evaluate the sensitivity of the reward model to the overall performance, we measure the correlations between reward scoring and GPT-4 scoring (described in section 4.1) on Arabic Vicuna-80. Following the pairwise comparison setting in GPT-4 scoring, we also calculate the performance ratio for normalized (to [0, 10] as GPT-4 scoring) reward scores on model-chatbot pairs. The Pearson correlation and Spearman correlation are 0.57 and 0.61 respectively, and the results are shown in Figure 1a. We conclude that the reward model shows a positive correlation with GPT-4 evaluation on Arabic Vicuna, which indicates it can offer an effective signal on overall performance.

		Automatic eva	luation	Huma	ın Evalu	ation (v	vs. Turbo)
Comparison	ACVA	Arabic Vicuna-80	Arabic AlpacaEval	win	tie	loss	win or tie
AceGPT-7B-chat (w/o RLAIF)	42.48%	92.01% ± 1.3%	91.35% ± 0.08%	27.5%	29.2%	43.3%	56.7%
AceGPT-7B-chat	69.53%	94.82% ± 0.2%	93.81% ± 0.1%	27.5%	32.9%	39.6%	60.4%
AceGPT-13B-chat (w/o RLAIF)	74.18%	95.14% ± 1.0%	93.05% ± 0.2%	19.6%	37.5%	42.9%	57.1%
AceGPT-13B-chat	75.02%	100.88% ± 0.4%	97.95% ± 0.1%	16.3%	57.1%	26.7%	73.3%

Table 11: Experiments with/without RLAIF on Arabic Vicuna-80, Arabic AlpacaEval and ACVA.



(a) Correlations between reward model scoring and GPT-4 scoring on Arabic Vicuna-80.



(b) Reward score distributionon ACVA.

Figure 1: (a) Correlations between the reward model and GPT-4 and (b) reward distribution.

Localization of Reward model Then we evaluate the Arabic culture sensitivity of the reward model on the ACVA benchmark. Prompting with "Give me a fact about Arab culture, values, and laws" in Arabic, we calculate the reward scores of prompt-statement pairs for all statements from ACVA. The distribution of reward scores for yes/no statements is shown in Figure 1b. It demonstrates that reward scores for "yes" statements are higher than "no" statements overall, which suggests that our reward model has a cultural sensitivity.

5.3.2 Ablation

463

464

465

466

467

468

469

470

471

472

473

474

To empirically validate the contribution of RLAIF
on overall performance and localization to our
AceGPT models, we conduct ablation studies
across ACVA benchmarks, Arabic Vicuna-80 and
Arabic AlpacaEval, results are outlined in Table 11.

RLAIF improves localization RLAIF results in performance gains of 27.05% and 0.84% for AceGPT-7B and AceGPT-13B in ACVA respectively, despite not being explicitly trained for them. This suggests that RLAIF enhances alignment with Arabic culture and values. Notably, the improvement from RLAIF on the 7B model is much larger than that of 13B, partially because the 7b model is weaker and therefore has more space for improvement, while it may be in saturation in the 13B model. Another reason could be that the preference data in RLAIF are generated from AceGPT-7b and therefore the learned reward model fits better AceGPT-7b than AceGPT-13b. 480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

RLAIF improves instruction-following. The results show that RLAIF significantly enhances overall model performance on both Arabic Vicuna-80 and Arabic AlpacaEval, increasing AceGPT-7B's performance by 2.81% and 2.46%, while AceGPT-13B shows an improvement of 5.74% and 4.90%, respectively. By examining the "win or tie" metric, the 7B model shows an enhancement of 3.7% through RLAIF, while the 13B model shows a significant boost of 16.2%. This narrows the gap with Turbo. These enhancements across datasets underscore RLAIF's efficacy.

6 Conclusion

AceGPT addresses the "localization issue" in large language models by specifically catering to the distinct linguistic and cultural contexts of Arabic environments, leveraging further pre-training, instruction tuning, and reinforcement learning. It excels in multiple domains, including instruction-following and natural language understanding, setting a new standard among Arabic large language models. We contribute high-quality datasets and evaluation resources, highlighting the need for localizing large language models and introducing AceGPT as a pioneering solution for Arabic linguistic and cultural adaptation.

631

632

520 Limitation

Our AceGPT model exhibits several limitations. Its 521 vocabulary, is primarily focused on Arabic letters, 522 lacking further expansion, affecting Arabic text en-523 coding efficiency. Limited machine resources dur-524 ing pre-training restricted token allocation, suggest-525 526 ing untapped potential in Arabic content processing. We omitted reasoning/misinformation and bias testing in our evaluation, raising concerns about the 528 model's safety alignment and currently limiting its use to academic research rather than online de-530 ployment. Additionally, despite manual checks, 531 the cultural dataset requires enhancement in both 532 quality and quantity, which may affect the model's 534 practicality and adoption.

Ethical Statement

536

541

542

543

544

545

546

547

548

551

552

553

554

561

563

567

568

570

571

We mainly use public data to train our models. For the newly-collected data (e.g., Quora), we use Chat-GPT to filter the questions with any ethical issues. Additionally, we have adopted RL to aligned with human values.

References

- Asaad Alghamdi, Xinyu Duan, Wei Jiang, Zhenhai Wang, Yimeng Wu, Qingrong Xia, Zhefeng Wang, Yi Zheng, Mehdi Rezagholizadeh, Baoxing Huai, et al. 2023. Aramus: Pushing the limits of data and model scale for arabic natural language processing. *arXiv preprint arXiv:2306.06800*.
 - Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR*, abs/2204.05862.
 - Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
 - Sahil Chaudhary. 2023. Code alpaca: An instructionfollowing llama model for code generation. https: //github.com/sahil280114/codealpaca.
 - Zhihong Chen, Feng Jiang, Junying Chen, Tiannan Wang, Fei Yu, Guiming Chen, Hongbo Zhang, Juhao

Liang, Chen Zhang, Zhiyi Zhang, et al. 2023a. Phoenix: Democratizing chatgpt across languages. *arXiv preprint arXiv:2304.10453*.

- Zhihong Chen, Shuo Yan, Juhao Liang, Feng Jiang, Xiangbo Wu, Fei Yu, Guiming Hardy Chen, Junying Chen, Hongbo Zhang, Li Jianquan, Wan Xiang, and Benyou Wang. 2023b. MultilingualSIFT: Multilingual Supervised Instruction Fine-tuning.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%* chatgpt quality.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instructiontuned llm.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacafarm: A simulation framework for methods that learn from human feedback.
- Momchil Hardalov, Todor Mihaylov, Dimitrina Zlatkova, Yoan Dinkov, Ivan Koychev, and Preslav Nakov. 2020. EXAMS: A multi-subject high school examinations dataset for cross-lingual and multilingual question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 5427–5444. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. Openassistant conversations - democratizing large language model alignment. *CoRR*, abs/2304.07327.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. RLAIF: scaling reinforcement learning from human feedback with AI feedback. *CoRR*, abs/2309.00267.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao,

693

- 711
- 712

M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. arXiv preprint arXiv:2211.01786.

633

634

637

645

647

651

674

675

- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. arXiv preprint arXiv:2304.03277.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. CoRR, abs/1707.06347.
- Haitham Seelawi, Ibraheem Tuffaha, Mahmoud Gzawi, Wael Farhan, Bashar Talafha, Riham Badawi, Zyad Sober, Oday Al-Dweik, Abed Alhakim Freihat, and Hussein Al-Natsheh. 2021. ALUE: Arabic language understanding evaluation. In Proceedings of the Sixth Arabic Natural Language Processing Workshop, pages 173-184, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.
- Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, Osama Mohammed Afzal, Samta Kamboj, Onkar Pandit, Rahul Pal, et al. 2023. Jais and jais-chat: Arabic-centric foundation and instruction-tuned open generative large language models. arXiv preprint arXiv:2308.16149.
- Daria Soboleva, Al-Khateeb Faisal, Myers Robert Steeves Jacob R, Hestness Joel, and Dey Nolan. 2023. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. www.cerebras.net/ blog/slimpajama-a-627b-token-cleaned-anddeduplicated-version-of-redpajama.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. Learning to summarize from human feedback. CoRR, abs/2009.01325.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford_alpaca.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. CoRR, abs/2305.17926.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit

Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ NLP tasks. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022, pages 5085-5109. Association for Computational Linguistics.

- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. arXiv preprint arXiv:2304.12244.
- Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Jianquan Li, Guiming Chen, Xiangbo Wu, Zhiyi Zhang, Qingying Xiao, et al. 2023. Huatuogpt, towards taming language model to be a doctor. arXiv preprint arXiv:2305.15075.

716

10

A Localization Issues

A.1 Sample Questions for Localization

The sample questions for Arabic name entity comparison in Table 1 and 2 are as following Table 12.

	Table	12:	20	sample	questions
--	-------	-----	----	--------	-----------

ما هي الأسماء العشرة الأكثر شيوعًا؟
What are the ten most common names?
ما هي القصص التاريخية العشرة الأكثر شهرة؟
What are the ten most famous historical stories?
أسماء المدن العشرة الأكثر شيوعاً
Names of the ten most common cities
ما ہی اکثر عشر مدن یسمع عنہا الناس اکثر؟
What are the top ten cities people hear about the most?
ما هي أهم عشر منظمات اجتماعية عربية؟
What are the ten most important Arab social organizations?
أفضل عشرة أشخاص لهذا العام
Top Ten People of the Year
كيف يحتفل الناس بالعام الحديد؟
How do people celebrate the New Year?
لماذا يؤمن الناس بالدين؟
Why do people believe in religion?
من هم العلماء المثهورون الحدد في الوطن العربي؟
Who are the new famous scholars in the Arab world?
ما هي بعض المعالم السياحية الثهيرة؟ ما هي الأحداث المحلية؟
What are some famous tourist attractions? What are local
events?
من هم الكتاب والشعراء المثهورين في التاريخ ؟
Who are the famous writers and poets in history?
ما هي بعض الأطباق الشهية التي لا يتم الاستمتاع بها إلا خلال المهرجانات الكبرى؟
What are some of the delicacies that are only enjoyed during
major festivals?
ما هو تقاليدنا الأدبية؟
What is our literary tradition?
ما هي العادات والانشطة الاساسية المهمة في حفل الزفاف؟
What are the basic customs and activities that are important
at a wedding?
ما هي الكتب المقدسة والقديسين والاماكن المقدسة؟
What are the Holy Books, Saints, and Holy Places?
ما هي بعض الحكايات الشعبية؟
What are some folk tales?
كيف نشأت لغتنا؟
How did our language originate?
ما هي أهم المهرجانات؟ لماذا توجد هذه المهرجانات؟
What are the most important festivals? Why do these
festivals exist?
من هم الأشخاص الذين يجب ان نتذكرهم؟
Who are the people we should remember?
من هم بعض الفنانين المتميزين؟
Who are some of the distinguished artists?

A.2 Case Study

717

718

719

720

721

723

724

In this subsection, we analyze the performance of AceGPT by conducting a comparative analysis of its localization ability via case studies on the sampled 20 localization questions. Illustrated in Table 22, we observed a larger proportion of Arabic events in AceGPT. The first example in Table 22 aligns with the instance illustrated in Table 2. Both AceGPT and Turbo exhibit superior responses to the given query, significantly surpassing the answer provided by Jais. Specifically, AceGPT's understanding of a 'holy book' is not solely confined to the Bible; it demonstrates a nuanced acknowledgment that different regions, especially Arabic, have their respective sacred texts, reflecting a broad and inclusive comprehension of diverse religious traditions. This illustrates the advanced capability of AceGPT, akin to Turbo, in response generation for Arabic-speaking areas.

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

759

The second example exemplifies the capability of AceGPT to incorporate more Arabic elements when responding to historical questions. Specifically, AceGPT allocates a significant proportion of its responses, 4 out of 10, to Arabic historical figures. In contrast, Turbo only attributes 1 out of 10 responses to Arabs, while Jais exclusively presents choices associated with Western figures. This demonstrates that AceGPT has an inclination towards Arabic culture, emphasizing its capability to offer more Arabic culture-relevant responses in an Arabic context.

B Construction of ACVA

We employ a top-down approach for the construction of the Arabic Cultural and Value Alignment benchmark. First, we gathered over 50 topic keywords (see Table 13) representing various aspects of Arabic culture, including humanity, art, science, geography, history, manners, religion, and the influence between civilizations, sourced from several books on Arabic culture and values. Then, we query Turbo to generate 8000 data based on the given topic using the prompt shown below, where topic is the placeholder for the topic.

I am collecting some supervised fine tun-760 ing (sft) data about Arabic culture. It is 761 about the knowledge of Arabic culture 762 and manners. The data is some ques-763 tions in the Arabic language with an id 764 in the form of {"id": "1","label":"xx" "query":"xx" }. I will give you a topic in 766 Arabic culture. The "id" is the index of 767 the data. "label" is the topic I give you. 768 "query" is some question statement about 769 Arabic culture under that topic. The Data 770 should be of no repetition with a bal-771 anced proportion of true and false. Now 772 please generate 200 sft data in json in 773 arabic with the format under the topic of 774

5	topic
6	We further sample 50% topics to verify the rel-
7	evance of questions to Arabic cultures and values
8	and the accuracy of the Yes-No labels, which were
9	reviewed by Arabic speakers, leading to a 'Clean
0	set'.

77

783

784

791

794

803

804

806

809

823

C Preference Data for RLAIF

The data comprises two parts: Arabic preference data and open-source English preference data. Outputs for Arabic preference data are sampled from our fine-tuned 7B model with a temperature of 1. The open-source English preference data is incorporated to improve the generalization capability of the reward model and alleviate GPT4-preference hacking. We randomly sample 12K from three public human-annotated datasets - Anthropic helpfulness and harmlessness (Bai et al., 2022), OpenAI Summarize (Stiennon et al., 2020), and OpenAssistant Conversations (OASST1) (Köpf et al., 2023).

The core idea of preference labeling for Arabic preference data is to use a GPT-4 model with prompts as an automatic annotator to assess two responses generated by the same model for a given question. However, a significant challenge emerges as GPT-4 often shows a marked preference for the first response, around 80% of the time, with the exact percentage varying based on the specific prompt design. To counter this, we utilize an order-switch mechanism to ensure consistent preference data across two separate runs of GPT-4. In one run, two responses are placed arbitrarily, and in the other, their orders are switched. The prompt for labeling is shown below. instruction, response 1, and response 2 are the placeholders for the input instructions and the two generated responses.

810	[Instruction]
811	instruction
812	
813	[Response1]
814	response 1
815	
816	[Response2]
817	response 2
818	
819	[System]
820	A good response show
821	accurate, helpful, and

A good response should be relevant, accurate, helpful, and have enough details. Also, it should follow arabic culture, customs, beliefs and laws.

Which is better, Response1 or Re-	824
sponse2?	825
Do not explain your answer, just output	826
Response1 or Response2.	827
	828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

D Implementation of Training

D.1 Pre-Training

We employ the LLaMA2 framework for the pretraining process, capitalizing on a computational setup furnished with 24 Nvidia A100 80G GPUs. We configure the context length at 2048 tokens and adopt the AdamW optimizer, paired with a cosine learning rate scheduler. The learning rate is set at 1e-4. Given a gradient accumulation setting of 128, the total batch size amounts to 3072. Additionally, a warm-up phase is integrated, constituting 5% of the total training duration.

D.2 Supervised Fine-Tuning

We train for one epoch using a variety of datasets in Table 14. Native Arabic data like *Alpaca-Arabic-GPT4* and *Quora-Arabic-GPT4* are included thrice in the mixture, while datasets like *ShareGPT* and *Alpaca-Chinese-GPT4* are included once to minimize non-Arabic data ratio, totaling 629,293 data points.

Both AceGPT-7B and AceGPT-13B are finetuned with 8 Nvidia A100 80G GPUs. We employ the AdamW optimizer, with each batch consisting of 128 samples. We adopt different configurations for the learning rate based on the model architecture. For *AceGPT-7B*, the maximum learning rate is set to 5×10^{-5} , and for AceGPT-13B, it is 1×10^{-5} . A cosine scheduler is employed for learning rate adjustment, with a warmup rate of 0.03.

Following LLaMA2, we use the following form of system prompt:

[INST] ((SYS))	860
أنت مساعد مفيد ومحترم وصادق. أجب دائما بأكبر قدر ممكن من	861
المساعدة بينما تكون آمنا. يُحب ألا تتضمن إجاباتك أي محتوى ضار أو غير	862
أخلاقي أو عنصري أو جنسي أو سام أو خطير أو غير قانوني. يرجى التأكد	863
من أن ردودك غير متحيزة اجتماعيا وإيجابية بطبيعتها.	864
إذا كان السؤال لا معنى له أو لم يكن متماسكا من الناحية الواقعية، اشرح	865
السبب بدلا من الإجابة على شيء غير صحيح. إذا كنت لا تعرف إجابة	866
سؤال ما، فيرجى عدم مشاركة معلومات خاطئة.	867
< <sys>></sys>	868
[question] [INST]	869
The corresponding meaning in English is:	870

Table 13:	Topics	for	ACVA	construction
-----------	--------	-----	------	--------------

Country	Algeria, Bahrain, Comoros, Egypt modern, Iraq, Jordan, Kuwait, Lebanon, Libya, Mauritania, Morocco, Oman, Palestine, Qatar, Saudi Arabia, Somalia, Sudan, Syria, Tunisia, United Arab Emirates, Yemen
Civilization Relation	Influence From Ancient Egypt, Influence From Byzantium, Influence From China, Influence From Greece, Influence From Persia, Influence From Rome, Mesopotamia civilization
Science and Humanity	Arabic Astronomy, Arabic Math, Arabic Medicine, Arabic Physics and Chemistry, Arabic Literature, Arabic Music, Arabic Philosophy, Arab Empire, Arabic Architecture, Arabic Art, Arabic Calligraphy, Arabic Geography, Arabic History, Arabic Language Origin
Manners and Religion	Arabic Ceremony, Arabic Clothing, Arabic Culture, Arabic Food, Arabic Funeral, Arabic Ornament, Arabic Wedding, mindset, Special Expression, daily life, Influence From Islam, Islam branches and schools, Islam Education, Islamic law system

Table 14: Instruction Tuning Datasets; Datasets Constructed in This Work Are Highlighted in **bold**.

Data	Source	Numbers	
Data	questions	responses	TNUITDETS
Quora-Arabic-40K	collected from Quora	GPT-4	43,050
Alpaca (Peng et al., 2023) Alpaca-Chinese (Peng et al., 2023) Alpaca-Arabic	self-instruct (Taori et al., 2023) Turbo translated (Peng et al., 2023) GPT-4 translated from (Taori et al., 2023)	GPT-4	49,969 49,969 49,969
Code-Alpaca-Arabic	GPT-4 translated from (Chaudhary, 2023)	GPT-4	20,022
Evol-Instruct-Arabic	GPT-4 translated from (Xu et al., 2023)	GPT-4	69,997
ShareGPT	humans	ChatGPT	80,179

[INST] <<SYS>>

871

872

873

874

875

877

878

879

882

886

887

892

893

You are a helpful, respectful, and honest assistant. Always answer with the utmost assistance while being safe. Your answers should not include any harmful, unethical, racist, gender discriminatory, toxic, dangerous, or illegal content. Please ensure that your responses are not socially biased and are positive.

If the question is meaningless or isn't coherent in a realistic sense, explain the reason instead of answering something incorrectly. If you do not know the answer to a question, please refrain from sharing.

<<sys>>

[question] [INST]

D.3 Reward Model Training

The reward model is initialized with Ziya, an opensource 7B reward model ⁶. We use 8 Nvidia A100 80G GPUs for training. Each batch consists of 128 samples. We take two epochs with the AdamW optimizer. The maximum learning rate is set to 8e-6 and the warmup rate is set to 0.03 with cosine scheduler.

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

D.4 PPO

We implement PPO with DeepSpeed-Chat ⁷. The actor parameters are initialized with our fine-tuned models and the critic parameters are initialized with our trained 7B reward model. We sample 448 experiences with the mini-batch size of 224 ⁸, which is updated in only one epoch. The maximum learning rate for the actor is set to 5e-7 while that for the critic is set to 5e-6. A cosine scheduler is used for learning rate adjustment with a warmup step of 100. We set the KL penalty as 0.01. The policy gradient loss is clipped with the threshold as 0.2 while that for the value loss is 0.3. The reward is clipped to be [-5, 5]. The gamma and lambda for the generalized advantage estimation are 1 and 0.95 respectively.

Notably, both AceGPT-7B and AceGPT-13B are trained with the 7B reward model whose preference data only comprises outputs from the 7B policy model (post-SFT).

⁶https://huggingface.co/IDEA-CCNL/Ziya-LLaMA-7B-Reward

⁷https://github.com/microsoft/DeepSpeedExamples

⁸We use 7 GPUs, where there are 64 experiences and 32 samples in each mini-batch on each GPU.

916 917

918

919

920

922

926

928

929

930

931

932

933

934

935

936

937

942

943

947

951

953

955

957 958

959

961

963

E Implementation of Evaluation

E.1 Baselines and Benchmarks

We use the following baselines :

- LLaMA2 (Touvron et al., 2023), developed by Meta AI, are the most popular open-source large language models ranging in scale from 7 billion to 70 billion parameters. Our AceGPT models are also built upon LLaMA2-7B and -13B. We compare our AceGPT-base models to the corresponding size of LLaMA2.
 - Bloomz (Muennighoff et al., 2022) and Phoenix (Chen et al., 2023a,b): Bloomz is a classical family of multilingual models finetuned with multiple traditional NLP tasks. Phoenix are multilingual instruction following models using Bloomz as the backbone. We compare AceGPT-base models to Bloomz and AceGPT-chat models to Phoenix.
 - Jais (Sengupta et al., 2023) are concurrent open-source 13B Arabic-centric LLMs, including a foundation base model and an instruction-tuned model. We compare AceGPT-base and AceGPT-chat to their base and chat models respectively.
 - **Turbo** is the most popular and powerful closed-source multilingual LLM, second only to GPT-4. We compare both AceGPT-base and AceGPT-chat to it.

We use Arabic Vicuda-80, Arabic AlpacaEval, Arabic MMLU, Exams, ALUE, ACVA as our benchmarks, see Table 15.

In the scoring method for the Arabic MMLU dataset, the process begins by calculating the average score within each subcategory, this average represents the mean score for that particular subcategory. Subsequently, the average scores of all subcategories within a main category are aggregated and averaged again to yield the overall average score for that main category. Following this, the average scores of all main categories are collated and averaged to determine the overall average score for that category within the entire Arabic MMLU dataset. Finally, the aggregate average scores of all categories in the dataset are combined and averaged to arrive at the comprehensive average score across the entire Arabic MMLU dataset.

E.2 Evaluation on Instruction Following

We apply greedy decoding strategies for all models evaluated on Arabic Vicuna-80 and Arabic AlpacaEval benchmarks.

GPT-4 evaluation The temperature of the GPT-4 evaluator is set to 0.2, and the scoring is reiterated three times. The evaluation prompt is shown below. question, response 1, and response 2 are the placeholders for the input questions and the two generated responses. To mitigate position bias (Wang et al., 2023), we vary the order in which the model pair is presented in the prompt.

964

965

966

967

968

969

970

971

	_
[Question]	973
question	974
	975
[Assistant 1]	976
response 1	977
	978
[End of Assistant 1]	979
	980
[Assistant 2]	981
	982
response 2	983
	984
[End of Assistant 2]	985
	986
[System]	987
	988
We would like to request your feedback	989
on the two AI assistants in response to	990
the user question displayed above.	991
	992
Requirements: The response should be	993
helpful, relevant, accurate, and have	994
enough details. The response should	995
follow Arabic culture, customs, beliefs,	996
and laws.	997
	998
Please rate the performance of their	999
responses. Each assistant receives an	1000
overall score on a scale of 1 to 10,	1001
where a higher score indicates better	1002
performance.	1003
	1004
Please first output a single line contain-	1005
ing only two values indicating the scores	1006
for Assistant 1 and 2, respectively. The	1007
two scores are separated by a space. You	1008
should consider which response is more	1009
in line with the given requirements.	1010
	1011
In the subsequent line, please provide a	1012

Table 15: Evaluation Benchmark	ζS.
--------------------------------	-----

			D () (0'	T 6 1
Benchmark	Evaluation Aspects	Type of Evaluation	Dataset Size	Types of examples
Arabic Vicuna-80 Arabic AlpacaEval	Instruction following	Human & Automated	80 805	Freely-answered Questions
Arabic MMLU EXAMs	Knowledge Ability	Automated	14k 0.5k	Multiple-choice Questions
ALUE(see Appendix F.2)	Language Understanding	Automated	18k	Classification & Regression
ACVA-all ACVA-clean	Arabic Cultural and Value Alignment	Automated	9k 2.4k	Yes/no binary Questions

1013	comprehensive explanation of your eval-
1014	uation.

1015Human evaluationWe recruited 11 native peo-1016ple for annotation, including verification of the lo-1017calization dataset, calibration of translation results,1018and human evaluation, the backgrounds of these1019people can be found in Table 16. The evaluation1020interface is illustrated in Figure 2.

Table 16: Information of participants involved in theAceGPT testing

Name	Gender	Education	Language
Participant 1 Participant 2 Participant 3 Participant 4 Participant 5 Participant 6 Participant 7 Participant 8 Participant 9 Participant 10	female male female female female male female female female	PGDip PhD PGDip PGDip PGDip PGDip Master PGDip Master PGDip	Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native Arabic-Native
I articipant II	maie	THD	Thuble Hullve

E.3 Evaluation on Knowledge

1021

1022

1023

1024

1025

1026

1027 1028

1029

1030

1031

1032 1033

1034

1035

1037

There are two main differences in the MMLU evaluation between (Sengupta et al., 2023) and ours: (1) we translate MMLU into Arabic differently. The machine-translated version in (Sengupta et al., 2023) is facilitated through their in-house translation model while we leverage Turbo. Additionally, (Sengupta et al., 2023) further creates a humantranslated version. Unfortunately, both the humantranslated and machine-translated versions are not publicly available, which prevents us from evaluating on the same benchmark; (2) we adopt the widely accepted few-shot prompting setting commonly found in related literature for base models, while (Sengupta et al., 2023) opts the zero-shot setting. Due to these differences in translation methods and evaluation settings, the performance

metrics between the two works are not directly comparable.

1038

1039

1040

1041

1042

1043

1044

We benchmark Jais-13B-*base* and Jais-30B-*base* using our Turbo-translated MMLU dataset under the standard few-shot setting in Table 8. Moreover, we also benchmark Jais-13B-*chat* using the zeroshot setting in Table 18.

The evaluating template is shown as below:

The evaluating template is shown as below.	1045
• Few-shot	1046
فيما يلى أسئلة الاختيار من متعدد (مع الإجابات) حول	1047
[category]	1048
[exemplars]	1049
سؤال: [question]	1050
إجابة:	1051
• Zero-shot فيما يلي أسئلة الاختيار من متعدد حول [category]	1052 1053
	1054
	1055
سؤال: [question]	1056
من فضلك اختر إجابة واحدة من بين 'A, B, C, D' دون	1057
شرح.	1058
The corresponding meaning in English is:	
• Few-shot	1059 1060
Below are multiple choice	1061
questions (with answers) about	1062
[category]	1062
	1000
[exemplars]	1064
	1065
Question: [question]	1066
Answer:	1067
• Zero-shot	1068
Below are multiple choice ques-	1069
tions about [category]	1070

Figure 2: User interface of human annotation. Response positions are randomized to alleviate biases.

🏜 Human Evaluation				
notes: We will present you some questions and the answers from the two models, and please rate which of the two models has the better answer. Please note the following points: 1. general, non-specific responses should be scored low 2. If a question is supposed to be in Arabic, the model's responses in other languages should be greatly downgraded 3. Consider richness of model responses, logical clarity 4. Consider the professionalism and accuracy of the model				
ما هو طول القطعة الخطية إذا كانت نقاط النهاية لها (2، -2) و (10، 4)؟ Question: أو الم				
answer of model A	answer of model B			
:طول القطعة الخطية يمكن حسابه باستخدام الصيغة التالية	الحساب طول القطعة الخطية بين نقطتين، نستخدم الصيغة التالية			
$ f(x) = \sqrt{((x2 - x1)^2 + (y2 - y1)^2)}$	الطول = (x2-x1)*(y2-y1) + (y2-y2)*(x2-x2)			
.لكل من نقاط النهاية y هما قيم y2 و y1 لكل من نقاط النهاية، و x هما قيم x2 و x عيث أن	. هما نقطتي النهاية (x2، y2) و (x2، y1) حيث			
:هـي 10 و 10. نستخدم هذه القيم في الصيغة للحصول على x2 هـي 2 و -2، و قيم x1 في هذه الحالة، قيم	:بتطبيق هذه الصيغة على النقاط المعطاة، نحصل على			
$ f(x) = \sqrt{(10 - 2)^2 + (10 - (-2))^2)} = \sqrt{((8)^2 + (12)^2)} = \sqrt{(64 + 144)} = \sqrt{208}$	$(2-10)^{*}(2-4) + (2-4)^{*}(2-10) = 1 $			
. إذاً، طول القطعة الخطية هو / 208 وحدة	الطول= 4+4			
	الطول= 8			
	. إذاً، طول القطعة الخطية هو 8 وحدات			
, ,				
->model A is better	model B is better			
already answerd 80 your choice: draw				
previous question	next question			

	previous question	next question	
1071	[question]	سؤال: علاً مدرس العلوم في الدرسة الثانوية زجاجة سعتها ١ لتر	1098
1072	Please choose one answer from	بالنيتروجين النقي ويغلق الغطاء. الضغط هو ٢٠٠١ جوي، ودرجة حرارة	1099
1072	among 'A B C D' without avala	الغرفة هي ٢٥ درجة مئوية. ما هما المتغيران اللذان سيزيدان ضغط النظام	1100
1073	among A, B, C, D without expla-	مع الحفاظ على كل المتغيرات الأخرى ثابتة؟	1101
1074	nation.		1102

Thease choose one answer from
among 'A, B, C, D' without expla-
nation.

A specific example of five-shot prompting is:

	1075
فيما يلي أسئلة الاختيار من متعدد (مع الإجابات) حول طب جامعي	1076
سؤال: كيف يتم نقل الجلوكوز إلى خلية العضلات؟	1077
A. عبر ناقلات البروتين المسماة جموة ع.	1078
.B فقط في وجود الأنسولين.	1079
C. عبر الهكسوكيناز.	1080
.D عبر ناقلات حمض المونوكاربيليك.	1081
إجابة: А	1082
. يحتوي عدد كبير من عضلات الساق لدى العدائين المتحمسين للتحمل B.	1083
على ألياف من النوع اي	1084
C. جليكوجين الكبد مهم للحفاظ على تركيز الجلوكوز في الدم	1085
D. الأنسولين يعزز امتصاص الجلوكوز من جميع الأنسجة في الجسم	1086
إجابة: D	1087
سؤال: في اختبار جيني لرضيع حديث الولادة، يتم العثور على اضطراب	1088
جيني نادر ينتقل بشكل متنحي على قاعدة الصلة بالصبغي خ أي من	1089
العبارات التالية تعتبر صحيحة بشكل محتمل بخصوص مخطط هذا الاضطراب	1090
الجيني؟	1091
A. سيكون لدى جميع الأحفاد على الجانب الأمريكي المصابين	1092
بالاضطراب	1093
B. سيكون الإناث على مقربة من ضعف الذكور المصابين في هذه العائلة	1094
C. سيكون جميع البنات من ذوي الآباء المصابين مصابين بالمرض	1095
.D ستكون هناك توزيع متساوٍ للذكور والإناث المتأثرين بالمرض.	1096
إجابة: C	1097

A. زيادة درجة الحرارة، زيادة مولات الغاز	1102
B. زيادة درجة الحرارة، زيادة الحجم	1103
C. تقليل الحجم، تقليل درجة الحرارة	1104
D. تقليل مولات الغاز، زيادة الحجم	1105
إجابة: A	1106
سقال: ما هم الآثر الجانم التدقم لتكملة الكرياتين؟	1107
مون. لم يو در و ب بي محموع ما يمد ما و بي ون. A. ضعف العضلات.	1108
B. زيادة في كتلة الجمم.	1109
۔ .C. تشنجات العضلات.	1110
.D فقدان الكهرليتات.	1111
إجابة: B	1112
سؤال: ما هي الحينوم؟	1113
A. هو عدد الكروموسومات الموجودة في الخلايا الجسدية للكائن الحي.	1114
۔ B. هو تسلسل الحمض النووي الكامل لكائن حي.	1115
C. هي القائمة الكاملة للبروتينات التي يمكن للكائن الحي إنتاجها.	1116
D. هي تسلسل الحمض النووي الخاص بكروموسوم كامل.	1117
إجابة .	1118
A specific example of zero-shot prompting is:	
	. 1119

فيما يلي أسئلة الاختيار من متعدد حول المعرفة السريرية	1120
مد جيجا. ۸. ا	1121
, B. جيجا.	1122
. _{۲۲} جيجا.	1123
. _{۲٤} D.	1124
من فضلك اختر إجابة واحدة من بين 'A, B, C, D' دون شرح.	1125

1	1	26
1	1	27

E.4 Evaluation on Arabic Cultural and Value Alignment

The evaluation prompt for ACVA is

28	1128
9 • Few-shot فيما يلى أسئلة نعم أو لا (مع الإجابات) حول	1129 1130
[category]	1131
12	1132
3	1133
4 [exemplars]	1134
سؤال: [question]	1135
إجابة:	1136
7 • Zero-shot الاجاء التحكر في ما إذا كانت الجملة التالية محيجة أو لار إذا	1137 1138
» رب مناعم في قد به مناع بسه ملي في عام ماري. كانت صحيحة ، فالرجاء الرد بـ «نعم». إذا كانت غير صحيحة ،	1139
فالرجاء الرد بـ »لا	1140
	1142
سؤال: [question]	1143

The corresponding meaning in English is:

1145• Few-shot1146Below are multiple1147questions (with answer1148[category]1149Item [category]1150[exemplars]1151Question: [question]	
1146Below are multiple1147questions (with answer1148[category]1149[exemplars]1150[exemplars]1151Question: [question]	
1147questions (with answer1148[category]1149[exemplars]1150[exemplars]1151Question: [question]	choice
1148[category]1149[exemplars]1150[exemplars]1151Question: [question]	rs) about
11491150[exemplars]1151Question: [question]	
1150[exemplars]1151Question: [question]	
1151 Question: [question]	
1152 Answer:	
• Zero-shot	
1154 Please determine whe	ther the
1155 following sentence is true	or not. If
it is true, please respond w	vith 'Yes'
1157 If it is not true, please resp	ond with
'No'.	
1158	
1160 Question: [question]	

A specific example of five-shot prompting is
--

1161	A specific example of five shot prompting it
1162	فيما يلي أسئلة نعم أو لا (مع الإجابات) حول الجزائر
1163	سؤال: الجزائر تعتبر جزءًا من الاتحاد العربي.
1164	إجابة: نعم
1165	سؤال: العربية هي اللغة الرسمية في الحزائر.
1166	إجابة: نعم

Table 17: Average F1 on ACVA in the few shot setting. The best performance is in **bold** and the second best is underlined.

Model	All set	Clean set	
Bloomz	58.94%	60.91%	
Jais-13B-base	73.96%	75.80%	
Jais-30B-base	73.81%	<u>77.44%</u>	
AceGPT-7B-base AceGPT-13B-base	74.72% <u>78.62</u> %	70.32% 75.61%	
Turbo	80.12%	81.99%	
استقلالها من الاستعمار الإسباني.	مجحت في استعادة	سؤال: ثورة الجزائر نم	1167
-	-	إجابة: لا	1168
ا: بغية.	الحائم الأه	سةال: اللغة السمية ف	1169
	ي الشرائر علي الله	سون، منت مر ميم . إجابة: لا	1170
ازيغية.	لي الجزائر هي الأه	سؤال: اللغة الرسمية ف	1171
		إجابة: لا	1172
رئيسي على الصناعات التحويلية.	ئرى يعتمد بشكل	سؤال: الاقتصاد الجزا	1173
		إجابة:	1174
A specific exemple of	zoro chot	prompting is.	
	2010-51101	prompting is.	1175
ملة التالية صحيحة أم لا. إذا كانت	ما إذا كانت الح	»الرجاء التحكم في	1176
لت غیر صحیحة، فالرجاء الرد بـ »لا	ب ،،نعم،، إذا كا	صحيحة، فالرجاء الرد	1177
الفصول الأربعة والمواسم.	نوس دينية لتنظيم	سؤال: كان لديهم طق	1178
F More Experiment	ts of Ace	GPT	1179
Evaluation			1180
F.1 Supplementary Ex	xperimen	tal Results	1181
ACVA evaluation und	er the fe	w-shot setting.	1182
Table 17 demonstrates	the perfor	rmance of base	1183
models on ACVA. AceG	PT-30B-b	ase outperforms	1184
Jais-30B-base by 4.819	% in 'All	set', but fails	1185
slightly 1.83% behind it	in 'Clean	set'.	1186
Knowledge evaluation of	on the cha	at models. We	1187
evaluate chat models in	the zero.	-shot setting on	1188
Arabic MMLU and EXA	Ms. As i	llustrated in Ta-	1189
ble 18, Turbo consistentl	y outperf	orms other mod-	1190
els in both MMLU and H	EXAMs b	enchmarks. No-	1191
tably, Jais-13B-chat sho	wcases th	he superior per-	1192
formance, which is cons	sistent wi	th the results in	1193
(Sengupta et al., 2023).	Specifica	ally, its MMLU	1194
score stands at 37.11, tra	iling Cha	tGPT's score of	1195
46.07 by a mere 8.96 poin	ts. On the	EXAMs bench-	1196
mark, Jais-13B-chat scor	ed only 4	.79 points lower	1197
than Turbo. One possib	le reason	for Jais's good	1198
performance may be attr	ibuted to	traditional NLP	1199

1200task datasets in their SFT dataset such as Super-1201NaturalInstructions (Wang et al., 2022), which con-1202tains multiple-choice questions akin to the MMLU1203and EXAMs. Our model, in contrast, hasn't been1204trained on such data.

5 F.2 Evaluation on Arabic NLU Tasks

ALUE ALUE ⁹ is a popular online benchmark, which is similar to the GLUE benchmark but has a main focus on Arabic Language Understanding Evaluation. It includes traditional NLP tasks such as sentiment analysis, semantic matching, text relation classification, and dialect identification. It comprises 9 tasks as illustrated in Table 19.

Experiment setting We train our AceGPT-13B-1213 base on each task independently in a fully su-1214 pervised manner, resembling the approach of the 1215 top models on the leaderboard. Moreover, high-1216 1217 ranking models on the leaderboard adopt the grid search method on validation sets to select hyper-1218 parameters. Similarly, we employ a Bayesian ap-1219 proach for hyperparameter adjustment. For tasks 1220 providing predefined validation split, we utilize the 1221 given validation sets. Otherwise, we allocate 10% 1222 of the data from the training set for validation pur-1223 poses. For the DIAG task, which does not provide 1224 1225 training data, we use the model trained on XNLI to evaluate it. 1226

Experiment results and analysis Table 20 1227 presents our performance on the ALUE benchmark. 1228 AceGPT ranks second in terms of the average score 1229 in these nine datasets, right behind AraMUS ((Al-1230 ghamdi et al., 2023)), which has conducted exten-1231 sive pre-training in Arabic data. In future endeav-1232 ors, we plan to incorporate a richer set of Arabic 1233 pre-training corpora and supervised data to enhance 1234 the model's NLU capabilities. 1235

G Detailed Results on Human Evaluation

The results of the human evaluation corresponding to Table 7 for each annotator are shown in Table 21.

1238 1239

1236

1237

1206

1207

1208 1209

1210

1211

⁹https://www.alue.org/home

	Arabic MMLU						
Model	Average	STEM	Humanities	Social Sciences	Others (Business, Health, Misc)	EXAMs	
Phoenix Phoenix-multiple-langs Jais-13B-chat	29.65 17.37 <u>41.06</u>	27.06 16.77 <u>39.82</u>	28.35 15.65 <u>42.21</u>	31.66 18.51 <u>41.75</u>	31.54 18.54 <u>42.50</u>	31.60 16.48 <u>40.84</u>	
AceGPT-7B-chat AceGPT-13B-chat	31.22 36.65	27.63 34.72	31.99 34.19	31.37 40.02	33.88 37.68	29.73 40.35	
Turbo	46.07	44.17	35.33	61.26	43.52	45.63	

Table 18: Accuracy of chat models on Arabic MMLU and EXAMs. The best is in **bold** and the second is <u>underlined</u>.

Table 19: Summary of NLU Tasks and Metrics in ALUE benchmark

Task	Metric	Size	Test Set Ground Truth
MQ2Q (NSURL-2019 Shared Task 8)	F1-score	4000	private
OOLD (OSACT4 Shared Task-A)	F1-score	1000	private
OHSD (OSACT4 Shared Task-B)	F1-score	1000	private
SVREG (SemEval-2018 Task 1)	Pearson correlation	1000	private
SEC (SemEval-2018 Task 1)	Jaccard similarity score	1000	private
FID (IDAT@FIRE2019)	F1-score	1006	public
MDD (MADAR Shared Task Subtask 1)	F1-score	5200	public
XNLI (Cross-lingual Sentence Representations)	Accuracy	2490	public
DIAG (Diagnostic dataset)	Matthews correlation	1147	public

Table 20: Experimental results in ALUE (Seelawi et al., 2021) including online baselines. While the leaderboard calculates the 'scores' excluding Task DIAG, we also incorporate it to derive the 'Avg'.

Model	#Params	Avg	Score	MQ2Q	MDD	SVREG	SEC	FID	OOLD	XNLI	OHSD	DIAG
ARABIC-BERT	135M	63.5	67.1	85.7	59.7	55.1	25.1	82.2	89.5	61.0	78.7	19.6
ARABERTv0.1-base	135M	64.2	68.4	89.2	58.9	56.3	24.5	85.5	88.9	67.4	76.8	23.5
ARABIC-BERT	110M	68.6	69.3	89.7	59.7	58.0	26.5	84.3	89.1	67.0	80.1	19.0
CAMeLBERT-MIX	108M	66.7	70.4	89.4	61.3	69.5	30.3	85.5	90.3	56.1	80.6	11.8
AraT5-base	289M	67.6	71.1	91.3	63.8	65.9	30.5	82.3	88.8	68.2	77.9	15.4
ARBERT	163M	65.5	71.4	89.3	61.2	66.8	30.3	85.4	89.5	70.7	78.2	24.3
MARBERT	163M	63.9	72.2	83.3	61.9	75.9	36.0	85.3	92.1	64.3	78.9	12.3
JABER	135M	68.2	73.7	93.1	64.1	70.9	31.7	85.3	91.4	73.4	79.6	24.4
Char-JABER	136M	70.1	75.3	92.0	66.1	74.5	34.7	86.0	92.3	73.1	83.5	26.7
ALM-1.0	350M	70.3	75.8	94.5	65.1	70.1	35.3	86.0	91.7	77.7	85.7	30.2
SABER	369M	71.4	<u>77.3</u>	93.3	<u>66.5</u>	<u>79.2</u>	<u>38.8</u>	<u>86.5</u>	93.4	76.3	84.1	26.2
AraMUS	11B	74.0	79.8	95.2	67.5	80.4	41.6	87.2	95.5	83.2	87.4	<u>42.0</u>
AceGPT-13B-base	13B	72.8	76.6	94.9	63.3	72.4	36.8	85.1	94.2	81.0	85.4	42.2

Dataset	Comparison		win	tie	lose
	AceGPT-7B-chat vs. Jais-13B-chat	volunteer 1 volunteer 2 volunteer 3	66 65 67	3 9 4	11 6 9
Arabic Vicuna-80	AceGPT-7B-chat vs. Turbo	volunteer 1 volunteer 2 volunteer 3	26 40 0	0 0 79	54 40 1
	AceGPT-7B-chat (w/o RLAIF) vs. Turbo	volunteer 1 volunteer 2 volunteer 3	23 12 31	12 58 0	45 10 49
	AceGPT-13B-chat vs. Jais-13B-chat	volunteer 1 volunteer 2 volunteer 3	68 65 66	6 5 5	6 10 9
	AceGPT-13B-chat vs. Turbo	volunteer 1 volunteer 2 volunteer 3	14 21 4	35 28 74	31 31 2
	AceGPT-13B-chat (w/o RLAIF) vs. Turbo	volunteer 1 volunteer 2 volunteer 3	19 22 6	14 19 57	47 39 17
	AceGPT-7B-chat vs. Jais-30B-chat	volunteer 1 volunteer 2 volunteer 3	54 48 60	6 30 0	20 2 20
	AceGPT-13B-chat vs. Jais-30B-chat	volunteer 1 volunteer 2 volunteer 3	48 55 52	29 0 7	3 25 21
Arabic AlpacaEval	AceGPT-7B-chat vs. Jais-13B-chat	volunteer 1 volunteer 2 volunteer 3	515 619 146	196 54 632	94 132 27
	AceGPT-7B-chat vs. Turbo	volunteer 1 volunteer 2 volunteer 3	259 71 158	291 632 200	255 102 447
	AceGPT-13B-chat vs. Jais-13B-chat	volunteer 1 volunteer 2 volunteer 3	283 400 509	504 343 187	18 62 109
	AceGPT-13B-chat vs. Turbo	volunteer 1 volunteer 2 volunteer 3	216 47 346	326 664 84	263 94 375

Table 21: Details of human evaluations on Arabic Vicuna-80 and Arabic AlpacaEval.

Table 22: Examples of comparison of LLM Responses. Arabic-rooted terms in blue; non-Arabic terms in red.

User:

ما هي الكتب المقدسة والقديسين والأماكن المقدسة؟

(What are the holy books, saints, and holy places?)

AceGPT-7B-chat:

الكتب المقدسة هي الكتب التي تعتبر مقدسة أو من الله أو من الآلهة في الديانات المختلفة. تعتبر هذه الكتب الصدر الرئيسي للتعاليم والمعتقدات والشعائر في العديد من الديانات. تشعل الكتب المقدسة مثل الكتاب المقدس في المسيحية، القرآن في الإسلام، التوراة في الهبودية، والفيدا في الهندوسية. القديسين هم أشخاص يعتبرون من المؤمنين بأنهم قد حققوا مستوى عال من القداسة أو الإلهية، وفي بعض الأحيان يعتبرون أنهم وصلوا إلى حالة من النعمة الإلهية أو الوحدة مع الله. يعتبر القديسين شخصيات مقدسة ومحترمة في الديانات المسيحية والإسلامية واليهودية، والفيدا في المقدسة هي المواني يعتبرون أنهم وصلوا إلى حالة من النعمة الإلهية أو الوحدة مع الله. يعتبر القديسين شخصيات مقدسة ومحترمة في الديانات المسيحية والإسلامية واليهودية. الأماكن المقدسة هي الواقع التي تعتبر مقدسة أو مكرسة في الديانات المختلفة. تشمل هذه الأماكن المعابد، الكنائس، المباجد، الأديرة، الواقع الدينية، والأماكن التي يعتقد أن الأحداث الدينية المهمة قد وقعت فيها. مثل القدس في فلسطين، ومكة في الملكة السمودية، والفاتيكان في إيطاليا، ومصر والهند ولندن ونيويورك وغيرها من الموامة في محيم الها.

(Holy books are books that are considered sacred or from God or the gods in different religions. These books are considered the main source of teachings, beliefs, and rituals in many religions. It includes sacred books such as the Bible in Christianity, the Qur'an in Islam, the Torah in Judaism, and the Vedas in Hinduism. Saints are people who are considered by believers to have achieved a high level of holiness or divinity and are sometimes considered to have reached a state of grace. Divinity or unity with God. Saints are considered and respected figures in the Christian, Islamic, and Jewish religions. Holy places are sites that are considered sacred or consecrated in different religions. These places include temples, churches, mosques, monasteries, religious sites, and places where important religious events are believed to have occurred. Such as Jerusalem in Palestine, Mecca in Saudi Arabia, the Vatican in Italy, Egypt, India, London, New York and other important sites in the world.)

Jais-13B-chat:

الكتاب المقدس هو مجموعة من النصوص الدينية التي تعتبر مقدسة لدى المسيحيين. القديسون هم الأشخاص المقدسون الذين يعتبرون أنهم عاشوا حياة فاضلة بشكل استثنائي ويتم تبجيلهم كقديسين من قبل بعض الأديان. الأماكن المقدسة هي مواقع تعتبر مقدسة أو مكرسة من قبل بعض الأديان ، مثل الكنائس والعابد والمساجد (The Bible is a collection of religious texts that are considered sacred by Christians. Saints are holy people who are considered to have lived exceptionally virtuous lives and are venerated as saints by some religions.

Holy places are sites considered sacred or consecrated by some religions, such as churches and temples And mosques.)

Turbo:

الكتب المقدسة هي الكتب التي يعتقد أنها موحى بها من قبل الله وتحمل رسالته للبشرية. في الديانات السماوية الثلاث (اليهودية والمسيحية والإسلام)، تعتبر الكتب المقدسة هي: اليهودية: التوراة (العهد القدم)، وهي تنضن الكتب الخمسة الأولى من الكتاب المقدس اليهودي وتشمل التكوين والحروج واللاويين والعدد والتثنية. المسيحية: الكتاب المقدس السيعي يتكون من العهد القدم (التوراة) والعهد الجديد. العهد الجديد يتضمن الأناجيل الأربعة (متى ومقس ولوقا ويوحنا) ورسائل الرسل وسفر الرؤيا. الإسلام: القرآن الكريم، وهو الكتاب المقدس للمسلمين ويعتقد أنه وحي من التوراة) والعهد الجديد. العهد الجديد يتضمن الأناجيل الأربعة (متى ومقس ولوقا ويوحنا) ورسائل الرسل وسفر الرؤيا. الإسلام: القرآن الكريم، وهو الكتاب المقدس للمسلمين ويعتقد أنه وحي من الله للنبي محمد. أما القديسين، فهم الأنحيان الذين يعتقد أنهم قد تحققوا درجة عالية من التقوى والقداسة في الحياة الدنيا ويعتبرون مثالاً يعتدى به للمؤمنين. في المسيحية، يتم تكريم من الله للنبي محمد. أما القديسين، فهم الأنحيان الذين يعتقد أنهم قد تحققوا درجة عالية من التقوى والقداسة في الحياة الدنيا ويعتبرون مثالاً يعتدى به للمؤمنين. في الميحية، يتم تكريم القديسين والقديسات الذين يعتقد أنهم قد حققوا هذه الدرجة من القداسة ويعتبرون وسطاء بين الله والبشر. أما الأماكن المقدسة، فهي الأماكن التي تحم ولم الأماكن المقدسة، الفهم المومية، ولمع الدرجة من القداسة ويعتبرون وسطاء بين الله والبشر. أما الأماكن المقدسة (معمد الموعنية)، ولمائمين الكمبة في مكة المرمة (الإسلام)، الحائط الغربي في القدس (اليهودية)، الكنيسة القديسة العائدة في بيت لحم (السيحية)، والسجد الأملين.

(The Holy Scriptures are books that are believed to be inspired by God and carry His message to humanity. In the three monotheistic religions (Judaism, Christianity, and Islam), the holy books are: Judaism: The Torah (Old Testament), which includes the first five books of the Jewish Bible and includes Genesis, Exodus, Leviticus, Numbers, and Deuteronomy. Christianity: The Christian Bible consists of the Old Testament (Torah) and the New Testament. The New Testament includes the four Gospels (Matthew, Mark, Luke, and John), the Epistles of the Apostles, and the Book of Revelation. Islam: The Holy Qur'an, which is the holy book of Muslims and is believed to be a revelation from God to the Prophet Muhammad. As for the saints, they are people who are believed to have achieved a high degree of piety and holiness in this world's life and who are considered an example for believers to follow. In Christianity, male and female saints who are believed to have achieved this degree of holiness are honored and considered mediators between God and humans. As for the holy places, they are places that carry sanctity and important religious history for believers. Some famous holy sites include: the Kaaba in Mecca (Islam), the Western Wall in Jerusalem (Judaism), the Holy Returning Church in Bethlehem (Christianity), and Al-Aqsa Mosque in Jerusalem (Islam).)

User:

(Who are the famous writers and poets in history?)

AceGPT-7B-chat

هناك العديد من الكتاب والشعراء المهورين في التاريخ، ولكن للأسف، يصعب توفير قائمة شاملة لهم جيمًا بسبب تنوع الأدب والشعر في العالم. ولكن عكننا ذكر بعض الأسماء المعروفة مثل: ٢. ويليام شكسير: شاعر وكاتب مسرحي إنجليزي مثهور عرف بأعماله الخالدة مثل روميو وجولييت ومكبث. ٢. جابريل غارسيا ماركيز: كاتب وروائي وصحفي وناشط سياسي كولومبي مثهور عرف بروايتيه مئة عام من العزلة والحب في أوقات الكوليرا. ٣. جون ميلتون: شاعر وفيلسوف إنجليزي مثهور عرف بقصيرت ومكبث. ٢ بروايتيه مئة عام من العزلة والحب في أوقات الكوليرا. ٣. جون ميلتون: شاعر وفيلسوف إنجليزي مثهور عرف بقصيدته أفقدت الجنة. ٤. فروكتيوس مار: كاتب وفيلسوف فرنسي مثهور عرف بروايته الضحى البيضاء. ٥. نحيب محفوظ: روائي وكاتب قصص قصيرة مصري مثهور عرف بروايتيه بين القصرين وقصر الشوق. ٢. ليو تولستوي: روائي وكاتب قصص قصيرة روسي مثهور عرف بروايتيه حرب وسلام وأنا كارينينا. ٦. جان بول سارتر: كاتب وفيلسوف فرنسي مثهور عرف بأعماله الفلسفية والروائية مثل الوحدة. ٨. جبران خليل جبران: شاعر وفيلسوف فرنسي مثهور عرف بروايتيه حرب وسلام وأنا كارينينا. ٦. جان بول سارتر: كاتب وفيلسوف فرنسي مثهور عرف بأعماله الفلسفية والروائية مثل الوحدة. ٨. جبران خليل جبران: شاعر وفيلسوف وناقد لبناني مثهور عرف بأعماله الشعرية والناكيرينا. ٣. جان بول سارتر: كاتب وفيلسوف فرنسي مثهور عرف بروايتيه الحب في زمن الكوليرا والامي وكاتب قصص قصيرة رواي وكاتب قصص قصيرة روسي مثهور عرف بروايتيه حرب وسلام وأنا كارينينا. ٣. جان بول سارتر: كاتب وفيلسوف فرنسي مثهور عرف بروايتيه الحب في زمن الكوليرا والأسود يليق بك. ٢. بومان منان شاعر وكاتب ألماني مثهور عرف بأعماله الشعرية والنائين. ٣. محمد حسن عجمي: شاعر مصري مثهور عرف بروايتيه الحب في زمن الكوليرا والأسود يليق بك. ٢. ال

(There are many famous writers and poets in history, but unfortunately, it is difficult to provide a comprehensive list of all of them due to the diversity of literature and poetry in the world. But we can mention some Well-known names such as 1. William Shakespeare: A famous English poet and playwright knew for his immortal works such as "Romeo and Juliet" and "Macbeth". 2. Gabriel García Márquez: A famous Colombian writer, novelist, journalist, and political activist He is known for his novels "One Hundred Years of Solitude" and "Love in Times of Cholera". 3. John Milton: A famous English poet and philosopher knew for his poem "Paradise Lost". 4. Fructius Mahler: Writer And a famous French philosopher are known for his novel "The White Doha". 5. Naguib Mahfouz: A famous Egyptian novelist and short story writer known for his novels "Between the Two Palaces" and "The Palace of Longing". 6. Leo Tolstoy: A famous Russian novelist and short story writer known for his novels "War and Peace" and "Anna Karenina". 7. Jean-Paul Sartre: A famous French writer and philosopher known for his philosophical and fictional works such as "Loneliness". 8 Gibran Khalil Gibran: A famous Lebanese poet, philosopher, and critic known for his poetic and prose works such as "The Prophet". 9. Muhammad Hassan Ajami: A famous Egyptian poet known for his novels "Love in the Time of Cholera" and "Black Suits You." 10. Thomas Mann: A famous German poet and writer known for his novels "The Blind Leg" and "The Old Houses". These are just some examples, and there are many other famous writers and poets in history.)

Jais-13B-chat

هناك العديد من الكتاب والشعراء الشهورين في التاريخ ، ما في ذلك ويليام شكسير ، وجون ميتون ، وإميلي ديكنسون ، وإميلي مرنت ، ويليام بليك ، وجون دون ، وشارلوت يرونتي ، وصمويل تايلور كولريدج ، وجون كيتس ، ويدسي بيش شيلي ، وويليام ووردزورث ، وجون ستيوارت (There are many famous writers and poets in history, including William Shakespeare, John Milton, Emily Dickinson, Emily Bronte, William Blake, John Donne, Charlotte Brontë, Samuel Taylor Coleridge, John Keats, Percy Bysshe Shelley, William Wordsworth, and John Stuart)

Turbo:

هناك العديد من الكتاب والشعراء المهورين في التاريخ، ومن بينهم: ١. وليم شكمبير: كاتب إنجليزي فهير، ويعتبر أعظم كاتب في اللغة الإنجليزية وأحد أعظم كتاب المسرح في التاريخ. ٢. جوته: شاعر وكاتب ألماني، ويعتبر أحد أبرز شعراء الأدب العالي. ٣. هوميروس: شاعر يوناني قديم، ويعتبر مؤلف الأممال الأدبية الكلاسيكية الإلياذة والأوديسة. ٤. دانتي الأليغيري: شاعر إيطالي وكاتب الكوميديا الإلهية، ويعتبر أحد أعظم الأعمال الأدبية في التاريخ. ٥. شارل بودلير: شاعر فرنسي، ويعتبر أحد أبرز شعراء الحركة الرومانسية في القرن التاسع عشر. ٦. وليم وكاتب إنجليزي، ويعتبر أحد أبرز شعراء الأدبية في التاريخ. ٥. شارل بودلير: شاعر فرنسي، ويعتبر أحد أبرز شعراء الحركة الرومانسية في القرن التاسع عشر. ٦. وليم وردزوورث: شاعر إنجليزي، ويعتبر أحد أبرز شعراء الأحيال الأدبية في الترن التاسع عشر. ٢. جون ميلتون: شاعر إنجليزي، ويعتبر أحد أبرز شعراء الأركة الرومانسية في القرن التاسع عشر. ٦. وليم شاعر أمريكي، ويعتبر أحد أبرز شعراء المركة الومانسية في القرن التاسع عشر. ٣. جون ميلتون: شاعر إنجليزي، ويعتبر أحد أعظم شعراء الأدب الإلجليزي السابع عشر. ٢. ويرت فروست: شاعر أمريكي، ويعتبر أحد أبرز شعراء القرن السابع عشر. ٢. جون ميلتون: شاعر إنجليزي، ويعتبر أحد أعظم شعراء المرك شاعر أمريكي، ويعتبر أحد أبرز شعراء القرن العارين. ٩. جور بيرنارد شو: كاتب ومسرحي إيرلندي، ويعتبر أحد أعظم كتاب المسرح في القرن المابع عشر. ٢. روبرت فروست: على جائزة نوبل في الأدن العشرين. ٥. جورج بيرنارد شو: كاتب ومسرحي إيرلندي، ويعتبر أحد أعظم كتاب المسرح في القرن العشرين. ٢. خيب محفوظ: كاتب مصري وحائز (There are many famous writers and poets in history, among them: 1. William Shakespeare: a famous English writer, who is considered the greatest writer in the English language and one of the greatest playwrights in history. 2. Goethe: a German poet and writer, who is considered One of the most prominent poets of world literature. 3. Homer: An ancient Greek poet, and is considered the author of the classic literary works "The Iliad" and "The Odyssey". 4. Dante Alighieri: An Italian poet and writer of the "Divine Comedy". ", and it is considered one of the greatest literary works in history. 5. Charles Baudelaire: A French poet, and is considered one of the most prominent poets of the Romantic movement in the nineteenth century. 6. William Wordsworth: An English poet, and is considered one of the most prominent poets of the Romantic movement in the nineteenth century. 8. Robert Frost: An American poet, considered one of the most prominent poets of the twentieth century. 9. George Bernard Shaw: An Irish writer and playwright, and is considered one of the greatest playwrights of the twentieth century. 10. Naguib Mahfouz: An Egyptian writer and winner of the Nobel Prize in Literature, and is considered one of the most prominent writers of Arabic literature in the twentieth century.)