

ALLaM: A Series of Large Language Models for Arabic and English

Anonymous ACL submission

Abstract

In this work, we present ALLaM: Arabic Large Language Model, a series of large language models to support the ecosystem of Arabic Language Technologies (ALT). ALLaM is carefully trained, considering the values of *language alignment* and *transferability* of knowledge at *scale*. The models are based on an autoregressive decoder-only architecture and are pretrained on a mixture of Arabic and English texts. We illustrate how the second-language acquisition via vocabulary expansion can help steer a language model towards a new language without any major catastrophic forgetting in English. Furthermore, we highlight the effectiveness of using translation data and the process of knowledge encoding within the language model’s latent space. Finally, we show that effective alignment with human preferences can significantly enhance the performance of a large language model (LLM) compared to less aligned models of a larger scale. ALLaM achieves state-of-the-art performance in various Arabic benchmarks, including MMLU Arabic, ACVA, and Arabic Exams. Our aligned models improve both in Arabic and English from its base aligned models.

1 Introduction

Language modeling has significantly progressed from its humble origins, transitioning from fundamental probabilistic methods to complex neural priors. The foundational work by Shannon (1951) on the information theory of language laid the groundwork for predicting the next word in a sequence, which was initially tackled by Bengio et al. (2003) in neural space. The field experienced a substantial leap with the introduction of LSTMs (Hochreiter and Schmidhuber, 1997) in language model (LM) (Peters et al., 2018), which could capture longer dependencies in LMs but lacked scaling capability. The emergence of scalable and distributed architectures like Transformers (Vaswani et al., 2017), the

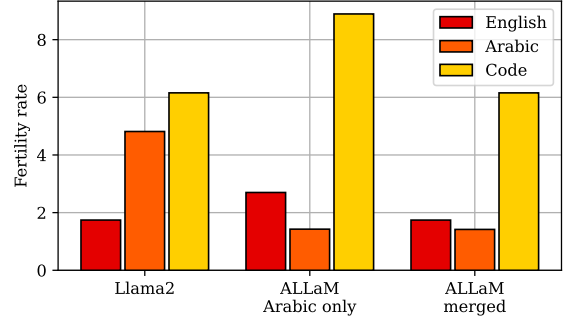


Figure 1: Comparison of Fertility Rates of LLaMa-2 and ALLaM tokenizers. The chart illustrates the fertility rates across three models: LLaMa-2, ALLaM Arabic only, and ALLaM merged with LLaMa-2 tokenizer, with datasets in English, Arabic and Code.

potential for precisely (Kaplan et al., 2020; Hoffmann et al., 2022) compressing web-scale data has resonated in recent years with the advancements of *Generative Pretraining* (Radford et al., 2018; Brown et al., 2020; Anil et al., 2023).

With the release of ChatGPT (OpenAI, 2022), followed by the introduction of more frontier class models Gemini (Google, 2024), Claude (Anthropic, 2022), Reka (Ormazabal et al., 2024), Mistral (Mistral, 2024), Llama-3 (Meta, 2024) and recently released Qwen-2 (Alibaba, 2024), generative models have experienced a significant leap from previous models (Laskar et al., 2023), raising potential implications of Artificial General Intelligence (Hendrycks and Mazeika, 2022; Marcus, 2022). This advancement has spurred discussions across various fields, including ethics, economics, and technology (Weidinger et al., 2021). Judging from the initial capabilities (Bubeck et al., 2023), the potential of these frontier models are reinventing the way humans interact with machines, impacting social norms, productivity, trends, and culture on a broader scale (Zhou et al., 2024). However, most of these frontier-class models are primarily trained on English or a few languages and often lack integration of localized regional cultures

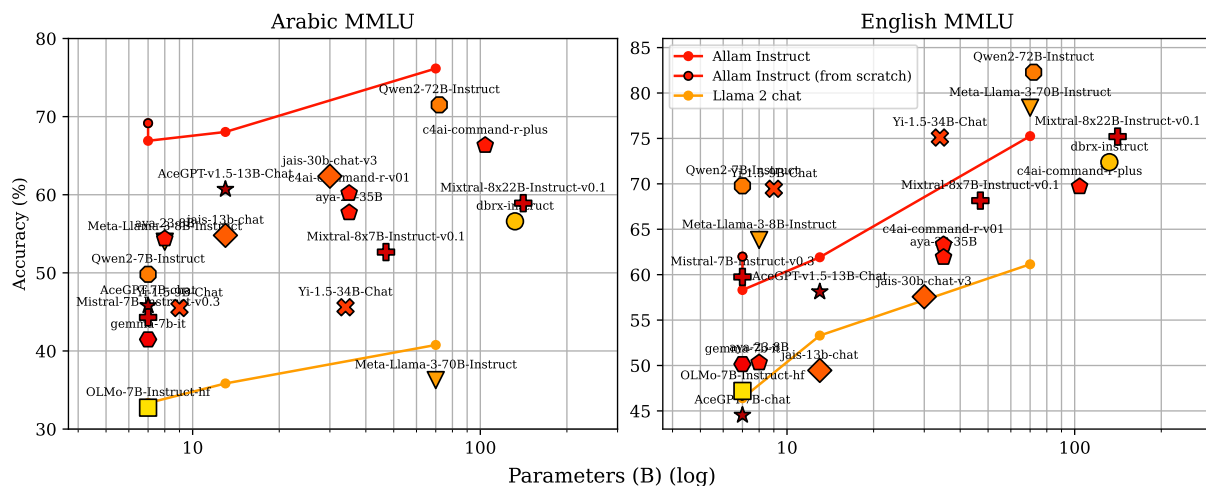


Figure 2: Performance of Various Models on Arabic (Koto et al., 2024) and English (Hendrycks et al., 2020) MMLU Benchmarks. ALLaM shows impressive improvement from it’s base model Llama 2.

and norms (Naous et al., 2024), risking *slow, irreversible manipulation* of regional identities and potentially leading to cultural homogenization.

The significant training costs of LLMs and their environmental impact have become major concerns in recent years (Strubell et al., 2019). The vast computational resources required to train LLMs contribute to substantial carbon emissions (Luccioni and Hernandez-Garcia, 2023). Governments¹ and non/for-profit organizations (Dodge et al., 2022; Google, 2021; Amazon, 2021), are increasingly aware of these issues. This awareness has led to discussions about the ethical implications of AI development and the need for sustainable practices concerning “*When and how to scale the training of these models*” To address these concerns, instead of scaling *fast*, we have opted to continue training from a *well-documented, strong*, but potentially *under-trained* pre-trained model rather than starting from a randomly initialized model. We initialize our model from **Llama 2** (Touvron et al., 2023) weights. This approach offers several key advantages that align with both our technical goals and our commitment to sustainable practices.

Technically, *continue pre-training* a model in a new language can aid in understanding **Second Language Acquisition** (SLA) (Swain and Lapkin, 1995), popularized by Bari et al. (2020) in NLP and recently adopted by Nguyen et al. (2023). This process involves the challenging task of incorporating an additional distribution without *compromising the source*. For instance, if a pre-trained model was initially trained in English, expanding to an addi-

tional language presents the specific challenge of addressing tokenization issues. Figure 1 gives an overview of ALLaM tokenizers. We expand the vocabulary of Llama 2 tokenizer from an Arabic-only ALLaM tokenizer. With the vocabulary expanded model, we continue pre-train our model for additional 1.2 Trillion tokens² on English and Arabic data mixtures and show impressive improvement over the Llama 2 base model. Finally, we apply these learnings to pre-train and align a 7B parameter model from scratch³, showing impressive improvements across the range of 7B parameter open models. In general, our contributions are listed below:

- We present the ALLaM model series, a collection of large language models developed specifically for Arabic and English languages, with the goal of supporting the cultural values of the *Arab World*. We train four models at three different scales: 7B, 13B, and 70B model initialized by Llama weights and a 7B model from scratch.
- Unlike recent trends, we explain our training methodologies and the thought process behind the decision-making involved in training the LLM. We provide necessary ablation studies for most of our crucial decisions.
- Our model achieves state-of-the-art results in Arabic as well as improving overall English performance of the original LLaMa-2 model. Check figure 2 for a quick overview.

¹<https://www.cnrs.fr/en/update/jean-zay-supercomputer-recycling-its-heat>

²For ALLaM-70B model, we train on 600B tokens

³Model initialized with random weights

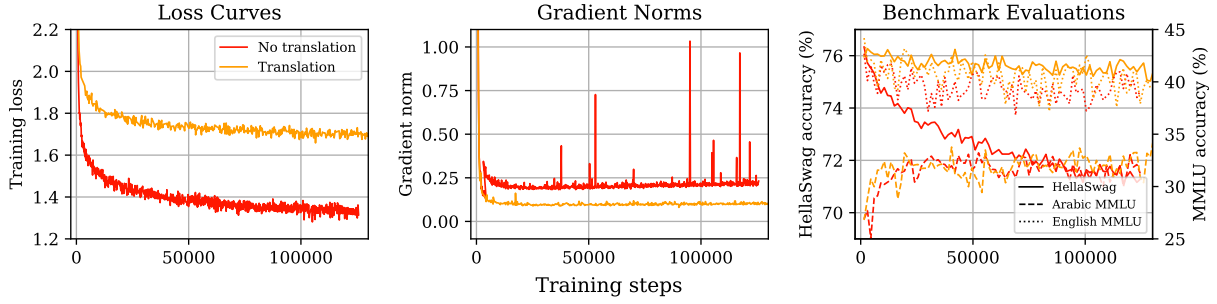


Figure 3: Effect of Arabic Translation Data in Pretraining.

2 Pretraining

Pretraining language models on trillions of natural language tokens represents the bulk of cost required to build an effective language model. This large investment of time and compute precludes experimentation or ablation for every decision. Thus, before starting to train ALLaM from random initialization, or “scratch”, we experiment in the continue-pretraining regime. As the name implies, Continue pretraining is the practice of warm-starting a pre-training experiment from an already pretrained LM.

2.1 Pretraining Data

Starting from a Llama 2 pretrained model, we continue pretraining the ALLaM-7B and ALLaM-13B models on 1.2T tokens, covering both English and Arabic languages. For the ALLaM-70B model, we train on 600B tokens. We included English data in our mixture to avoid degrading the performance of our model on English. For English, we harnessed subsets from Dolma-v1 (Soldaini et al., 2024) and Pile (Gao et al., 2021) datasets e.g., Dolma CC, The Stack (Kocetkov et al., 2022) and PeS2o, and PubMed, DM-Math (Saxton et al., 2019) and Stack-Exchange (Soboleva et al., 2023).

Our Arabic pretraining data include inhouse crawled diverse sources covering Web documents, news articles, books (literature, religion, law and culture, among others), Wikipedia (over 1M articles), and audio transcripts (books and news)⁴. To ensure high quality Web data, we applied the following processing steps: (i) Drop documents with language identification score < 95%, (ii) Drop short documents that are less than 30 words, (iii) Drop documents with duplicate URLs, high ratio of spam and stop words, (iv) Drop duplicate documents (using exact matching; although we experimented with fuzzy matching but we found it to be harsh and given that the Arabic data is scarce we

⁴We are currently working on systematic auditing of our pretraining data. Right now we do not have any timeline or visibility when or if we can share our data for research.

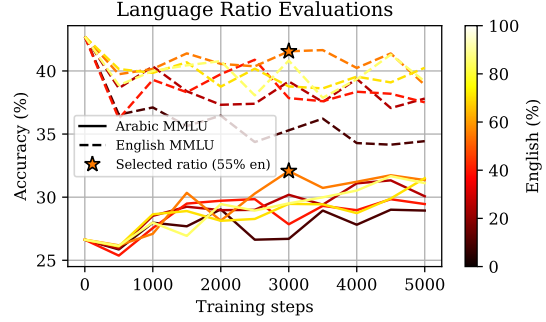


Figure 4: To find the right Arabic/English language mixture that acquires Arabic knowledge while still retaining English, we conducted an ablation over 20B tokens, in which We found that a 55/45 English/Arabic ratio achieves the best trend in performance, as measured via English and translated Arabic MMLU.

opted not to use fuzzy matching for this version).

Additionally, we extended our Arabic data with translated English content using an in-house machine translation system. We translated the following English datasets from Dolma: Wikipedia, books, C4 and peS2o, which also are part of our English data, the hypothesis is that this will improve English-Arabic language alignment, leading to a better Arabic model. Figure 3 demonstrates the impact of Arabic translation dataset in the pretraining data mixture. While models trained without translation data exhibit lower training loss, those trained with translation data show more stable training, as evidenced by fewer spikes in gradient norms. Incorporating Arabic translation data in the pretraining dataset mitigates catastrophic forgetting in English. In total, we curate 500B arabic tokens⁵.

Data Mixture. To build a performant model in both English and Arabic, we conducted experiments to figure out an optimal language mix. Fig 4 gives an overview of data-mixture experiments on our curated English-Arabic corpus. We conducted the experiments with the same sampling ratio (Table 1) and data order. We observe best trend in performance with 55 : 45 English:Arabic data mix.

⁵Token counted by our merged tokenizer.

Domain	English	Arabic		Overall
		Natural	Translated	
Web	31%	71%	65%	48%
Books	9%	13%	12%	11%
Wiki	—	0.70%	0.61%	0.3%
News	—	14%	—	3%
Science	16%	—	22%	14%
Code	39%	—	—	21%
Math	5%	—	—	2.5%
Other	—	1.3%	0.39%	0.2%
Lang Mix	55%	22.5%	22.5%	100%
Tokens	660 B	270 B	270 B	1200 B

Table 1: ALLaM Pretraining data mixture. We upsample data to match the mixture rates when needed. (Each column sums to 100%)

Table 1 shows the language and category mixing distributions for English, Arabic Natural, Arabic Translated and final mix. As depicted, and following mainstream work, Web data constitutes the highest ratio with 71%, 65% and 48% of the AR Natural, AR Translated and Final, respectively. We limited the contribution of Web English data to 31%, as Llama 2 base model was trained on Web data already and increasing its ratio might degrade performance. We ensured that high-quality sources such as books, news articles and code are well-represented in our mixture.

2.2 Continued pretraining

Open-source and open weight models present an attractive option to conduct pretraining experiments cheaply, however, they also present challenges since most such models do not natively support Arabic or other languages. We develop a simple approach to enhance any language model with capabilities in new languages (e.g. language expansion). The approach relies on two steps: (i) tokenizer augmentation and (ii) expanded vocabulary learning. We demonstrate that this approach leads to minimal degradation of capabilities in the original language.

Tokenization To calculate the fertility of our tokenizers, we subsample the entire training corpus and use this subsample as test dataset.

Existing open-weight language models (e.g., Llama 2) tokenize Arabic (and other languages) poorly, often splitting words down to the character level or even relying on byte-fallback mechanisms for tokenization. This results in inefficient training, as the pretraining corpus size is inflated, and unoptimized inference, since the model must generate more tokens per word. Additionally, the context length is reduced because it is based on a fixed number of tokens. To address these issues, we use

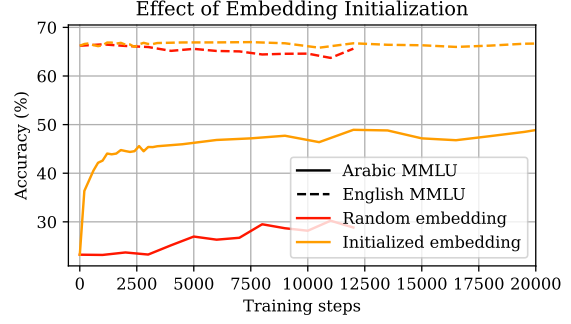


Figure 5: Effect of *Random initialization* vs *embedding initialization* during the start of continue pre-training.

a corpus of text in the target language to train a tokenizer specialized in that language. We then merge the original tokenizer with the language-specific tokenizer. Merging is accomplished by adding all tokens from the language-specific tokenizer that do not exist in the original tokenizer. As shown in Figure 1, this effectively reduces the fertility rate in the target language of the merged tokenizer to the level of the language-specific tokenizer.

Newly added tokens in the merged tokenizer have no associated embedding representations in the pretrained language model’s weights. To learn these representations, we experiment with two approaches: (i) random initialization and (ii) initialization from combined representations of tokens in the original tokenizer. Approach (ii) is accomplished by tokenizing the vocabulary of the new tokenizer using the original tokenizer. The associated representations of this tokenization are then averaged and assigned as the vector representation of the new token. Since we work with tokenizers with byte-fallback, such a tokenization is guaranteed to exist. Figure 5 provides an overview of our initialization method. Initializing the new embeddings from the combination of previously learned 2T token trained embeddings gives a significant boost to the learning of a new language. Figure 1 gives an overview of our tokenizers.

Learning rate In all of our continued pretraining experiments, we used the final learning rate of the pretrained language model (usually $3e-5$). We experimented with approaches to gradually increase the learning rate and then decay it but found limited success. Such models typically exhibited catastrophic forgetting, indicated by significant drops in performance in the source language. We also considered optimizer state warmup (as open-weight models typically do not include the optimizer states) but found this had little effect on

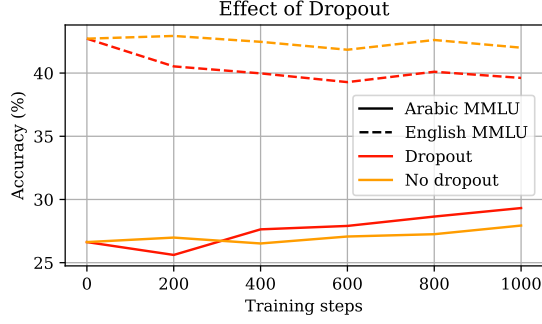


Figure 6: Effect of *Dropout* during the start of continue pre-training experiments.

performance. Figure 6 provides an overview of adding dropout during continued pretraining. We observe that adding dropout helps the Arabic language, as it acts as a regularizer for the new distribution. However, Llama 2 was pretrained on 2T tokens without any dropout, and adding dropout negatively impacts the source distribution. Considering this trade-off, we decided not to add dropout in the continue pretraining stage.

2.3 Pretraining from scratch

We were able to curate 500B Arabic tokens. Following (Hoffmann et al., 2022; Touvron et al., 2023), training a high-quality English model from scratch requires a substantial amount of tokens. Even when pretraining from random initialization, we find it beneficial to start training with a high-resource language (en) and then continue pretraining to Arabic.

In pre-training from scratch, selecting and identifying good training dynamics requires spending a lot of tokens, as different evaluations start to discriminate at different stages⁶. This may require extensive ablation studies to determine the optimal setup. Our initial experiments with 1B parameter models show that training with two languages can sometimes degrade the performance in English or result in slow learning of both language distributions. We also hypothesize that low-resource languages can dilute in the large volume of high-resource language data, even with careful tuning.

On the contrary, our continued pretraining from scratch recipe retains the natural English distribution without catastrophic forgetting, effectively transferring knowledge from one distribution to another. Judging by this trade-off, we decide to first achieve a good English distribution before applying the same approach for large-scale language alignment. The only difference here is that there is

⁶For example, a 7 billion parameter model begins to show discrepancies in MMLU at the 1 trillion tokens range.

no need for vocabulary expansion.

3 Alignment

Building effective LLMs requires ensuring they perform well and adhere to ethical standards and user expectations. This alignment process is crucial, especially for models used in diverse linguistic and cultural contexts.

Supervised Finetuning (Section 3.1) refines a pre-trained model using a carefully selected dataset relevant to specific tasks and domains. Preference training (Section 3.2), on the other hand, aligns the model’s outputs with human values and preferences by prioritizing responses that meet user expectations and ethical guidelines. Together, these methods create reliable and ethically sound LLMs for real-world use.

3.1 Supervised Finetuning Training

Data. Our Supervised Finetuning (SFT) data is curated from a diverse array of sources. For English, we primarily use public web content as our main source, offering a broad range of high-quality and especially diverse prompts. In contrast, our Arabic data comes from a combination of public and proprietary sources to ensure comprehensive coverage and relevance. We utilize classifiers, human and/or generative models (Ding et al., 2023) to *identify/interact* if the text can be considered suitable for supervised finetuning and/or if we can generate an SFT dataset from any context. To gather data from the source, we collect seed websites or data sources, which involves utilizing domain experts, prompt librarians, local institutes specializing in areas such as Arabic language, history, and politics, the use of permissible commercial LLMs to generate data, and machine translation models to convert rich English SFT data into Arabic. Our datasets cover various domains and capabilities, ensuring the model’s proficiency in handling tasks across education, history, Arabic linguistics, politics, religion, computer science, and other fields. The entire collection is named as Ultra-Instinct, which is not *human generated* rather *human driven*.

Quality Is All You Need. Unlike Zhou et al. (2023); AI et al. (2024) we hypothesized that scaling SFT data can unlock diverse capability as well as improve responsiveness to the prompts. Initially we crawled public web for supervised finetuned samples. The first version (v1) of Ultra-Instinct includes 6M samples each from English and Arabic, while the second version (v2), is a reduced version

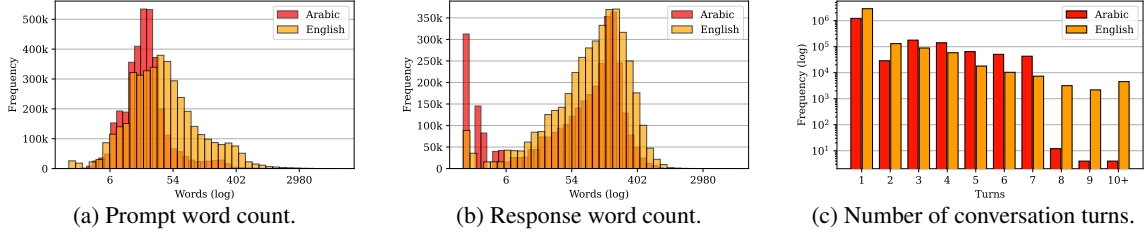


Figure 7: SFT data distributions.

Quality Metric	V1		V2	
	Prompt	Response	Prompt	Response
Word length	146.94	97.19	60.81	136.47
Lexical diversity	76.34	75.25	85.29	69.53

Table 2: Comparison of average word length and lexical diversity for (v1) and (v2) in prompts and responses.

with half the number of samples. For v1, we did not implement rigorous quality checks or extensive data removal. In contrast, v2 underwent strict quality checks and random human assessments. Our quality checks for v2 included (i) Assessments based on instruction/response word length, (ii) Lexical and semantic diversity, exact and near-exact lexical deduplication, (iii) The removal of low quality machine-translated Arabic data from English sources, and ensuring diversity in questions and commands. For detailed metrics on instruction and response lengths and lexical diversity, see Table 2.

Figure 7a and Section 3.1 shows the distribution of the prompts and responses in v2, respectively.

Version	MMLU			Exams (ar)	ACVA	ETEC
	Huang et al. (2023)	Koto et al. (2024)	en			
Ultra-Instinct v1	51.0	68.0	63.8	56.8	79.8	66.8
Ultra-Instinct v2	51.39	68.49	63.3	56.8	76.66	65.91

Table 3: Comparative results of Ultra Instinct versions, v1 and v2, across various evaluation datasets.

To extrinsically evaluate the impact of higher quality SFT data, we trained two 13B models using v1 and v2 datasets. Despite v2 containing 50% less data, both versions performed equally well on English and Arabic evaluation benchmarks. This reduction in data volume led to faster training times and reduced costs without compromising performance. Table 3 provides a detailed comparison of the 13B model results on Ultra-Instinct v1 and v2. Ultra-Instinct contains a large amount of multi-turn conversations. Figure 7c shows the distribution of “# of turn” from Ultra-Instinct.

While training the SFT model, we encountered an issue with the tokenizer. Llama 2 tokenizer was trained using sentencepiece⁷, which breaks

the beginning and end of sequence token with multiple tokens, adversely affecting long multi-turn conversations. To address this issue, we patched sentencepiece using the huggingface LlamaTokenizer wrapper. During many stages of training we saw that having 1% of noisy text (i.e., empty response) can visibly affect the model.

3.2 Preference Training

After SFT, models are able to converse in multi-turn conversations. However, they are not fully aligned with human preferences. For example, our SFT models were terse and had limited guardrails. To circumvent these issues, we performed preference tuning with human verified samples via Direct Preference Optimization (DPO) (Rafailov et al., 2024).

The DPO inputs we utilized were sourced from early model testers and a manually curated selection of domains, such as questions related to ethics or model ownership. DPO training necessitates both negative and positive output samples to train a reward model. We relied on the testers’ feedback to identify the positive outputs. In the absence of positive outputs, we generated and verified positively aligned outputs. While (Tunstall et al., 2023) utilized preference data from AI Feedback (AIF) at *scale*, we adopt a more cautious approach in creating DPO data. We generate a smaller volume of DPO data, ensuring it is fully reviewed, edited and/or re-written by humans. From our initial experiments with small toy datasets, we observed visible issues even with **0.1%** of noisy seed DPO data. However, after **scaling**, there is a possibility that the model can ignore this noisy text.

There are two approaches for generating negative outputs: (i) on-policy: we use the generations of the model we are tuning as negative outputs, and (ii) off-policy: we use another, roughly similar, model to generate the negative outputs. We did not verify that the negative outputs were worse than the positive; we ensured that the positive outputs were of the highest quality, such that they were almost always better than the negative outputs.

⁷<https://github.com/google/sentencepiece>

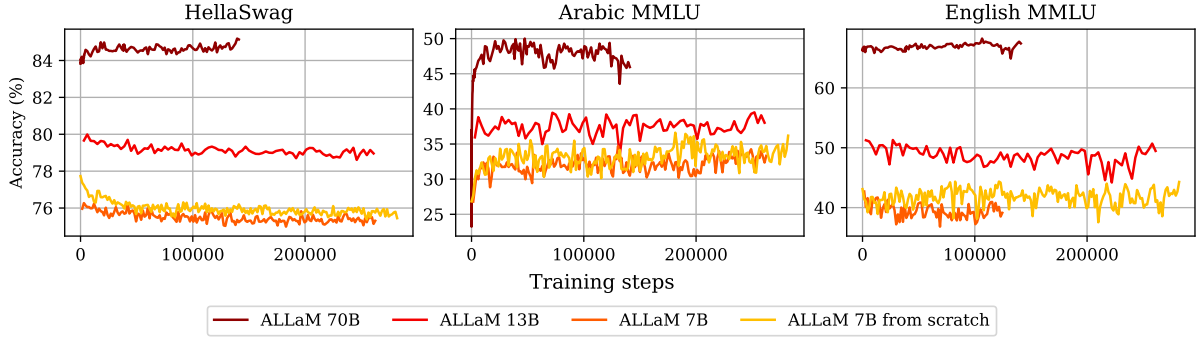


Figure 8: Benchmark evaluations throughout ALLaM model training. Using HellaSwag as a proxy for language understanding, seems that smaller models’ performance reduce when introducing Arabic, while larger models (70b) have enough capacity to improve simultaneously in English and Arabic. Arabic language acquisition is rapid in all models, as indicated by Arabic MMLU.

Compared to pretraining and SFT, the model is most sensitive to the DPO data. Therefore, we ensured the highest quality data are collected and verified. In early DPO models, we did not verify all the samples, and found that moderately noisy samples resulted in broken models that repeat generations, or output incoherent text.

4 Evaluation

In this section, we dive deeper into the evaluation of our model and report the results of our validations of ALLaM 7B, 13B and 70B models, as well as relevant models such as GPT-4, Command-R+ (Gomez, 2024), Jais-30B (Sengupta et al., 2023), and others. Our evaluation mechanisms integrate three key aspects: (i) automatic evaluations, (ii) LLM-based evaluations, (iii) human evaluations.

Limitations We start discussing evaluation by stating current limitations. Recently (Alzahrani et al., 2024) showed that multiple choice or cloze test based evaluation can be tricky and flip the benchmark. In addition to that MT-bench uses LLM as a judge and reportedly has high contamination possibility. Additionally, doing human evaluation is time consuming and requires training human evaluators. In this work, we try to ensure robust validation and attain a balanced assessment of the quantitative metrics and qualitative effectiveness and relevance of models in various applications and domains.

4.1 Automatic Evaluations

Figure 8 shows the *continuous evaluation* of our pretraining. Table 4 and 5 give an overview of the performance of ALLaM-instruct models compared to the relevant models. More detailed results can be found in Table 8, and 9. In Arabic benchmarks, we can see that ALLaM 70B scores are the best in

five (MMLU arabic both versions, Exams, ETEC, araTruthfulQA) out of the eight benchmark sets. For the remaining benchmarks: araSwag Jais 30B v3 scored the best (for this dataset, it is not publicly available but the authors shared with us the training and dev set and we are reporting on the dev set); ACVA ALLaM 7B scored the best and for araMath LLama3 70B scored the best with ALLaM 70B scoring second best. In English benchmarks, we can see a high competition between ALLaM 70B and LLama 3 70B, where LLama 3 70B scored the best in seven (MMLU, MMLU-Pro, Ethics, TruthfulQA, ARC, MixEval (hard - standard)) out of the nine benchmark sets and ALLaM 70B scoring second best in five of these (MMLU, MMLU-Pro, ARC, MixEval (hard-standard)). For the Ethics benchmark ALLaM 13B scored second best and for TruthfulQA Mistral 7B scored second best. As for the remaining two benchmarks AGIEval and HellaSwag ALLaM 70B scored the best.

4.2 LLM-based Evaluations

MT Bench (Zheng et al., 2024) consists of 80 multi-turn questions to evaluate models’ capabilities and complex instruction-following. In addition to the English version, we created an Arabic version of MT Bench developed via human translation and localization. GPT-4 serves as the LLM judge, scoring responses as recommended in (Zheng et al., 2024). Model performance is compared turn by turn, with results shown in Table 6, where ALLaM 70B achieves the best Arabic performance.

4.3 Human Evaluation

Finally we perform human evaluations to gather voting and calculate ELO scores. We developed an Arabic multi-turn dataset that covers seven domains: Arabic linguistics, history, health, politics,

			araSwag	ACVA	MMLU (ar)		Exams (ar)	ETEC	araTruthfulQA	araMath
			10-shot	5-shot	Koto et al. (2024)		5-shot	0-shot	0-shot	5-shot
					0-shot	0-shot				
ALLaM-Instruct	7B	49.28	80.33		66.9	49.6	52.7	62.95	36.4	36.5
AceGPT-Chat	7B	43.4	59.35		45.8	33.58	35.57	36.05	37.9	22.5
Llama 2-Chat	7B	24.44	52.46		33.33	26.45	25.33	26.69	29.9	21.5
Mistral-Instruct-v0.3	7B	30.59	60.7		44.3	34.06	31.1	34.41	30.3	26.0
Llama 3-Instruct	8B	33.99	75.21		53.98	41.49	44.32	49.42	34.0	38.3
ALLaM-Instruct	13B	54.77	78.59		<u>68.11</u>	51.03	<u>54.93</u>	65.59	37.5	46.8
Llama 2-Chat	13B	25.75	60.14		35.84	28.73	22.91	30.44	31.4	22.3
Jais-Chat	13B	<u>77.12</u>	70.68		54.8	41.43	46.93	48.68	31.6	25.3
ALLaM-Instruct	70B	57.91	79.01		75.92	62.23	58.47	78.38	38.4	<u>56.8</u>
Jais-Chat-v3	30B	88.37	70.05		62.37	30.15	51.21	38.53	37.3	32.5
Llama 2-Chat	70B	30.72	59.49		40.77	32.86	28.68	30.6	32.3	25.5
Llama 3-Instruct	70B	45.75	<u>80.26</u>		36.27	<u>60.11</u>	58.47	<u>71.41</u>	37.7	59.70

Table 4: Comparison of Arabic benchmarks for various instruct models.

			AGIEval	MMLU	MMLU-Pro	Ethics	TruthfulQA	ARC	HellaSwag	MixEval	
			0-shot	Average	CoT 5-shot	0-shot	0-shot	Challenge	0-shot	Hard	Standard
				0-shot				0-shot		5/0-shot (base/ft)	5/0-shot (base/ft)
ALLaM-Instruct	7B	47.09	58.31	27.78	69.8	42.11	51.45	75.2	28.9	67.6	—
AceGPT-Chat	7B	26.33	44.53	—	53.38	49.34	42.32	70.92	—	—	—
Llama 2-Chat	7B	35.55	46.4	22.87	58.88	45.32	44.28	75.52	30.8	61.7	—
Mistral-Instruct-v0.3	7B	42.22	59.75	36.33	73.59	59.65	58.7	82.88	36.2	70.0	—
Llama 3-Instruct	8B	44.35	63.82	41.32	68.07	51.72	56.83	75.81	45.6	75.0	—
ALLaM-Instruct	13B	48.42	61.8	34.05	<u>76.47</u>	57.69	55.89	81.14	37.2	72.8	—
Llama 2-Chat	13B	37.73	53.3	27.19	70.52	43.95	50.17	79.66	—	—	—
Jais-Chat	13B	31.45	49.46	—	64.92	39.66	46.84	77.6	—	—	—
ALLaM-Instruct	70B	65.67	<u>75.43</u>	<u>48.61</u>	76.16	58.78	<u>59.56</u>	84.97	<u>51.60</u>	<u>83.5</u>	—
Jais-Chat-v3	30B	36.78	57.57	26.45	68.03	42.34	51.02	78.91	—	—	—
Llama 2-Chat	70B	46.0	61.15	35.16	68.5	52.77	54.27	82.14	38.0	74.6	—
Llama 3-Instruct	70B	<u>63.78</u>	78.38	59.52	77.09	61.79	64.33	82.49	55.90	84.00	—

Table 5: Comparison of English benchmarks for various instruct models.

Model	English			Arabic		
	Avg.	Turn 1	Turn 2	Avg.	Turn 1	Turn 2
AceGPT 13B-chat	5.44	6.76	4.12	6.33	7.01	5.64
ALLaM 13B Instruct	7.34	7.67	7.01	7.57	7.9	7.23
ALLaM 70B Instruct	7.44	7.91	6.96	8.19	8.4	7.97
Jais 13B Chat	4.18	4.39	3.96	4.72	5.07	4.36
Jais 30B Chat v1	3.89	4.13	3.64	3.54	4.13	2.95
Jais 30B Chat v3	5.86	6.25	5.47	6.28	6.78	5.78
Cohere Command R+	7.41	7.63	7.18	7.97	8.28	7.65
Cohere Command R	6.99	7.19	6.79	7.47	7.82	7.12
DBRX Instruct	7.16	7.33	6.98	7.83	8.19	7.46
GPT 3.5 Turbo	7.55	7.79	7.31	8.12	8.39	7.84

Table 6: MT Bench scores for Arabic and English. The scores represent the average GPT judge score over the 80 samples ranging from 0 to 10.

coding, entertainment, and ethics, each domain contains ten questions with two turns. Each comparison was evaluated by three evaluators, and we calculated the majority voting among them. In cases of disagreement, a fourth evaluator was used to break the tie. ALLaM 13B win rate was always higher than its loss rate compared with other models. Figure 9 shows the ELO scores of the human evaluations. ELO scoring had two configuration, the default scoring rewards the good model with 1 point, the tie (good and both-bad) with 0.5 points, penalizing the bad model. The custom configuration, however, penalizes the bad model and both models if both models provided bad responses. From the figure, GPT-4 achieved the highest score, followed by ALLaM 13B with the second highest score, outperforming (or matching) larger models such as CommandR+.

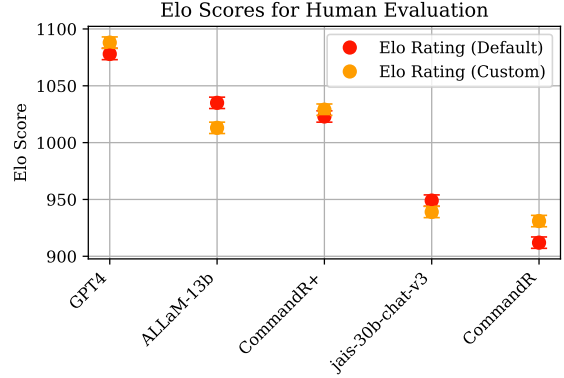


Figure 9: ELO Scores for Human Evaluation Across Various Models

5 Conclusion

ALLaM model series mark a significant leap in Arabic Language Technologies (ALT) by achieving state-of-the-art performance across various Arabic benchmarks and enhancing English performance. Through careful training that emphasizes language alignment and transferability, our models demonstrate effective second-language acquisition without catastrophic forgetting. The strategic use of translation data, knowledge encoding, and alignment with human preferences have been crucial in this success. Our openly available models on the [redacted](#) aim to support and enrich the cultural and technological landscape of the Arab World, fostering further advancements in LLMs.

6 Limitations

The model was trained on data that may potentially include toxic language, unsafe content, and societal biases originally sourced from the internet, leading to the possible amplification of these biases and toxic responses, particularly when prompted with toxic inputs. Although the model underwent concise safety training during the alignment phase, more community feedback is needed to iteratively improve the model. Additionally, inherent uncertainties in generative models mean that trials cannot encompass every possible use case, making it impossible to predict the model’s responses in all contexts. This can occasionally result in inaccurate, biased, or socially unacceptable outputs, even if the prompt itself is not explicitly offensive. Developers must conduct thorough safety evaluations and make specific adjustments to ensure the model is suitable for its intended purposes. Furthermore, the output generated by this model should not be considered a statement from the model’s creators or any affiliated organization.

7 Ethical Statement

While conducting and presenting this research, we are committed to upholding the highest ethical standards. We recognize the potential impact of large language models on society and the importance of ensuring their responsible development and deployment. Our work adheres to principles of fairness, transparency, and inclusivity, striving to mitigate biases and ensure diverse representation in our training data. We are mindful of privacy concerns and have taken steps to anonymize and secure data used in our research. Additionally, we acknowledge the potential for misuse of language technologies and advocate for their ethical application, promoting beneficial use cases while being vigilant about unintended consequences. Our models are made openly available to foster collaboration and further research, with the aim of contributing positively to the advancement of language technologies and supporting the cultural and technological growth of the Arabic-speaking world.

8 Risk Statement

The deployment and use of LLMs in various applications pose significant risks, including data privacy and security concerns due to the inadvertent inclusion of sensitive information in training datasets. LLMs can perpetuate or amplify biases, resulting

in unfair treatment and discrimination in critical decision-making processes. They can also generate convincing but inaccurate content, spreading misinformation and potentially influencing public opinion negatively. Over-reliance on LLMs may diminish human judgment, and the models’ susceptibility to adversarial attacks can compromise system integrity. To mitigate these risks, we follow robust governance, continuous monitoring, and iterative improvements. We also adhere to best practices in data handling and model training, fostering transparency and accountability in LLM development.

References

- Muhammad Abdul-Mageed, Abdelrahim Elmadany, Alcides Inciarte, Md Tawkat Islam Khondaker, et al. 2023. Jasmine: Arabic gpt models for few-shot learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16721–16744.
01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. [Yi: Open foundation models by 01.ai](#). *Preprint*, arXiv:2403.04652.
- 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*.
- Reem Alghamdi, Zhenwen Liang, and Xiangliang Zhang. 2022. [ArMATH: a dataset for solving Arabic math word problems](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 351–362, Marseille, France. European Language Resources Association.
- Alibaba. 2024. Qwen2 technical report.
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. 2024. [When benchmarks are targets: Revealing the sensitivity of large language model leaderboards](#). *Preprint*, arXiv:2402.01781.
- Amazon. 2021. [Sustainability in the cloud](#).
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng

637	Chen, Eric Chu, Jonathan H. Clark, Laurent El	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	697
638	Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	698
639	rav Mishra, Erica Moreira, Mark Omernick, Kevin	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	699
640	Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao,	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	700
641	Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez	Gretchen Krueger, Tom Henighan, Rewon Child,	701
642	Abrego, Junwhan Ahn, Jacob Austin, Paul Barham,	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	702
643	Jan Botha, James Bradbury, Siddhartha Brahma,	Clemens Winter, Christopher Hesse, Mark Chen,	703
644	Kevin Brooks, Michele Catasta, Yong Cheng, Colin	Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin	704
645	Cherry, Christopher A. Choquette-Choo, Aakanksha	Chess, Jack Clark, Christopher Berner, Sam Mc-	705
646	Chowdhery, Clément Crepy, Shachi Dave, Mostafa	Candlish, Alec Radford, Ilya Sutskever, and Dario	706
647	Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz,	Amodei. 2020. Language models are few-shot learn-	707
648	Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu	ers . <i>Preprint</i> , arXiv:2005.14165.	708
649	Feng, Vlad Fienber, Markus Freitag, Xavier Gar-	Sébastien Bubeck, Varun Chandrasekaran, Ronen El-	709
650	cía, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-	dan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Pe-	710
651	Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua	ter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg,	711
652	Howland, Andrea Hu, Jeffrey Hui, Jeremy Hur-	Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro,	712
653	witz, Michael Isard, Abe Ittycheriah, Matthew Jagiel-	and Yi Zhang. 2023. Sparks of artificial general in-	713
654	ski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun,	telligence: Early experiments with gpt-4 . <i>Preprint</i> ,	714
655	Sneha Kudugunta, Chang Lan, Katherine Lee, Ben-	arXiv:2303.12712.	715
656	jamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li,	Christopher Clark, Kenton Lee, Ming-Wei Chang,	716
657	Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu,	Tom Kwiatkowski, Michael Collins, and Kristina	717
658	Frederick Liu, Marcello Maggioni, Aroma Mahendru,	Toutanova. 2019. BoolQ: Exploring the surprising	718
659	Joshua Maynez, Vedant Misra, Maysam Moussalem,	difficulty of natural yes/no questions. In <i>Proceedings</i>	719
660	Zachary Nado, John Nham, Eric Ni, Andrew Nys-	of NAACL-HLT 2019 .	720
661	trom, Alicia Parrish, Marie Pellat, Martin Polacek,	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	721
662	Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif,	Ashish Sabharwal, Carissa Schoenick, and Oyvind	722
663	Bryan Richter, Parker Riley, Alex Castro Ros, Au-	Tafjord. 2018. Think you have solved question an-	723
664	rko Roy, Brennan Saeta, Rajkumar Samuel, Renee	swering? try arc, the ai2 reasoning challenge. <i>ArXiv</i> ,	724
665	Shelby, Ambrose Slone, Daniel Smilkov, David R.	abs/1803.05457.	725
666	So, Daniel Sohn, Simon Tokumine, Dasha Valter,	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	726
667	Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang,	Jacob Hilton, Reiichiro Nakano, Christopher Hesse,	727
668	Pidong Wang, Zirui Wang, Tao Wang, John Wiet-	and John Schulman. 2021. Training verifiers to solve	728
669	ing, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting	math word problems . <i>Preprint</i> , arXiv:2110.14168.	729
670	Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven	Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient	730
671	Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav	and effective text encoding for chinese llama and	731
672	Petrov, and Yonghui Wu. 2023. Palm 2 technical	alpaca. <i>arXiv preprint arXiv:2304.08177</i> .	732
673	report . <i>Preprint</i> , arXiv:2305.10403.	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi	733
674	Anthropic. 2022. The claude 3 model family: Opus,	Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun,	734
675	sonnet, haiku .	and Bowen Zhou. 2023. Enhancing chat language	735
676	Wissam Antoun, Fady Baly, and Hazem Hajj. 2020.	models by scaling high-quality instructional conver-	736
677	Aragpt2: Pre-trained transformer for arabic language	sations. <i>arXiv preprint arXiv:2305.14233</i> .	737
678	generation. <i>arXiv preprint arXiv:2012.15520</i> .	Jesse Dodge, Taylor Prewitt, Remi Tachet Des Combes,	738
679	M Saiful Bari, Shafiq Joty, and Prathyusha Jwalapu-	Erika Odmark, Roy Schwartz, Emma Strubell,	739
680	ram. 2020. Zero-resource cross-lingual named entity	Alexandra Sasha Luccioni, Noah A. Smith, Nicole	740
681	recognition. In <i>Proceedings of the aaai conference</i>	DeCario, and Will Buchanan. 2022. Measuring the	741
682	<i>on artificial intelligence</i> , volume 34, pages 7415–	carbon intensity of ai in cloud instances . <i>Preprint</i> ,	742
683	7423.	arXiv:2206.05229.	743
684	Loubna Ben Allal, Niklas Muennighoff, Lo-	Moussa Kamal Eddine, Nadi Tomeh, Nizar Habash,	744
685	gesh Kumar Umapathi, Ben Lipkin, and	Joseph Le Roux, and Michalis Vazirgiannis. 2022.	745
686	Leandro von Werra. 2022. A framework	Arabart: a pretrained arabic sequence-to-sequence	746
687	for the evaluation of code generation mod-	model for abstractive summarization. <i>arXiv preprint</i>	747
688	els. https://github.com/bigcode-project/	<i>arXiv:2203.10945</i> .	748
689	bigcode-evaluation-harness .	AbdelRahim Elmadany, Muhammad Abdul-Mageed,	749
690	Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and	et al. 2022. Arat5: Text-to-text transformers for ara-	750
691	Christian Janvin. 2003. A neural probabilistic lan-	bic language generation. In <i>Proceedings of the 60th</i>	751
692	guage model . <i>J. Mach. Learn. Res.</i> , 3:1137–1155.	<i>annual meeting of the association for computational</i>	752
693	Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng	<i>linguistics (Volume 1: Long papers)</i> , pages 628–647.	753
694	Gao, and Yejin Choi. 2020. Piqa: Reasoning about		
695	physical commonsense in natural language. In <i>Thirty-</i>		
696	<i>Fourth AAAI Conference on Artificial Intelligence</i> .		

754	Leo Gao, Stella Biderman, Sid Black, Laurence Gold-	localizing large language models in arabic. <i>arXiv</i>	808
755	ing, Travis Hoppe, Charles Foster, Jason Phang,	<i>preprint arXiv:2309.12053</i> .	809
756	Horace He, Anish Thite, Noa Nabeshima, Shawn		
757	Presser, and Connor Leahy. 2021. The pile: An	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke	810
758	800gb dataset of diverse text for language modeling.	Zettlemoyer. 2017. Triviaqa: A large scale distantly	811
759	<i>CoRR</i> , abs/2101.00027.	supervised challenge dataset for reading comprehen-	812
		sion. In <i>Proceedings of the 55th Annual Meeting of</i>	813
760	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman,	<i>the Association for Computational Linguistics</i> , Van-	814
761	Sid Black, Anthony DiPofi, Charles Foster, Laurence	couver, Canada. Association for Computational Lin-	815
762	Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li,	guistics.	816
763	Kyle McDonell, Niklas Muennighoff, Chris Ociepa,		
764	Jason Phang, Laria Reynolds, Hailey Schoelkopf,	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B.	817
765	Aviya Skowron, Lintang Sutawika, Eric Tang, An-	Brown, Benjamin Chess, Rewon Child, Scott Gray,	818
766	ish Thite, Ben Wang, Kevin Wang, and Andy Zou.	Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.	819
767	2023. A framework for few-shot language model	Scaling laws for neural language models . <i>Preprint</i> ,	820
768	evaluation .	arXiv:2001.08361.	821
769	Aidan Gomez. 2024. Introducing Command R+: A	Mohammad Abdullah Matin Khan, M Saiful Bari,	822
770	Scalable LLM Built for Business .	Xuan Long Do, Weishi Wang, Md Rizwan Parvez,	823
		and Shafiq Joty. 2023. xcodeeval: A large scale mul-	824
771	Google. 2021. Carbon free energy for google cloud	tilingual multitask benchmark for code understand-	825
772	regions .	ing, generation, translation and retrieval . <i>Preprint</i> ,	826
		arXiv:2303.03004.	827
773	Google. 2024. Gemini: A family of highly capable	Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia	828
774	multimodal models . <i>Preprint</i> , arXiv:2312.11805.	Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine	829
		Jernite, Margaret Mitchell, Sean Hughes, Thomas	830
775	Momchil Hardalov, Todor Mihaylov, Dimitrina	Wolf, Dzmitry Bahdanau, Leandro von Werra, and	831
776	Zlatkova, Yoan Dinkov, Ivan Koychev, and Preslav	Harm de Vries. 2022. The stack: 3 tb of permissively	832
777	Nakov. 2020. Exams: A multi-subject high	licensed source code . <i>Preprint</i> , arXiv:2211.15533.	833
778	school examinations dataset for cross-lingual and		
779	multilingual question answering. <i>arXiv preprint</i>	Fajri Koto, Haonan Li, Sara Shatnawi, Jad Dough-	834
780	<i>arXiv:2011.03080</i> .	man, Abdelrahman Boda Sadallah, Aisha Alraeasi,	835
		Khalid Almubarak, Zaid Alyafeai, Neha Sengupta,	836
781	Dan Hendrycks, Collin Burns, Steven Basart, Andrew	Shady Shehata, et al. 2024. Arabicmmlu: Assessing	837
782	Critch, Jerry Li, Dawn Song, and Jacob Steinhardt.	massive multitask language understanding in arabic.	838
783	2021a. Aligning ai with shared human values. <i>Pro-</i>	<i>arXiv preprint arXiv:2402.12840</i> .	839
784	<i>ceedings of the International Conference on Learning</i>		
785	<i>Representations (ICLR)</i> .	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	840
		field, Michael Collins, Ankur Parikh, Chris Alberti,	841
786	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	Danielle Epstein, Illia Polosukhin, Matthew Kelcey,	842
787	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	Jacob Devlin, Kenton Lee, Kristina N. Toutanova,	843
788	2020. Measuring massive multitask language under-	Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob	844
789	standing. <i>arXiv preprint arXiv:2009.03300</i> .	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natu-	845
		ral questions: a benchmark for question answering	846
790	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	research. <i>Transactions of the Association of Compu-</i>	847
791	Arora, Steven Basart, Eric Tang, Dawn Song, and	<i>tational Linguistics</i> .	848
792	Jacob Steinhardt. 2021b. Measuring mathematical		
793	problem solving with the math dataset. <i>NeurIPS</i> .	Imad Lakim, Ebtesam Almazrouei, Ibrahim Abualhaol,	849
		Merouane Debbah, and Julien Launay. 2022. A holis-	850
794	Dan Hendrycks and Mantas Mazeika. 2022. X-risk	tic assessment of the carbon footprint of noor, a very	851
795	analysis for ai research . <i>Preprint</i> , arXiv:2206.05862.	large arabic language model. In <i>Proceedings of Big-</i>	852
		<i>Science Episode# 5–Workshop on Challenges & Per-</i>	853
796	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long	<i>spectives in Creating Large Language Models</i> , pages	854
797	short-term memory . <i>Neural computation</i> , 9:1735–	84–94.	855
798	80.		
799	Jordan Hoffmann, Sebastian Borgeaud, Arthur Men-	Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur	856
800	sch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-	Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty,	857
801	ford, Diego de Las Casas, Lisa Anne Hendricks,	and Jimmy Xiangji Huang. 2023. A systematic study	858
802	Johannes Welbl, Aidan Clark, et al. 2022. Train-	and comprehensive evaluation of chatgpt on bench-	859
803	ing compute-optimal large language models . <i>arXiv</i>	mark datasets . <i>Preprint</i> , arXiv:2305.18486.	860
804	<i>preprint arXiv:2203.15556</i> .		
		Aitor Lewkowycz, Anders Andreassen, David Dohan,	861
805	Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao	Ethan Dyer, Henryk Michalewski, Vinay Ramasesh,	862
806	Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen	Ambrose Slone, Cem Anil, Imanol Schlag, Theo	863
807	Alharthi, Bang An, Ziche Liu, et al. 2023. Acegpt,	Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy	864

865	Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models .	918
866		919
867	Peiqin Lin, Shaoxiong Ji, Jörg Tiedemann, André FT Martins, and Hinrich Schütze. 2024. Mala-500: Massive language adaptation of large language models. <i>arXiv preprint arXiv:2401.13303</i> .	920
868		921
869	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	922
870		923
871		924
872		925
873		926
874		927
875		928
876		929
877	Alexandra Sasha Luccioni and Alex Hernandez-Garcia. 2023. Counting carbon: A survey of factors influencing the emissions of machine learning . <i>Preprint</i> , arXiv:2302.08476.	930
878		931
879		932
880		933
881	Gary Marcus. 2022. Is chatgpt really a “code red” for google search?	934
882		935
883	Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date .	936
884		937
885	Mistral. 2024. Au large .	938
886	El Moatez Billah Nagoudi, Muhammad Abdul-Mageed, AbdelRahim Elmadany, Alcides Alcoba Inciarte, and Md Tawkat Islam Khondaker. 2022. Jasmine: Arabic gpt models for few-shot learning. <i>arXiv preprint arXiv:2212.10755</i> .	939
887		940
888		941
889		942
890		943
891	Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei Xu. 2024. Having beer after prayer? measuring cultural bias in large language models . <i>Preprint</i> , arXiv:2305.14456.	944
892		945
893		946
894		947
895		948
896	Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. 2023. Seallms – large language models for southeast asia . <i>Preprint</i> , arXiv:2312.00738.	949
897		950
898		951
899		952
900		953
901	Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and Yang You. 2024. Mixeval: Deriving wisdom of the crowd from llm benchmark mixtures . <i>Preprint</i> , arXiv:2406.06565.	954
902		955
903		956
904		957
905		958
906	OpenAI. 2022. Chatgpt: Optimizing language models for dialogue .	959
907		960
908		961
909	Aitor Ormazabal, Che Zheng, Cyprien de Masson d’Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugénie Lamprecht, Hai Pham, Isaac Ong, Kaloyan Aleksiev, Lei Li, Matthew Henderson, Max Bain, Mikel Artetxe, Nishant Relan, Piotr Padlewski, Qi Liu, Ren Chen, Samuel Phua, Yazheng Yang, Yi Tay, Yuqi Wang, Zhongkai Zhu, and Zhihui Xie. 2024. Reka core, flash, and edge: A series of powerful multimodal language models . <i>Preprint</i> , arXiv:2404.12387.	962
910		963
911		964
912		965
913		966
914		967
915		968
916		969
917		970
	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations . <i>Preprint</i> , arXiv:1802.05365.	971
		972
	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.	
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Szymon Ruciński. 2024. Efficient language adaptive pre-training: Extending state-of-the-art large language models for polish. <i>arXiv preprint arXiv:2402.09759</i> .	
	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale. <i>arXiv preprint arXiv:1907.10641</i> .	
	David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. 2019. Analysing mathematical reasoning abilities of neural models . <i>Preprint</i> , arXiv:1904.01557.	
	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. <i>arXiv preprint arXiv:2308.16149</i> .	
	Claude E Shannon. 1951. Prediction and entropy of printed english. <i>Bell system technical journal</i> , 30(1):50–64.	
	Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama .	
	Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. Dolma: an open corpus of three trillion tokens for language model pretraining research . <i>CoRR</i> , abs/2402.00159.	
	Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. Energy and policy considerations for deep learning in nlp . <i>Preprint</i> , arXiv:1906.02243.	

973	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. <i>arXiv preprint arXiv:2210.09261</i> .	1031
974		1032
975		1033
976		1034
977		1035
978		
979	Merrill Swain and Sharon Lapkin. 1995. Problems in Output and the Cognitive Processes They Generate: A Step Towards Second Language Learning . <i>Applied Linguistics</i> , 16(3):371–391.	1036
980		1037
981		1038
982		1039
983	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashii Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	1040
984		1041
985		1042
986		1043
987		1044
988		1045
989		1046
990		
991		1047
992		1048
993		1049
994		1050
995		1051
996		
997		1052
998		1053
999		1054
1000		1055
1001		1056
1002		
1003		
1004		
1005		
1006	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Cl��mentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of lm alignment. <i>arXiv preprint arXiv:2310.16944</i> .	
1007		
1008		
1009		
1010		
1011		
1012	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . <i>Preprint</i> , arXiv:1706.03762.	
1013		
1014		
1015		
1016	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. <i>arXiv preprint arXiv:2406.01574</i> .	
1017		
1018		
1019		
1020		
1021		
1022	Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2021. Ethical and social risks of harm from language models . <i>Preprint</i> , arXiv:2112.04359.	
1023		
1024		
1025		
1026		
1027		
1028		
1029		
1030		
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> .	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. Agieval: A human-centric benchmark for evaluating foundation models . <i>Preprint</i> , arXiv:2304.06364.	
	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment . <i>Preprint</i> , arXiv:2305.11206.	
	Xuhui Zhou, Zhe Su, Tiwalayo Eisape, Hyunwoo Kim, and Maarten Sap. 2024. Is this the real life? is this just fantasy? the misleading success of simulating social interactions with llms . <i>Preprint</i> , arXiv:2403.05020.	

1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090

A Appendix	
Contents	
1 Introduction	1
2 Pretraining	3
2.1 Pretraining Data	3
2.2 Continued pretraining	4
2.3 Pretraining from scratch	5
3 Alignment	5
3.1 Supervised Finetuning Training . .	5
3.2 Preference Training	6
4 Evaluation	7
4.1 Automatic Evaluations	7
4.2 LLM-based Evaluations	7
4.3 Human Evaluation	7
5 Conclusion	8
6 Limitations	9
7 Ethical Statement	9
8 Risk Statement	9
A Appendix	14
B Related Work	14
C Alignment Details	15
C.1 SFT Data Details.	15
C.2 SFT Training Details.	15
C.3 DPO Training Details.	15
C.4 DPO vs PPO	15
D Evaluation Details	15
D.1 Human Evaluations	15
D.2 Automatic Evaluation Frameworks	16
D.3 Dataset List	16
D.4 Detailed Results	16
E Intended Use	16
F Writing Help	16
G Computational Budget and Infra	16
H Training Framework	19

B Related Work	1091
The most prominent Arabic-focused LLMs are:	1092
1. Jais (Sengupta et al., 2023): 13b and 30b base and chat models trained from scratch using a combination of natural and translated Arabic data along with English and code data.	1093 1094 1095 1096
2. AceGPT (Huang et al., 2023): 7b and 13b base and chat models trained from Llama 2 <i>without</i> vocabulary expansion. They also highlight the dangers of using translated data on LLM localization.	1097 1098 1099 1100 1101
While Jais and AceGPT are the most prominent ones, early open models such as AraGPT (Antoun et al., 2020), AraT5 (Elmadany et al., 2022), AraBART (Eddine et al., 2022), and Noon ⁸ models that utilized limited resources to serve Arabic and fueled the ambition to pursue Arabic focused models.	1102 1103 1104 1105 1106 1107 1108
Other closed models such as Noor (Lakim et al., 2022), Jasmine (Abdul-Mageed et al., 2023), and Aramus (Alghamdi et al., 2023) are worth mentioning to show the interest in serving a language with over 400 million speakers worldwide.	1109 1110 1111 1112 1113
Language adaptation of open models to other languages has been investigated in many research papers, some focus on languages written Latin scripts, which lessens the need for vocabulary expansion, such as Polish (Ruciński, 2024), in their work they adapted Mistral 7B. Mala-500 is another effort to expand to 534 languages, they expanded the vocabulary to 260K tokens, and further pretrained Llama 2 using LoRA adaptors (Lin et al., 2024), they used significantly less data for each language, and the evaluation of the approach was limited to measuring perplexity, and automatic classification benchmarks. (Cui et al., 2023) introduced Chinese Language adaptation of Llama and Alpaca models, where the vocabulary was increased to 50K tokens, then continued to pretrain the models and finally finetune them.	1114 1115 1116 1117 1118 1119 1120 1121 1122 1123 1124 1125 1126 1127 1128 1129 1130
It is worth noting that low-resource no longer means low in data, more significantly it also means low in compute, our work has certainly benefited from the open-source community, and our direction is to provide ALLaM models as open-source pending final checks and approvals.	1131 1132 1133 1134 1135 1136

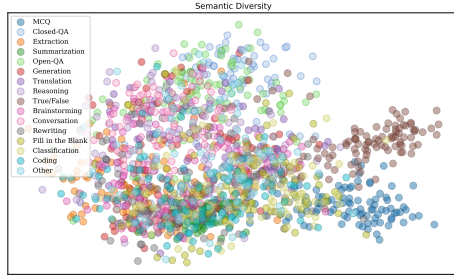


Figure 10: The semantic diversity of the prompts capabilities in Arabic v2 SFT data.

C Alignment Details

C.1 SFT Data Details.

In SFT data, we ensured that the prompts cover a sufficiently diverse embedding space, Figure 10 shows the diversity in capabilities for Arabic (v2) SFT data.

C.2 SFT Training Details.

We trained our base model, which was trained on 3.2 trillion (2T Llama + 1.2T ALLaM) tokens, for 3 epochs using Ultra-Instinct-v2 with a learning rate of $5e-6$ and a batch size of 1024. For assistant training, the model was not supposed to generate the prompt; therefore, we masked out our prompt tokens when calculating the loss. Ultra-Instinct-v2 contains a substantial number of multi-turn conversations. To train on these multi-turn conversations, we performed turn-augmentation. Figure 11 visually explains the process of turn augmentation.

C.3 DPO Training Details.

For DPO, we used 512 batch size with $KL_{penalty}$ and learning rate $9e-7$ decayed to $5e-7$ using Cosine Annealing learning rate scheduler.

Khan et al. (2023) demonstrated that model outputs can vary significantly depending on the sampling mechanism used. Building on this insight, we generate 10 additional samples for each instance by employing different temperature and nucleus sampling techniques. These additional samples are utilized to produce rejected samples, ensuring that our model provides more grounded responses and generalizes well across various sampling mechanisms. We then train the model for a single epoch using all the generated samples.

⁸<https://huggingface.co/Naseej/noon-7b>

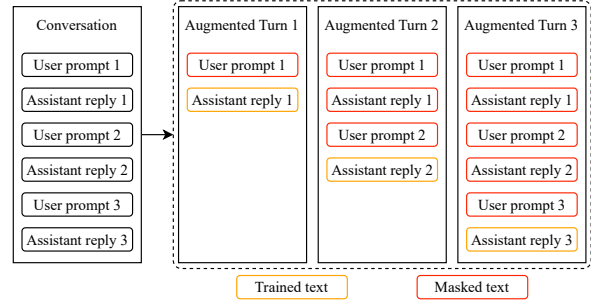


Figure 11: Augmentation process for conversations: The original conversation (left) is expanded into multiple turns (right), with user prompts and assistant replies marked for training (red) and masking (orange) to enhance the model’s language understanding and multi-turn response generation capabilities.

Model	Elo rating	
	Default	Custom
GPT-4	1,078	1,088
ALLaM Instruct 13b	1,035	1,013
Command R+	1,026	1,029
Jais 30b Chat v3	949	939
Command R	912	931

Table 7: Elo rating from human evaluations on Arabic prompts.)

C.4 DPO vs PPO

One of the fundamental differences between DPO and PPO is that PPO is always on-policy with an external *Reward Model*. In our experience with DPO, we did not encounter any **significant issues** with off-policy experiments. Additionally, DPO allows for faster iteration and easier understanding of the training dynamics. The decision to use DPO over PPO was based on logistical constraints rather than a performance comparison of the algorithms. Given our compute setup and time constraints, we chose to proceed with DPO. We plan to explore PPO in future iterations of our alignment efforts.

D Evaluation Details

D.1 Human Evaluations

In the human evaluation, we have presented two models to the human for the multi-turn dataset, and we asked the human to provide their rating with a consent that their preference (without any personal identification) will be used to further improve the language model. The shared instruction with the evaluators were the following: choose model x as the winner if it had the best answer, tie if both models good, and both-bad if both models had bad

responses. The model’s response was considered good if it was answering the questions correctly, have a coherent and natural language, is grammatically correct, the response is in the right language (if asked in Arabic the response should be in the same language except if the question was specifying answering in other languages.), and if the answer is aligned with human values and don’t contain hate or any bias. The Elo scores for human evaluations are detailed in Table 7.

D.2 Automatic Evaluation Frameworks

Most evaluations were completed using the Language Model Evaluation Harness framework (Gao et al., 2023) with the following exceptions: HumanEval was evaluated using BigCode Evaluation Harness (Ben Allal et al., 2022). MMLU-Pro, MixEval, and Arabic MMLU (Koto et al., 2024) were evaluated using the repositories of the dataset creators.

D.3 Dataset List

The evaluation pipeline covers Arabic and English benchmarks grouped into the categories listed below:

1. Multi-domain: MixEval (Ni et al., 2024), MMLU-Pro (Wang et al., 2024), and BBH (Suzgun et al., 2022).
2. Reasoning and Commonsense: HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), WinoGrande (Sakaguchi et al., 2019), and AraSwag (Nagoudi et al., 2022).
3. World Knowledge and Language Understanding: MMLU (Hendrycks et al., 2020), ARC Easy and Challenge (Clark et al., 2018), TriviaQA (Joshi et al., 2017), BoolQ (Clark et al., 2019), NQ Open (Kwiatkowski et al., 2019), AGIEval (Zhong et al., 2023), Exams-Ar (Hardalov et al., 2020), MMLU Arabic (tr) (Huang et al., 2023), MMLU Arabic (MBZU) (Koto et al., 2024), and ETEC (in-house curated).
4. Safety and Alignment: Hendrycks Ethics (Hendrycks et al., 2021a), ACVA (Huang et al., 2023), TruthfulQA (Lin et al., 2022), and AraTruthfulQA (in-house curated).
5. Conversation: MT Bench (Zheng et al., 2024), and Arabic domain capability dataset (in-house curated).

6. Math: Minerva MATH (Lewkowycz et al., 2022; Hendrycks et al., 2021b), GSM8K (Cobbe et al., 2021) and araMath (in-house curated).

The following benchmarks were developed and processed in-house: ETEC is a 1891 multiple choice questions covering different exams performed by the Education and Training Evaluation Commission at KSA⁹. Additionally AraMath is a subset that focuses testing the model performance on Arabic math problems, it consists of 600 test samples that were post-processed and prepared from the AraMath dataset (Alghamdi et al., 2022). The dataset AraTruthfulQA is a dataset created using similar methodology to TruthfulQA dataset. It comprise a total of 541 samples, 285 samples were translated directly from TruthfulQA using GPT-4, then it was carefully validated and aligned to Arabic culture by human labelers. Additionally, 256 questions were curated by humans to ensure their contextual relevance and cultural appropriateness. As for MT-Bench Arabic version, we have used GPT-4 to translate the original dataset then it was reviewed and aligned to Arabic culture by human evaluators.

D.4 Detailed Results

Follow Table 8 for arabic and Table 9 for English evaluation results.

E Intended Use

ALLaM is specifically designed to expedite the research and development of ALT through Large Language Models (LLM). It serves as one of the foundational elements for building product offerings as well as facilitating experimental initiatives.

F Writing Help

We prompt ALLaM-13B-Instruct to perform grammatical check of the content.

G Computational Budget and Infra

From different stage of training we had access from 128 A100 GPUs to 1024 A100 GPUs. We trained on GPU cluster with infiniband connections to enable high-speed communication between nodes. The all-reduce test on the cluster ranges around 1200-1400 Gbps (node-node interconnect (RoCE)). The entire training period of the models are estimated around 5M GPU hours.

⁹<https://etec.gov.sa/home>

		araSwag	ACVA	MMLU (ar)		Exams (ar)	ETEC	araTruthfulQA	araMath
				Koto et al. (2024)	Huang et al. (2023)				
		10-shot	5-shot	0-shot	0-shot	5-shot	0-shot	0-shot	5-shot
Pretrained									
ALLaM-Base (from scratch)	7B	52.68	68.46	44.45	36.28	42.09	41.7	29.4	25.5
ALLaM-Base	7B	51.63	66.18	41.52	34.42	38.55	36.58	29.9	11.5
AceGPT	7B	46.8	59.54	36.33	27.18	32.22	25.42	30.1	19.3
Llama 2	7B	25.62	62.93	33.61	26.64	23.09	27.85	25.7	24.8
Mistral-v0.3	7B	30.33	53.81	40.81	32.1	31.47	32.45	27.0	16.3
OLMo-1.7	7B	24.44	57.8	30.97	25.7	25.7	27.17	23.5	16.8
OLMo	7B	22.09	56.07	31.41	24.98	28.31	23.1	26.2	31.7
Qwen2	7B	40.26	78.74	52.91	47.16	46.0	55.23	29.9	51.2
Gemma	7B	25.36	54.82	46.33	26.04	22.91	25.48	24.0	39.3
Llama 3	8B	38.95	71.54	47.62	38.88	44.69	42.86	29.9	43.8
ALLaM-Base	13B	54.90	77.81	51.48	40.29	47.3	44.4	28.5	17.3
Yi-1.5	9B	28.76	61.19	46.36	34.11	34.82	40.01	24.0	44.8
AceGPT-v1.5	13B	48.89	73.47	42.24	33.18	40.6	33.56	30.3	18.8
Llama 2	13B	28.63	64.52	35.83	30.0	28.86	31.13	26.2	13.8
Jais	13B	49.28	60.76	32.2	29.23	33.33	27.96	28.7	28.5
ALLaM-Base	70B	59.35	79.67	59.21	49.34	53.82	55.97	33.5	38.7
Jais-v1	30B	54.51	68.25	37.6	32.94	43.39	34.04	29.6	19.3
Jais-v3	30B	53.86	70.49	45.19	38.31	50.28	45.61	30.5	25.2
Qwen1.5	32B	37.78	73.63	55.94	48.67	49.53	57.4	34.0	45.3
Yi-1.5	34B	32.16	65.25	42.93	36.26	33.71	36.21	23.7	52.0
Mixtral-8x7B-v0.1	47B	38.43	75.64	51.25	39.74	44.32	44.61	25.5	39.8
Llama 2	70B	34.38	51.16	44.79	37.1	37.99	39.38	26.6	32.3
Llama 3	70B	54.51	74.17	36.67	59.39	55.31	64.27	31.4	53.70
Qwen1.5	72B	44.84	76.0	61.38	54.44	54.0	62.84	34.90	51.8
Qwen2	72B	51.76	68.7	69.94	65	56.98	75.16	36	62.3
DBRX	132B	47.58	72.38	53.24	47.2	47.11	51.96	26.8	49.3
Mixtral-8x22B-v0.1	141B	45.1	77.21	53.6	45.92	48.42	53.96	29.8	51.0
Fine-tuned									
ALLaM-Instruct (from scratch)	7B	50.98	79.59	69.16	51.38	52.89	67.34	30.7	42.2
ALLaM-Instruct	7B	49.28	80.33	66.9	49.6	52.7	62.95	36.4	36.5
AceGPT-Chat	7B	43.4	59.35	45.8	33.58	35.57	36.05	37.9	22.5
Llama 2-Chat	7B	24.44	52.46	33.33	26.45	25.33	26.69	29.9	21.5
Mistral-Instruct-v0.3	7B	30.59	60.7	44.3	34.06	31.1	34.41	30.3	26.0
OLMo-Instruct	7B	25.36	58.74	32.74	26.5	24.77	27.33	29.6	36.5
Qwen2-Instruct	7B	37.78	79.3	49.82	48.07	47.3	56.18	35.1	51.3
Gemma-it	7B	25.62	58.03	41.48	23.15	22.91	23.73	34.8	37.0
Llama 3-Instruct	8B	33.99	75.21	53.98	41.49	44.32	49.42	34.0	38.3
Aya-23	8B	51.11	73.65	54.37	36.39	43.76	42.28	31.6	32.0
ALLaM-Instruct	13B	54.77	78.59	68.11	51.03	54.93	65.59	37.5	46.8
Yi-1.5-Chat	9B	29.8	67.57	45.5	36.02	31.47	43.6	28.7	47.8
AceGPT-Chat-v1.5	13B	49.41	64.93	60.7	37.92	40.04	42.81	36.4	22.5
Llama 2-Chat	13B	25.75	60.14	35.84	28.73	22.91	30.44	31.4	22.3
Jais-Chat	13B	77.12	70.68	54.8	41.43	46.93	48.68	31.6	25.3
ALLaM-Instruct	70B	57.91	79.01	75.92	62.23	58.47	78.38	38.4	56.8
Jais-Chat-v1	30B	80.52	71.14	60.4	43.99	48.6	48.52	32.9	26.0
Jais-Chat-v3	30B	88.37	70.05	62.37	30.15	51.21	38.53	37.3	32.5
Qwen1.5-Chat	32B	37.39	78.86	57.25	50.62	48.23	59.73	39.0	43.0
Yi-1.5-Chat	34B	30.85	65.96	45.6	35.47	35.2	40.22	25.3	49.8
CommandR	35B	55.42	78.34	60.19	48.38	50.65	55.44	33.8	47.2
Aya-23	35B	55.56	79.69	57.71	47.78	51.77	56.18	33.8	43.8
Mixtral-8x7B-Instruct-v0.1	47B	37.91	77.27	52.66	41.09	42.64	49.37	32.5	39.7
Llama 2-Chat	70B	30.72	59.49	40.77	32.86	28.68	30.6	32.3	25.5
Llama 3-Instruct	70B	45.75	80.26	36.27	60.11	58.47	71.41	37.7	59.70
Qwen1.5-Chat	72B	46.8	80.49	64.99	54.32	53.26	62.32	42.30	45.7
Qwen2-Instruct	72B	51.9	79.98	71.51	66.18	58.66	75.16	47.70	61.70
CommandR+	104B	59.35	80.37	66.33	52.98	52.89	62.1	37.0	50.2
DBRX-instruct	132B	45.75	76.46	56.6	46.73	48.79	53.17	30.5	48.8
Mixtral-8x22B-Instruct-v0.1	141B	43.79	76.45	58.92	46.74	49.72	55.55	35.1	46.0

Table 8: Comparison of Arabic benchmarks for Various Models.

Table 9: Comparison of English benchmarks for Various Models.

Model	STEM	MMLU			Other	Average	ELIZA	Winogrande	TrivialQA	ARC		BoolQ	Holding	NQ Open	MATH	H3H	Macro		Human Eval								
		Overall	STEM	Humanities						Overall	Overall						Science	Overall		Overall	Overall	Overall	Overall	Overall	Overall	Overall	Overall
MiniGPT-4	7B	25.46	35.41	40.04	44.73	45.02	44.30	31.31	50.91	65.43	35.36	80.58	67.42	43.52	65.49	76.36	27.03	15.05	10.15	30.39	48.3	—					
	7B	24.19	34.73	37.43	45.23	46.96	49.71	20.43	44.15	60.91	35.22	79.01	73.96	45.05	71.77	76.17	42.1	13.96	10.88	41.28	5.52	10.88	41.28	30.2	17.28	27.2	
	7B	25.32	34.32	38.83	44.12	47.06	41.23	20.56	44.15	60.96	38.96	79.11	74.58	46.25	77.74	75.97	52.31	14.99	5.4	13.87	30.63	—	—	—	—	—	
	7B	39.55	39.51	42.54	54.26	56.61	49.29	33.57	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	30.68	39.31	42.01	54.15	56.96	49.25	35.67	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	22.79	28.29	27.04	28.75	28.23	27.05	—	48.25	69.38	35.55	79.43	68.77	40.36	72.48	75.05	29.69	11.61	1.9	5.23	29.66	—	—	—	—	—	
	7B	38.0	33.89	33.67	71.14	68.49	61.16	34.34	66.09	73.03	45.49	81.23	69.51	53.34	82.81	80.72	—	—	0.15	—	—	—	—	—	—	—	
	7B	35.6	35.89	34.02	74.57	71.14	68.49	61.16	34.34	66.09	73.03	45.49	81.23	69.51	53.34	82.81	80.72	—	—	0.15	—	—	—	—	—	—	
	7B	24.14	14.12	45.78	57.67	49.49	40.68	35.36	62.03	73.82	33.32	56.36	80.79	59.25	56.77	80.13	55.33	19.36	8.02	17.47	47.89	—	—	—	—	—	
	7B	37.8	62.88	60.01	79.11	73.61	69.37	40.58	63.35	73.01	46.67	89.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	
MiniGPT-4-Base	V1.1	32.5	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2	32.5	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-1	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-2	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-3	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-4	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-5	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-6	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-7	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-8	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
MiniGPT-4-Base	7B	25.46	35.41	40.04	44.73	45.02	44.30	31.31	50.91	65.43	35.36	80.58	67.42	43.52	65.49	76.36	27.03	15.05	10.15	30.39	48.3	—					
	7B	24.19	34.73	37.43	45.23	46.96	49.71	20.43	44.15	60.91	35.22	79.01	73.96	45.05	71.77	76.17	42.1	13.96	10.88	41.28	5.52	10.88	41.28	30.2	17.28	27.2	
	7B	25.32	34.32	38.83	44.12	47.06	41.23	20.56	44.15	60.96	38.96	79.11	74.58	46.25	77.74	75.97	52.31	14.99	5.4	13.87	30.63	—	—	—	—	—	
	7B	39.55	39.51	42.54	54.26	56.61	49.29	33.57	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	30.68	39.31	42.01	54.15	56.96	49.25	35.67	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	22.79	28.29	27.04	28.75	28.23	27.05	—	48.25	69.38	35.55	79.43	68.77	40.36	72.48	75.05	29.69	11.61	1.9	5.23	29.66	—	—	—	—	—	
	7B	38.0	33.89	33.67	71.14	68.49	61.16	34.34	66.09	73.03	45.49	81.23	69.51	53.34	82.81	80.72	—	—	0.15	—	—	—	—	—	—	—	
	7B	35.6	35.89	34.02	74.57	71.14	68.49	61.16	34.34	66.09	73.03	45.49	81.23	69.51	53.34	82.81	80.72	—	—	0.15	—	—	—	—	—	—	
	7B	24.14	14.12	45.78	57.67	49.49	40.68	35.36	62.03	73.82	33.32	56.36	80.79	59.25	56.77	80.13	55.33	19.36	8.02	17.47	47.89	—	—	—	—	—	
	7B	37.8	62.88	60.01	79.11	73.61	69.37	40.58	63.35	73.01	46.67	89.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	
MiniGPT-4-Base	V1.1	32.5	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2	32.5	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-1	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-2	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-3	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-4	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-5	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-6	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-7	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
	V1.2+T5-8	33.0	37.8	37.8	52.61	52.61	52.61	35.3	60.35	63.35	35.3	80.79	79.12	54.69	85.81	77.78	54.32	17.87	29.96	61.82	70.97	—	—	—	—	—	—
MiniGPT-4-Base	7B	25.46	35.41	40.04	44.73	45.02	44.30	31.31	50.91	65.43	35.36	80.58	67.42	43.52	65.49	76.36	27.03	15.05	10.15	30.39	48.3	—					
	7B	24.19	34.73	37.43	45.23	46.96	49.71	20.43	44.15	60.91	35.22	79.01	73.96	45.05	71.77	76.17	42.1	13.96	10.88	41.28	5.52	10.88	41.28	30.2	17.28	27.2	
	7B	25.32	34.32	38.83	44.12	47.06	41.23	20.56	44.15	60.96	38.96	79.11	74.58	46.25	77.74	75.97	52.31	14.99	5.4	13.87	30.63	—	—	—	—	—	
	7B	39.55	39.51	42.54	54.26	56.61	49.29	33.57	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	30.68	39.31	42.01	54.15	56.96	49.25	35.67	66.19	69.95	35.67	86.15	75.61	45.13	88.14	87.38	55.06	11.41	13.98	47.42	11.41	13.98	47.42	30.2	—	—	
	7B	22.79	28.29	27.04	28.75	28.23	27.05	—	48.25	69.38	35.55	79.43	68.77	40.36	72.48	75.05	29.69	11.61	1.9	5.23	29.66	—	—	—	—	—	
	7B	38.0	33.89	33.67	71.14	68.49	61.16	34.34	66.09	73.03	45.49	81.23	69.51	53.34	82.81	80.72	—	—	0.15	—	—	—	—	—	—	—	
	7B	35.6	35.89	34.02	74.57	71.14	68.49	61.16	34.34	66.09	73.03	45.4															

H Training Framework

At the start of the project we forked [Megatron-LM](#) and did our own customizations. During training, we utilized data, tensor, and pipeline parallelism to efficiently manage the large-scale model computations. By leveraging these parallelism techniques, we achieved significant improvements in training speed and model scalability. Our modifications also included improving data iterators, adding metadata in the checkpoints, custom data pipelines etc. Depending on how many GPUS, nodes, batchsize, overlapping strategy and parallelism, our TFlops varies in between 120 to 167. We trained our model on bf16 mixed-precision.