

GenieBlue: Integrating both Linguistic and Multimodal Capabilities for Large Language Models on Mobile Devices

Xudong Lu^{*1,2†}, Yinghao Chen^{*1}, Renshou Wu^{*1}, Haohao Gao¹, Xi Chen¹, Xue Yang³, Xiangyu Zhao³, Aojun Zhou², Fangyuan Li¹, Yafei Wen¹, Xiaoxin Chen¹, Shuai Ren^{1‡}, Hongsheng Li²✉
¹vivo AI Lab ²CUHK MMLab ³Shanghai Jiao Tong University
 {luxudong@link, hsli@ee}.cuhk.edu.hk
 shuai.ren@vivo.com

Abstract

Recent advancements in Multimodal Large Language Models (MLLMs) have enabled their deployment on mobile devices. However, challenges persist in maintaining strong language capabilities and ensuring hardware compatibility, both of which are crucial for user experience and practical deployment efficiency. In our deployment process, we observe that existing MLLMs often face performance degradation on pure language tasks, and the current NPU platforms on smartphones do not support the MoE architecture, which is commonly used to preserve pure language capabilities during multimodal training. To address these issues, we systematically analyze methods to maintain pure language capabilities during the training of MLLMs, focusing on both training data and model architecture aspects. Based on these analyses, we propose **GenieBlue**, an efficient MLLM structural design that integrates both linguistic and multimodal capabilities for LLMs on mobile devices. GenieBlue freezes the original LLM parameters during MLLM training to maintain pure language capabilities. It acquires multimodal capabilities by duplicating specific transformer blocks for full fine-tuning and integrating lightweight LoRA modules. This approach preserves language capabilities while achieving comparable multimodal performance through extensive training. Deployed on smartphone NPUs, GenieBlue demonstrates efficiency and practicality for applications on mobile devices.

1. Introduction

Recent advancements in Large Language Models (LLMs) have significantly improved people’s daily lives [1, 29, 34, 67, 80], particularly through multimodal models (MLLMs) that seamlessly integrate information from different sources such as text, images, and videos [3, 4, 7, 12, 48, 58, 65]. As the scope of LLM and MLLM applications continues to expand, efficient deployment on smartphones is gaining

	Model	MATH	AlignBench	MT-Bench
Base LLM	Qwen2.5-3B	61.74	6.00	5.81
MLLM	InternVL2.5-4B	55.20	5.18	4.94
Drop (%)		10.59	13.67	14.97
Base LLM	Qwen2.5-3B	61.74	6.00	5.81
MLLM	Qwen2.5-VL-3B	58.92	5.38	4.72
Drop (%)		4.57	10.33	18.76
Base LLM	Qwen1.5-7B	22.02	5.40	5.77
MLLM	Wings-Qwen1.5-8B	13.96	4.86	4.56
Drop (%)		36.60	10.00	20.97
Base LLM	BlueLM-3B	38.94	5.67	5.42
MLLM	GenieBlue-3B	38.94	5.67	5.42
Drop (%)		0	0	0

Table 1. We assess the pure language capabilities of several representative MLLMs alongside their corresponding LLMs. The evaluation reveals that these MLLMs typically exhibit a performance drop exceeding 10% across all three datasets. In contrast, our proposed GenieBlue does not sacrifice any pure language ability.

increasing attention [14, 15, 52, 78, 82] due to their ability to enhance user privacy and support offline functionality.

In the practical process of deploying LLMs and MLLMs on smartphones, we inevitably face the storage and memory limitations inherent to these devices. Therefore, we aim to deploy a single model that can efficiently handle both pure language tasks and multimodal tasks simultaneously [7, 12]. Currently, various MLLMs suitable for on-device deployment have emerged, such as Qwen2.5-VL-3B [7], MiniCPM-V-2 [82], and InternVL2.5-4B [12], etc. These small models can achieve performance comparable to larger counterparts while having fewer parameters, making them ideal for on-device deployment. However, during the practical deployment of MLLMs on smartphone NPU (neural processing unit), we encounter the following issues:

Issue (1) MLLMs still cannot achieve satisfactory pure language capabilities currently:

Current MLLMs, while excelling in multimodal tasks, still perform moderately on pure language tasks, especially in subjective language tasks, where they still exhibit significant performance gaps compared to corresponding pure language models. We here carry out a pilot study

*Equal contribution ✉Corresponding author ‡Project lead †Intern at vivo.

to showcase this phenomenon. We evaluate the pure language capabilities of several representative MLLMs alongside their corresponding LLMs. Both Qwen2.5-VL-3B [7] and InternVL2.5-4B [12] are based on the Qwen2.5-3B [80] language model. Additionally, Wings [88] introduces a method to train MLLMs without causing text-only forgetting. Therefore, we also assess the NLP metrics of the provided Wings-Qwen1.5 checkpoint based on Qwen1.5-7B. We select three datasets for evaluation: MATH [31], which consists of challenging mathematical reasoning problems; AlignBench [45], a subjective dataset for evaluating LLMs’ human alignment in Chinese; and MT-Bench [91], a subjective benchmark for assessing multi-turn conversational capabilities. For the evaluation of AlignBench and MT-Bench, we leverage Google Gemini 1.5 Pro [65] as the judge LLM. As shown in Tab. 1, these MLLMs generally suffer a drop of more than 10% across the three datasets.

Remark: For the deployment of LLMs on mobile devices, we prioritize the performance of subjective language tasks. On-device models in smartphone environments frequently engage in more nuanced, subjective tasks in daily usage, such as text refinement, call summarization, etc.

Issue (2) Mainstream smartphone NPU platforms currently do not support deploying MoE structures:

Currently, model structural improvements designed to integrate both multimodal and pure language capabilities typically rely on the Mixture of Experts (MoE) architecture, e.g., CogVLM [71], Wings [88]. While the MoE architecture reduces the number of activated parameters during model inference, it still necessitates loading the entire original model into memory during initialization, which is not ideal for practical smartphone deployment given the limited memory available. As of now, NPU platforms of MediaTek and Qualcomm SoCs, e.g., MediaTek Dimensity 9400 and Qualcomm Snapdragon 8 Elite, do not support the deployment of MoE architectures.

Based on issue 1), despite extensive research into data and training methodologies for MLLMs, maintaining the pure language capability of MLLMs remains challenging. Based on issue 2), for end-side scenarios, the design of model architectures must also account for the constraints imposed by deployment environments. Inspired by these two challenges, this paper systematically analyzes how to maintain pure language capabilities during the training of MLLMs from both training data and model architecture aspects, emphasizing end-side deployment considerations.

From the training data perspective, we train LLMs using representative open-source MLLM datasets [69], consisting of 2.5M samples for pre-training and 7M for fine-tuning. Our findings reveal a significant decline in pure language capabilities. We then augment the fine-tuning dataset with an additional 2M samples of pure language data and re-

train the MLLM. This modification demonstrates moderate benefits for objective NLP tasks, but it yields only minimal improvements in subjective tasks due to the currently limited volume of high-quality training data available for human preference alignment [8, 89]. From these observations, we conclude that simply increasing training data is insufficient to address the decline in pure language capabilities at the current stage. Therefore, we explore the design of model structures considering hardware limitations of mobile NPUs. In this paper, we introduce **GenieBlue**, which integrates linguistic and multimodal capabilities for LLMs on mobile devices through efficient structural designs.

Specifically, to preserve the pure language capabilities of the original LLM, we freeze all LLM parameters during the multimodal training process. We then copy the transformer block every n th block for full parameter training and add LoRA [32] modules to the remaining blocks. During inference, we adopt a non-shared base deployment approach. In the LLM inference process, we utilize the originally frozen model. For MLLM inference, we replace the original transformer blocks (every n th block) with the fully trained ones and incorporate the trained LoRA parameters.

After extensive data training, GenieBlue achieves multimodal capabilities comparable to those of fully fine-tuned MLLMs without sacrificing any pure language capabilities. We also deploy GenieBlue on the NPU of real smartphones, demonstrating its efficiency and practicality for edge computing applications on mobile devices. The contributions of our work can be summarized as follows:

- 1) We examine the deployment of MLLMs on smartphones, identifying performance degradation in text-only tasks and highlighting the limitations of current NPU platforms that do not support the deployment of MoE models.
- 2) We analyze how to maintain pure language performance during the training of MLLMs from both the training data and model structure perspectives. Then, we introduce GenieBlue, which integrates both linguistic and multimodal capabilities for LLMs on mobile devices through efficient and more hardware-friendly model structural designs.
- 3) We train GenieBlue with large amounts of multimodal datasets, achieving multimodal capabilities comparable to fully fine-tuned MLLMs without compromising any pure language abilities. We also support the deployment of GenieBlue on actual smartphone NPUs, demonstrating efficient performance in real-world mobile environments.

2. Related Works

2.1. On-device LLMs and MLLMs

In recent years, beyond exploring scaling laws and training models with larger numbers of parameters on extensive datasets [7, 12, 29, 44], a promising research direction has emerged: enabling smaller LLMs and MLLMs to achieve performance comparable to larger models [14,

Type	#Samples	Datasets
General QA	840k	UltraFeedback [22], UltraChat [23], NoRobots [60], LIMA [93], SlimOrca [43], WizardLM-Evol-Instruct-70K [76], Llama-3-Magpie-Pro [77], Magpie-Qwen2-Pro [77], Firefly [81], Dolly [19], OpenAI-Summarize-TLDR [9], Know-Saraswati-CoT [36]
Code	360k	Code-Feedback [92], Glaive-Code-Assistant [27], XCoder-80K [73], Evol-Instruct-Code [55]
Mathematics	830k	GSM8K-Socratic [17], NuminaMath-TIR [38], NuminaMath-CoT [39], InfinityMATH[87], MathQA [2], MetaMathQA [83]

Table 2. We expand the Cambrian-7M dataset with 2M pure text data training samples, primarily sourced from the InternVL2.5 paper [12].

BlueLM-3B	#Samples	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG
MLLM Tasks	7M	74.81	68.32	74.60	55.30	62.35	67.91	60.06	66.19
	7M+2M	74.03	69.36	74.63	56.70	58.04	68.24	62.34	66.19
BlueLM-3B	#Samples	DROP	GPQA	GSM8K	MATH	MMLU	AlignBench	MT-bench	AVG
LLM Tasks	-	81.57	29.46	86.13	38.94	74.13	5.67	5.42	60.16
	7M	62.49	23.21	66.11	19.26	57.50	3.87	3.92	43.78
	7M+2M	64.67	28.80	69.90	30.60	57.67	3.84	3.92	47.03
Qwen2.5-3B	#Samples	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG
MLLM Tasks	7M	77.20	67.36	68.84	54.70	61.05	68.19	57.72	65.01
	7M+2M	76.98	68.48	64.25	56.20	62.09	69.43	55.54	64.71
Qwen2.5-3B	#Samples	DROP	GPQA	GSM8K	MATH	MMLU	AlignBench	MT-bench	AVG
LLM Tasks	-	70.82	30.30	74.75	61.74	66.31	6.00	5.81	60.29
	7M	69.38	20.71	68.54	31.46	63.46	4.61	4.54	49.29
	7M+2M	71.45	27.78	69.37	40.18	64.34	4.36	4.34	51.45

Table 3. We fully fine-tune BlueLM-V-3B from scratch (with SigLIP [86] and BlueLM-3B [52]/Qwen2.5-3B [80]) using Cambrian 2.5M pre-training data and 7M fine-tuning data. We also conduct fine-tuning by adding 2M text-only data to the Cambrian-7M fine-tuning dataset. The inclusion of text-only data does not cause obvious degradation in MLLM performance and partially improves the accuracy on objective NLP tasks, but does not help with subjective NLP tasks (#Samples denotes the number of fine-tuning data samples).

15, 33, 52, 68, 82, 94]. Small models with strong performance are more suitable for edge deployment scenarios, especially given the constraints of memory and computational resources on mobile devices like smartphones. With the advancement of this area of research, various small language models (SLMs) have been created [53], including the Qwen series models [7, 70, 79, 80], InternLM series models [12, 26, 54, 66], and the MiniCPM series models [33, 82]. In addition to exploring methods for training SLMs with high performance, recent research has also focused on how to more effectively deploy these models on edge devices [64, 78].

2.2. Language Capability Maintenance of MLLMs

The maintenance of original pure language capabilities during the training of MLLMs is a critical issue [7, 71, 88]. This is particularly significant in scenarios where memory and storage are limited on edge devices, emphasizing the importance of having a model that can efficiently handle both pure language and multimodal tasks. There are now basically two types of approaches used to maintain pure language capabilities during multimodal training. The first is to increase the amount of language data during the multimodal training process [7, 12, 48]. However, as demonstrated by the experiments in Sec. 1, the current approach provides limited assistance in restoring language capabilities.

The second approach is to carefully design the model structure [54, 71, 88]. Most existing methods utilize MoE architectures, which separate the “experts” that process text from those that handle other modal information. However, mainstream NPU platforms currently do not support the deployment of MoE structures. Recently, RL methods (e.g., DPO [59]) have been utilized to align models with human preference, such as Qwen2.5-VL [7]. However, these methods still do not fully restore the language capabilities of the model (see Tab. 1 and Tab. 8). Additionally, most mainstream MLLMs still rely on pre-training and fine-tuning strategies, such as InternVL 2.5 [12], DeepSeek-VL2 [75], Ovis [51], and LLaVA-OneVision [37]. Therefore, we discuss the pre-training and fine-tuning approach in this paper.

2.3. Benchmarks for Evaluating LLMs

The benchmarks for assessing LLMs can now be broadly categorized into two types: objective benchmarks and subjective benchmarks. Objective benchmarks are mainly designed to directly evaluate the knowledge capabilities of LLMs, encompassing areas such as general knowledge [16, 30, 74], mathematics and science [18, 31, 61], coding proficiency [5, 11], etc. Subjective benchmarks, on the other hand, are characterized by their reliance on human judgment and interpretation [45, 90], often requiring creativity and nuanced understanding rather than mere factual accuracy.

racy [40–42]. For on-device deployment (e.g., on smartphones), LLMs do not necessarily need to master complex knowledge but rather require better instruction following abilities, prioritizing a stronger ability in subjective tasks.

3. Text Capability Maintenance for MLLMs

In this section, we explore how to maintain the pure language capabilities during the training of MLLMs from both the training data (Sec. 3.1) and model structure perspectives (Sec. 3.2). Based on our analyses, we propose GenieBlue (Sec. 3.3), an efficient and hardware-friendly model structural design for MLLMs that combines both linguistic and multimodal capabilities, specifically tailored for LLMs/MLLMs on the NPUs of mobile devices.

3.1. Training Data Perspective

Approach Analysis: To preserve pure language capabilities during the MLLM training process, the most straightforward and commonly used method is to add text-only data to the MLLM’s training dataset. Currently, both InternVL2.5 [12] and Qwen2.5-VL [7] utilize this approach. However, this method presents some challenges. Firstly, it is difficult to collect a large amount of high-quality text-only instruction-tuning data, especially for subjective NLP tasks. Secondly, adding substantial amounts of text-only data during MLLM training will lead to longer training time.

Quantitative Experiments: To validate the effectiveness of this approach, we fully fine-tune an MLLM from scratch using a ViT and an LLM. Specifically, we utilize the BlueLM-V-3B architecture, which is tailored for end-side smartphone deployment, with SigLIP [86] as the ViT and BlueLM-3B [52]/Qwen2.5-3B [80] as the LLM. We follow the training recipe of Cambrian-1 [69], using the provided 2.5M alignment data for pre-training and the 7M data¹ for fine-tuning. For comparison, we add another 2M pure-text data samples to the fine-tuning dataset, primarily sourced from the InternVL2.5 paper [12], as shown in Tab. 2. We select 7 LLM benchmarks and 7 MLLM benchmarks for evaluation. For multimodal capabilities, we choose AI2D_{test} [35], ChartQA_{test} [56], DocVQA_{val} [57], OCR-Bench [46], RealWorldQA [21], ScienceQA_{val} [49] and TextVQA_{val} [63]. For pure language capabilities, we choose DROP_{val} [25], GPQA Diamond [62], GSM8K_{test} [17], MATH_{test} [31], MMLU_{test} [30], AlignBench [45] and MT-Bench [91]. The first five LLM benchmarks assess objective language capabilities, while the last two evaluate subjective language abilities. The evaluation results are shown in Tab. 3. We come across two observations:

Finding (1) Adding pure-text datasets has little impact on the MLLM performance:

¹Cambrian-7M dataset contains around 1.5M pure-text data samples.

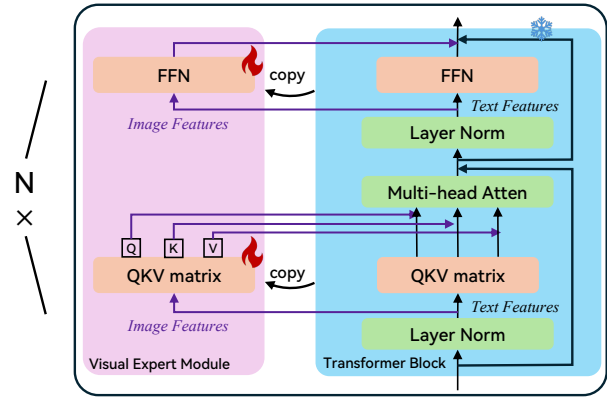


Figure 1. CogVLM [71] replicates an identical visual expert module alongside each transformer block to handle multimodal inputs.

After adding a pure language dataset containing 2M training samples, we find that the multimodal capabilities of the trained MLLM remain virtually unchanged. This phenomenon indicates that incorporating a certain amount of pure text data during the training of an MLLM does not significantly affect its multimodal performance.

Finding (2) Adding pure text data leads to a moderate improvement in the performance of objective NLP tasks but does not assist with subjective tasks:

As can be seen from Tab. 3, the incorporation of multimodal data (7M) leads to a significant decline in both the objective and subjective language performance of the original LLM. To address this issue, we refer to InternVL2.5 [12] and integrate an additional 2M pure text samples for training. As there is still a lack of sufficient high-quality open-source training data for human alignment [8], the newly added pure-text data partially restores the performance for objective NLP tasks and provides almost no help for subjective NLP tasks. This indicates that maintaining the pure-language capabilities of LLMs by adding additional pure-language data currently remains a challenging endeavor.

3.2. Model Structure Perspective

Approach Analysis: Based on the analyses in Sec. 3.1, we conclude that maintaining NLP performance during the training of MLLMs by increasing pure text data is currently challenging. Consequently, another research direction focuses on the design of MLLM architectures, aiming to enhance NLP capabilities through architectural innovations rather than solely relying on additional pure text data. Representative works in this area include CogVLM [71] and Wings [88], both of which utilize the MoE structure.

However, during our deployment journey, we still observe that Wings [88] leads to a significant decline in pure language capabilities. As noted in the experiment presented in Sec. 1, there is an average drop of over 20% in NLP performance, which is unacceptable for our deployment pur-

BlueLM-3B	#Param	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG	Retention (%)
Full-Finetune	3161.26M	74.03	69.36	74.63	56.70	58.04	68.24	62.34	66.19	-
LoRA	458.06M	68.23	61.24	66.17	48.70	55.56	68.57	56.97	60.78	91.82
CogVLM-Post	1005.69M	67.81	60.80	66.49	51.00	57.12	67.00	58.58	61.26	92.55
CogVLM-Pre	1005.69M	69.04	64.28	70.23	51.50	52.29	67.67	60.42	62.20	93.98
CogVLM-Skip	1005.69M	70.01	66.36	71.97	54.60	56.34	68.91	59.37	63.94	96.60
Qwen2.5-3B	#Param	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG	Retention (%)
Full-Finetune	3527.81M	76.98	68.48	64.25	56.20	62.09	69.43	55.54	64.71	-
LoRA	456.84M	65.35	54.32	55.84	48.10	55.56	72.72	58.40	58.61	90.58
CogVLM-Post	1146.75M	68.72	60.48	65.14	51.30	48.89	64.76	59.85	59.88	92.53
CogVLM-Pre	1146.75M	68.88	62.12	67.95	52.30	53.73	72.87	57.36	62.17	96.08
CogVLM-Skip	1146.75M	69.30	65.92	71.10	54.10	50.59	69.48	59.62	62.87	97.16

Table 4. Evaluation results on MLLM benchmarks. We fine-tune all the models using the 9M dataset, comparing full fine-tuning, LoRA fine-tuning, and CogVLM fine-tuning. **Post**, **Pre**, and **Skip** means adding the visual expert module to the last quarter of the layers, the first quarter of the layers, and at every quarter interval of the layers. Apart from full fine-tuning, other methods can maintain pure language capability consistent with the original LLM during inference through the use of the non-shared base deployment strategy. CogVLM-Skip achieves the best MLLM performance retention. We also provide the trainable parameter numbers (#Param) during MLLM training.

poses. Regarding CogVLM [71], it replicates an identical visual expert module alongside each transformer block to handle multimodal inputs while keeping the original LLM frozen during training, as shown in Fig. 1. This design ensures that the performance of the original LLM remains unchanged during inference. However, this design still has two shortcomings. **1)**, during deployment, both the LLM and all corresponding visual expert modules need to be loaded into memory simultaneously, doubling the model’s memory requirements. **2)**, as analyzed in Sec. 1, current smartphone NPU platforms do not yet support the deployment of MoE models. This results in deployment issues for CogVLM on real mobile devices.

Quantitative Experiments: To ensure the completeness of our work, we evaluate the MLLM performance of the models after training using the CogVLM approach with both BlueLM-3B and Qwen2.5-3B LLMs. To address the memory issues that arise during deployment, we integrate a visual expert module into one-quarter of the layers. We experiment with adding visual expert modules to the last quarter of the layers, the first quarter of the layers, and at every quarter interval of the layers [6]. For other transformer blocks, we add LoRA² weights to the attention modules and feed-forward modules. We compare the three CogVLM-based methods with full fine-tuning and full-LoRA training. To provide more insights, we also list the trainable parameters (including ViT and projector layer) during MLLM training. The results are shown in Tab. 4.

Finding (3) Compared to full fine-tuning, LoRA and CogVLM methods lead to a decrease in the multimodal performance of the trained MLLM:

Due to limitations in the number of trainable parameters, both LoRA and CogVLM methods fall short of the multimodal performance achieved by full fine-tuning. Nevertheless, they typically reach over 90% of the performance seen

²In all experiments, we set the LoRA rank to 8.

with full fine-tuning. Besides, CogVLM outperforms LoRA in MLLM performance. It is important to note that full fine-tuning has a significant negative impact on the performance of pure-text tasks (Tab. 3), while LoRA and CogVLM do not influence the pure language performance through the use of the non-shared base deployment strategy (Sec. 3.3).

Finding (4) For CogVLM, the addition of visual expert modules at every quarter interval of the layers results in the best MLLM performance:

Incorporating visual experts at every quarter interval of the layers results in over 96% accuracy retention for the MLLM compared to full fine-tuning. Since CogVLM’s training approach does not affect pure-text performance, we have decided to design GenieBlue based on this method.

3.3. GenieBlue

Based on the analyses from both data (Sec. 3.1) and structural (Sec. 3.2) perspectives, we propose to integrate linguistic and multimodal capabilities into the training of MLLMs through structural design. In this subsection, we provide a detailed illustration of the GenieBlue structure.

Approach Analysis: We modify from the CogVLM [71] structure, particularly paying attention to the limitations of NPUs on the MoE architecture. The main idea behind CogVLM is to separate the processing of text tokens and multimodal tokens. It employs an MoE architecture where different experts handle text and visual tokens. In contrast, our design principle focuses on bypassing the MoE structure by selecting separate model weights for LLM/MLLM deployment, thereby maintaining the original LLM architecture unchanged during the multimodal inference process.

The framework of GenieBlue is shown in Fig. 2. To save model storage on smartphones, we replicate the transformer blocks at every quarter interval throughout the layers of the LLM while integrating LoRA modules into the remaining transformer blocks. During multimodal training, we freeze

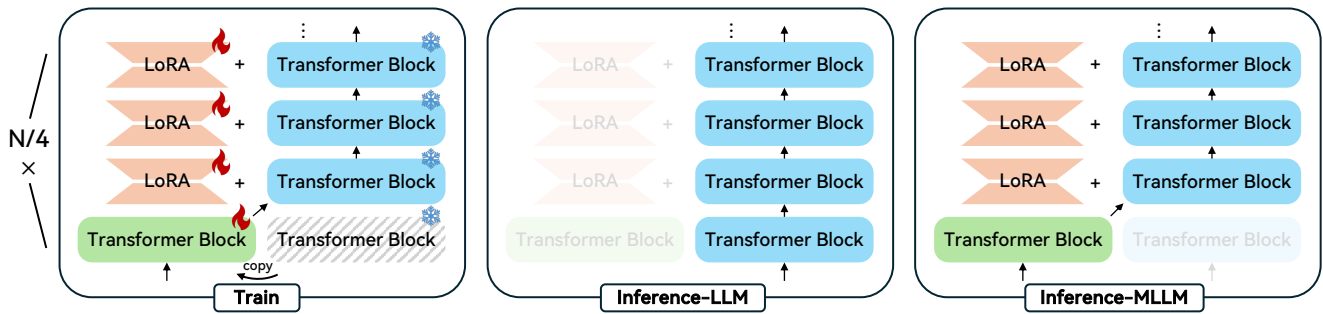


Figure 2. Overview of GenieBlue. We replicate the transformer blocks at every quarter interval of the layers in the LLM and incorporate LoRA modules into the other transformer blocks. During multimodal training, we freeze the original LLM while fully training the replicated transformer blocks and the added LoRA parameters. For pure-text inference, we utilize the original LLM. For multimodal inference, we replace the original blocks with the trained transformer blocks at every quarter interval and add LoRA to the remaining transformer blocks. This non-shared base approach avoids the MoE structure while decoupling the inference processes of the LLM and MLLM.

BlueLM-3B	#Param	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG	Retention (%)
Full-Finetune	3161.26M	74.03	69.36	74.63	56.70	58.04	68.24	62.34	66.19	-
CogVLM-Skip	1005.69M	70.01	66.36	71.97	54.60	56.34	68.91	59.37	63.94	96.60
GenieBlue-Post	1005.73M	68.49	61.68	67.78	49.80	55.42	69.96	61.59	62.10	93.82
GenieBlue-Pre	1005.73M	72.90	66.20	71.11	46.50	58.30	73.20	60.03	64.03	96.74
GenieBlue-Skip	1005.73M	73.67	69.32	74.26	55.30	57.39	68.34	60.37	65.52	98.99
Qwen2.5-3B	#Param	AI2D	ChartQA	DocVQA	OCRBench	RealWorldQA	ScienceQA	TextVQA	AVG	Retention (%)
Full-Finetune	3527.81M	76.98	68.48	64.25	56.20	62.09	69.43	55.54	64.71	-
CogVLM-Skip	1146.75M	69.30	65.92	71.10	54.10	50.59	69.48	59.62	62.87	97.16
GenieBlue-Post	1146.79M	67.29	59.80	60.70	49.30	56.47	75.35	59.88	61.26	94.66
GenieBlue-Pre	1146.79M	69.01	58.44	56.65	43.90	58.04	75.01	62.19	60.46	93.44
GenieBlue-Skip	1146.79M	72.99	63.04	62.74	53.90	57.39	71.05	61.68	63.26	97.76

Table 5. Evaluation results on MLLM benchmarks after training with the 9M fine-tuning dataset. Similar to the experiment setting of CogVLM, we replicate transformer blocks at the last, first, and every interval quarter of layers. Results show that GenieBlue-Skip demonstrates the best MLLM performance, yielding over 97% retention in MLLM performance compared to full fine-tuning.

the original LLM, allowing ViT, the replicated transformer blocks, and the added LoRA parameters to be fully trained.

For pure-text inference, we utilize the original, unmodified LLM to perform all calculations. In contrast, for multimodal inference, we replace the original blocks with the trained transformer blocks at every quarter interval and incorporate LoRA into the remaining transformer blocks. This non-shared base strategy effectively avoids the MoE structure and decouples the inference processes of the LLM and MLLM. During actual NPU deployment, we only need to replace the weights and adapt the LoRA module. This makes deployment simple and efficient.

Quantitative Experiments: We compare our proposed GenieBlue against full fine-tuning and the CogVLM methods with both BlueLM-3B and Qwen2.5-3B LLMs, using the 2.5M pre-training data and 9M fine-tuning data. For a fair comparison with CogVLM, we replicate transformer blocks at the last (Post), first (Pre), and every interval (Skip) quarter of layers. The results are shown in Tab. 5.

Finding (5) For GenieBlue structure, GenieBlue-Skip achieves the best multimodal performance, GenieBlue-Skip also outperforms CogVLM-Skip:

Similar to the results of CogVLM, replicating transformer blocks at every interval quarter of layers achieves better multimodal performance. Besides, we find that GenieBlue-Skip outperforms CogVLM-Skip. This could possibly be attributed to CogVLM’s approach of incorporating visual expert modules. In CogVLM’s design, text features and image features are rigidly separated and processed separately for QKV and FFN calculations. Although CogVLM considers the fusion of text and image features during multi-head attention, this fusion is not as effective as completely sharing weights, which limits better integration throughout the entire MLLM inference process.

Non-shared Base Deployment Strategy: By splitting the LLM and MLLM inference process, deploying GenieBlue with the non-shared base strategy (as shown in Fig. 2) can maintain the pure language capabilities of the original LLM. To validate the importance of this approach, we evaluate GenieBlue’s performance on LLM benchmarks, comparing the shared and non-shared base deployment strategies. The shared base deployment strategy refers to unifying the inference processes of LLM and MLLM into the single deployment mode depicted on the right of Fig. 2. Specifically, during the inference of pure language tasks, we also leverage the fully trained transformer blocks and incorpo-

BlueLM-3B	Shared Base	DROP	GPQA	GSM8K	MATH	MMLU	AlignBench	MT-bench	AVG	Retention (%)
BlueLM-3B	-	81.57	29.46	86.13	38.94	74.13	5.67	5.42	60.16	-
Full-Finetune	-	64.67	28.80	69.90	30.60	57.67	3.84	3.92	47.03	78.18
LoRA	✓	79.71	29.02	84.46	39.08	69.76	4.62	4.61	56.33	93.63
GenieBlue-Post	✓	78.64	28.13	85.37	37.08	70.77	4.51	4.65	55.94	92.98
GenieBlue-Pre	✓	76.95	29.24	74.98	35.66	65.26	4.61	4.71	53.61	89.12
GenieBlue-Skip	✓	75.36	29.02	76.27	38.16	67.78	4.66	4.76	54.40	90.42
GenieBlue	✗	81.57	29.46	86.13	38.94	74.13	5.67	5.42	60.16	100.00

Qwen2.5-3B	Shared Base	DROP	GPQA	GSM8K	MATH	MMLU	AlignBench	MT-bench	AVG	Retention (%)
Qwen2.5-3B	-	70.82	30.30	74.75	61.74	66.31	6.00	5.81	60.29	-
Full-Finetune	-	71.45	27.78	69.37	40.18	64.34	4.36	4.34	51.45	85.33
LoRA	✓	53.94	23.74	70.96	43.94	66.00	4.17	4.36	49.13	81.48
GenieBlue-Post	✓	60.97	20.20	72.48	43.10	64.84	4.31	4.90	50.53	83.81
GenieBlue-Pre	✓	61.65	27.27	72.25	42.94	66.99	4.45	4.69	51.79	85.90
GenieBlue-Skip	✓	67.98	28.28	69.90	42.64	65.97	4.62	4.70	52.57	87.19
GenieBlue	✗	70.82	30.30	74.75	61.74	66.31	6.00	5.81	60.29	100.00

Table 6. Comparison of pure language capabilities using the shared base versus non-shared base deployment strategies, trained with 9M fine-tuning data. The non-shared base approach can maintain the pure text capabilities of the original LLM. In the shared-base strategy, training with BlueLM-3B indicates that the fewer trainable parameters involved in multimodal training, the better the retention of pure text capabilities. However, the LoRA-trained MLLM based on Qwen2.5-3B achieves the worst pure-text performance.

rate the LoRA module. Additionally, we provide the NLP performances of BlueLM-3B/Qwen2.5-3B, the fully fine-tuned models, and the models trained entirely with LoRA. The results are shown in Tab. 6.

***Finding (6)** Deploying with the non-shared base strategy results in significantly better pure-text capabilities compared to the shared base strategy:*

Undoubtedly, using the shared base deployment strategy leads to a loss of pure language capabilities, demonstrating the importance of the non-shared base deployment method. Another interesting finding is that, intuitively, with the same MLLM training data, having fewer trainable parameters results in less loss of the model’s pure language performance. Training with BlueLM-3B aligns with this intuition. However, the LoRA-trained MLLM based on Qwen2.5-3B achieves the worst pure-text performance. A plausible explanation for this phenomenon lies in the inherent mechanism of LoRA, which imposes low-rank matrices onto original weights rather than directly training the base parameters. The limited number of adapter parameters may hinder effective integration with the pre-existing model parameters, resulting in suboptimal parameter fusion and consequently injuring the LLM performance.

4. Training and Deployment Recipe

After analyzing from both training data and model structure perspectives in Sec. 3, we determine the model structure (GenieBlue-Skip) and deployment approach (non-shared base deployment strategy). In this section, we introduce the detailed training (Sec. 4.1) and deployment details (Sec. 4.2) of the final GenieBlue model.

4.1. Training Recipe

We employ the GenieBlue-Skip structure and strictly adhere to the training recipe and training data of BlueLM-V-3B [52]. Specifically, our training process consists of two stages. In the first stage, we pre-train the MLP projection layer while keeping the ViT and LLM frozen, using the 2.5M pre-training data. In the second stage, we fine-tune the GenieBlue-Skip model (ViT, projector, replicated transformer blocks, and the added LoRA parameters) with 645M fine-tuning data [52] while keeping the original LLM frozen. We use SigLIP as the ViT and BlueLM-3B as the LLM. During training, we set the LoRA rank to 8.

4.2. Deployment Recipe

We deploy GenieBlue on the NPU of the iQOO 13 smartphone, which is equipped with the Qualcomm Snapdragon 8 Elite (Gen 4) SoC. We leverage the Qualcomm QNN SDK³ for model deployment. For the ViT and projector layer, we employ W8A16 quantization. For the LLM, we adopt W4A16 quantization. Regarding the added LoRA parameters, we utilize a W8A16 quantization scheme. Currently, we support the single-patch ViT inference. It is important to note that the Snapdragon 8 Elite’s NPU platform does not support the deployment of MoE structures.

5. Performance of GenieBlue

Through extensive data training and NPU deployment, in this section, we evaluate the MLLM (Sec. 5.1) and LLM (Sec. 5.2) capabilities of GenieBlue, as well as its deployment efficiency on smartphone NPUs (Sec. 5.3).

³<https://www.qualcomm.com/developer/software/neural-processing-sdk-for-ai>

Model	#Params	AVG	MMBench	MMStar	MMMUS	MathVista	HallusionBench	AI2D	OCRBench	MMVet
BlueLM-V-3B [52]	3.2B	66.1	82.7	62.3	45.1	60.9	48.0	85.3	82.9	61.8
Ovis2-2B [51]	2.46B	65.2	76.9	56.7	45.6	64.1	50.2	82.7	87.3	58.3
Qwen2.5-VL-3B [7]	3.75B	64.5	76.8	56.3	51.2	61.2	46.6	81.4	82.8	60.0
SAIL-VL-2B [24]	2.1B	61.0	73.7	56.5	44.1	62.8	45.9	77.4	83.1	44.2
InternVL2.5-2B-MPO [72]	2B	60.9	70.7	54.9	44.6	53.4	40.7	75.1	83.8	64.2
GenieBlue	3.2(+0.55)B	64.2	78.2	59.4	47.6	58.0	46.3	83.1	82.9	58.1
InternVL2-8B [13]	8B	64.1	79.4	61.5	51.2	58.3	45.0	83.6	79.4	54.3

Table 7. Performance on MLLM benchmarks under the same evaluation settings as OpenCompass benchmark ($\leq 4B$, with InternVL2-8B for reference). GenieBlue retains over 97% accuracy of BlueLM-V-3B while outperforming InternVL2-8B on average. [†]The total number of parameters in the replicated transformer blocks and LoRA modules is 0.55B.

	#Params	DROP	GPQA	GSM8K	MATH	MMLU	AlignBench	MT-bench	AVG	Retention (%)
BlueLM-3B	2.7B	81.57	29.46	86.13	38.94	74.13	5.67	5.42	60.16	-
GenieBlue	3.2(+0.55)B	81.57	29.46	86.13	38.94	74.13	5.67	5.42	60.16	100.00
Qwen2.5-3B	3.1B	70.82	30.30	74.75	61.74	66.31	6.00	5.81	60.29	-
Qwen2.5VL-3B	3.75B	72.72	24.24	70.43	58.92	65.07	5.38	4.72	56.05	92.98

Table 8. Evaluation results on representative LLM benchmarks, including both objective and subjective benchmarks. GenieBlue retains 100% performance of the original LLM, whereas Qwen2.5VL-3B exhibits some degradation.

Model	Context (token)	Load Time (s)	ViT Time (s)	Input Speed (token/s)	Output Speed (token/s)	Storage (GB)	Memory (GB)
BlueLM-V-3B	2048	0.51	0.4	1515.15	33.00	1.77	1.73
GenieBlue	2048	0.80	0.4	1666.67	31.00	1.92	2.10

Table 9. Deployment efficiency comparison between GenieBlue and BlueLM-V-3B on Qualcomm 8 Elite SoC in peak performance mode. GenieBlue results in a longer model loading time, slightly higher storage and memory usage, and a marginally slower token output speed.

5.1. MLLM Performance

After extensive data training, we evaluate our model using representative MLLM benchmarks, including MM-bench [47], MMStar [10], MMMU [85], MathVista [50], HallusionBench [28], AI2D [35], OCRBench [46], and MM-Vet [84], which are integrated into the OpenCompass benchmark suite [20]. We compare GenieBlue with other MLLMs that have fewer than 4B parameters, and the results are presented in Tab. 7. GenieBlue achieves MLLM accuracy slightly lower than Qwen2.5-VL-3B while retaining 97% performance of BlueLM-V-3B. Besides, GenieBlue slightly outperforms InternVL2-8B on average.

5.2. LLM Performance

The most significant feature of GenieBlue is that it does not lose LLM performance when deployed using the non-shared base deployment strategy. Here, we evaluate its LLM performance on representative benchmarks. For comparison, we select Qwen2.5VL-3B, which claims to maintain LLM performance without degradation from MLLM training by incorporating pure-text data. As demonstrated in Tab. 8, GenieBlue achieves no loss in LLM performance, while Qwen2.5VL-3B exhibits some performance degradation, especially in subjective tasks. This indicates that exploring model structure design is more effective for maintaining pure-text capabilities than simply increasing the amount of pure-text data currently.

5.3. Deployment Efficiency

We deploy GenieBlue with the non-shared base strategy on Qualcomm Snapdragon 8 Elite (Gen 4) SoC. Different from [52], we now support the 1-patch ViT inference. We here provide the MLLM deployment statistics in Tab. 9, comparing BlueLM-V-3B and GenieBlue. With the inclusion of additional LoRA parameters, GenieBlue incurs longer model loading times, slightly larger storage and memory requirements, and a marginally slower token output speed. However, a token output speed of 30 token/s is fully sufficient for daily use on mobile devices.

6. Conclusion

In this paper, we approach the challenge of maintaining pure language capabilities from a practical deployment perspective on mobile devices (smartphones), analyzing both training data and model structure to identify effective strategies. Based on the analyses, we propose GenieBlue, an efficient and hardware-friendly MLLM design that integrates linguistic and multimodal capabilities for mobile LLMs. By freezing the original LLM parameters during training and acquiring multimodal capabilities through duplicated transformer blocks and lightweight LoRA modules, GenieBlue maintains language performance while achieving competitive multimodal results. Deployed on smartphone NPUs, GenieBlue demonstrates its practicality and efficiency, making it a promising solution for edge computing applications on mobile devices. We hope that our work will provide valuable insights for future research in this field.

Acknowledgment

This project is funded in part by National Key R&D Program of China Project 2022ZD0161100, by the Centre for Perceptual and Interactive Intelligence (CPII) Ltd under the Innovation and Technology Commission (ITC)’s InnoHK, and in part by NSFC-RGC Project N_CUHK498/24.

References

- [1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. 1
- [2] Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. MathQA: Towards interpretable math word problem solving with operation-based formalisms. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. 3
- [3] Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 1, 2023. 1
- [4] Anthropic. Claude 3. <https://www.anthropic.com>, 2023. Large Language Model. 1
- [5] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021. 3
- [6] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 5
- [7] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 2, 3, 4, 8
- [8] Maosong Cao, Taolin Zhang, Mo Li, Chuyu Zhang, Yunxin Liu, Haodong Duan, Songyang Zhang, and Kai Chen. Condor: Enhance llm alignment with knowledge-driven data synthesis and refinement. *arXiv preprint arXiv:2501.12273*, 2025. 2, 4
- [9] CarperAI. openai summarize tldr dataset. https://huggingface.co/datasets/CarperAI/openai_summarize_tldr, 2023. 3
- [10] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024. 8
- [11] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 3
- [12] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 1, 2, 3, 4
- [13] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 8
- [14] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023. 1, 2
- [15] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv preprint arXiv:2402.03766*, 2024. 1, 3
- [16] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018. 3
- [17] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 3, 4
- [18] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 3
- [19] Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open instruction-tuned llm, 2023. 3
- [20] OpenCompass Contributors. OpenCompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>, 2023. 8
- [21] X.AI Corp. Grok-1.5 vision preview: Connecting the digital and physical worlds with our first multimodal model. <https://x.ai/blog/grok-1.5v>, 2024. 4
- [22] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun.

- Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023. 3
- [23] Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023. 3
- [24] Hongyuan Dong, Zijian Kang, Weijie Yin, Xiao Liang, Chao Feng, and Jiao Ran. Scalable vision language model training via high quality data curation. *arXiv preprint arXiv:2501.05952*, 2025. 8
- [25] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint arXiv:1903.00161*, 2019. 4
- [26] Zhangwei Gao, Zhe Chen, Erfei Cui, Yiming Ren, Weiyun Wang, Jinguo Zhu, Hao Tian, Shenglong Ye, Junjun He, Xizhou Zhu, et al. Mini-internvl: a flexible-transfer pocket multi-modal model with 5% parameters and 90% performance. *Visual Intelligence*, 2(1):1–17, 2024. 3
- [27] GlaiveAI. Glaive code assistant v3 dataset. <https://huggingface.co/datasets/glaiveai/glaive-code-assistant-v3>, 2024. 3
- [28] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*, 2023. 8
- [29] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 1, 2
- [30] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020. 3, 4
- [31] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021. 2, 3, 4
- [32] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2
- [33] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. MiniCPM: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024. 3
- [34] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024. 1
- [35] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 235–251. Springer, 2016. 4, 8
- [36] knowrohit07. know saraswati cot dataset. <https://huggingface.co/datasets/knowrohit07/know-saraswati-cot>, 2023. 3
- [37] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 3
- [38] Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. NuminaMath tir. [<https://huggingface.co/AI-MO/NuminaMath-TIR>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024. 3
- [39] Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. NuminaMath. [<https://huggingface.co/AI-MO/NuminaMath-CoT>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024. 3
- [40] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024. 4
- [41] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From live data to high-quality benchmarks: The arena-hard pipeline, 2024.
- [42] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023. 4
- [43] Wing Lian, Guan Wang, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknum". Slimorca: An open dataset of gpt-4 augmented flan reasoning traces, with verification, 2023. 3
- [44] Dongyang Liu, Renrui Zhang, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie Geng, Ziyi Lin, Peng Jin, Kaipeng Zhang, et al. Sphinx-x: Scaling data and parameters for a family of multi-modal large language models. *arXiv preprint arXiv:2402.05935*, 2024. 2
- [45] Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Zhuoer Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, et al. AlignBench: Benchmarking chi-

- nese alignment of large language models. *arXiv preprint arXiv:2311.18743*, 2023. 2, 3, 4
- [46] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of OCR in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023. 4, 8
- [47] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European Conference on Computer Vision*, pages 216–233. Springer, 2025. 8
- [48] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. DeepSeek-VL: Towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 1, 3
- [49] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022. 4
- [50] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 8
- [51] Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis: Structural embedding alignment for multimodal large language model. *arXiv:2405.20797*, 2024. 3, 8
- [52] Xudong Lu, Yinghao Chen, Cheng Chen, Hui Tan, Boheng Chen, Yina Xie, Rui Hu, Guanxin Tan, Renshou Wu, Yan Hu, et al. BlueLM-v-3b: Algorithm and system co-design for multimodal large language models on mobile devices. *arXiv preprint arXiv:2411.10640*, 2024. 1, 3, 4, 7, 8
- [53] Zhenyan Lu, Xiang Li, Dongqi Cai, Rongjie Yi, Fangming Liu, Xiwen Zhang, Nicholas D Lane, and Mengwei Xu. Small language models: Survey, measurements, and insights. *arXiv preprint arXiv:2409.15790*, 2024. 3
- [54] Gen Luo, Xue Yang, Wenhan Dou, Zhaokai Wang, Jifeng Dai, Yu Qiao, and Xizhou Zhu. Mono-internvl: Pushing the boundaries of monolithic multimodal large language models with endogenous visual pre-training. *arXiv preprint arXiv:2410.08202*, 2024. 3
- [55] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023. 3
- [56] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022. 4
- [57] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. DocVQA: A dataset for VQA on document images. In *WACV*, pages 2200–2209, 2021. 4
- [58] OpenAI. Hello GPT-4o, 2024. 1
- [59] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *NeurIPS*, 36, 2024. 3
- [60] Nazneen Rajani, Lewis Tunstall, Edward Beeching, Nathan Lambert, Alexander M. Rush, and Thomas Wolf. No robots. https://huggingface.co/datasets/HuggingFaceH4/no_robots, 2023. 3
- [61] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. *arXiv preprint arXiv:2311.12022*, 2023. 3
- [62] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. 4
- [63] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019. 4
- [64] Yixin Song, Zeyu Mi, Haotong Xie, and Haibo Chen. Powerinfer: Fast large language model serving with a consumer-grade gpu. *arXiv preprint arXiv:2312.12456*, 2023. 3
- [65] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1, 2
- [66] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. <https://github.com/InternLM/InternLM-techreport>, 2023. 3
- [67] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities, 2023. 1
- [68] Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chen-zhuang Du, Chu Wei, Congcong Wang, Dehao Zhang, Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Haotian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng, Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Junjie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu, Mengfan Dong, Mengnan Dong, Nuo Xu, Pengyu Cheng, Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu, Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe Wu, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y. Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen, Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen, Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan, Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin Yang, Zhiqi

- Huang, Zihao Huang, Zijia Zhao, and Ziwei Chen. Kimi-VL technical report, 2025. 3
- [69] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024. 2, 4
- [70] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 3
- [71] Weiha Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. CogVLM: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 2, 3, 4, 5
- [72] Weiyun Wang, Zhe Chen, Wenhai Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Jinguo Zhu, Xizhou Zhu, Lewei Lu, Yu Qiao, and Jifeng Dai. Enhancing the reasoning ability of multimodal large language models via mixed preference optimization. *arXiv preprint arXiv:2411.10442*, 2024. 8
- [73] Yejie Wang, Keqing He, Dayuan Fu, Zhuoma Gongque, Heyang Xu, Yanxu Chen, Zhexu Wang, Yujia Fu, Guanting Dong, Muxi Diao, et al. How do your code llms perform? empowering code instruction tuning with high-quality data. *arXiv preprint arXiv:2409.03810*, 2024. 3
- [74] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhramil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024. 3
- [75] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024. 3
- [76] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024. 3
- [77] Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *arXiv preprint arXiv:2406.08464*, 2024. 3
- [78] Zhenliang Xue, Yixin Song, Zeyu Mi, Le Chen, Yubin Xia, and Haibo Chen. Powerinfer-2: Fast large language model inference on a smartphone. *arXiv preprint arXiv:2406.06282*, 2024. 1, 3
- [79] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024. 3
- [80] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. 1, 2, 3, 4
- [81] Jianxin Yang. Firefly: A chinese conversational large language model. <https://github.com/yangjianxin1/Firefly>, 2023. 3
- [82] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024. 1, 3
- [83] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023. 3
- [84] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023. 8
- [85] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 8
- [86] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of ICCV*, pages 11975–11986, 2023. 3, 4
- [87] Bo-Wen Zhang, Yan Yan, Lin Li, and Guang Liu. Infinitymath: A scalable instruction tuning dataset in programmatic mathematical reasoning, 2024. 3
- [88] Yi-Kai Zhang, Shiyin Lu, Yang Li, Yanqing Ma, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, De-Chuan Zhan, and Han-Jia Ye. Wings: Learning multimodal llms without text-only forgetting. *arXiv preprint arXiv:2406.03496*, 2024. 2, 3, 4
- [89] Xiangyu Zhao, Shengyuan Ding, Zicheng Zhang, Haiyan Huang, Maosong Cao, Weiyun Wang, Jiaqi Wang, Xinyu Fang, Wenhai Wang, Guangtao Zhai, et al. Omnialign-v: Towards enhanced alignment of mllms with human preference. *arXiv preprint arXiv:2502.18411*, 2025. 2
- [90] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with

- mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023. [3](#)
- [91] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024. [2](#), [4](#)
- [92] Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang Yue. Open-codeinterpreter: Integrating code generation with execution and refinement. *arXiv preprint arXiv:2402.14658*, 2024. [3](#)
- [93] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023. [3](#)
- [94] Yichen Zhu, Minjie Zhu, Ning Liu, Zhiyuan Xu, and Yaxin Peng. Llava-phi: Efficient multi-modal assistant with small language model. In *Proceedings of the 1st International Workshop on Efficient Multimedia Computing under Limited*, pages 18–22, 2024. [3](#)