



automatic data collection and cleaning more prone to failure (Li et al., 2024).

In the multimodal era, this challenge is further amplified. Building and evaluating VLMs requires not only text corpora, but also large-scale, high-quality image–text alignment data, multimodal instruction datasets, and standardized benchmarks with reproducible experimental setups. However, in the Tibetan setting, publicly available resources for image–text alignment, multimodal instruction data, and systematic evaluation protocols remain scarce (Alam et al., 2025). As a result, progress in Tibetan VLM research and reproducibility has been slow, and it is difficult to conduct reliable and fair capability assessments under unified conditions. Figure 1 provides a motivating example of Tibetan multi-turn vision-language interaction, where the model is required to follow Tibetan instructions and perform fine-grained visual grounding. This also highlights the need for reliable training signals and standardized evaluation infrastructure tailored to Tibetan.

To fill this longstanding infrastructure gap, we propose a resource suite for Tibetan multimodal research, collectively named **FTibSuite**. FTibSuite consists of three components: (i) **FTibVLM**, a reproducible Tibetan vision–language model baseline built upon the strong open-source backbone **Qwen3-VL-8B-Instruct**, trained via a general staged adaptation pipeline consisting of continued pretraining, multimodal alignment, and multimodal instruction fine-tuning; (ii) **FTibData**, a Tibetan data collection that supports both training and instruction tuning; and (iii) **FTibBench**, a high-quality evaluation suite constructed by translating and adapting multiple mainstream multimodal benchmarks to the Tibetan setting, enabling systematic evaluation of Tibetan VLMs.

Because translating and adapting benchmarks can easily introduce noise, the credibility of the evaluation suite is particularly critical. To improve the quality and reliability of **FTibBench**, we adopt a hierarchical quality-control pipeline. We first use a large model for automatic verification and scoring to identify inconsistencies in translations, incorrect mappings between answers and options, and other high-risk issues such as numerical errors and negation mismatches. We then submit high-risk samples to manual review by Tibetan experts as the final safeguard, reducing potential systematic biases. To improve the quality and reliability of **FTibBench**, we adopt a hierarchical quality-control pipeline.

We first utilize **DeepSeek-V3** for automatic verification and scoring to identify inconsistencies and high-risk issues. We then submit these samples to manual review by Tibetan experts as the final safeguard, reducing potential systematic biases.

Experimental results show that this reproducible, data-and-benchmark-centric pipeline substantially improves Tibetan multimodal capabilities on top of the backbone baseline, and provides the first comprehensive and reproducible experimental evidence for systematic evaluation of Tibetan VLMs.

Our contributions are summarized as follows:

- We release **FTibVLM**, the first reproducible Tibetan VLM baseline built upon a strong open-source backbone model.
- We construct and open-source **FTibData**, a training data collection covering the key data types required throughout the full adaptation pipeline, including Tibetan text corpora for continual pretraining, Tibetan image–text data for multimodal alignment, and Tibetan instruction data for multimodal instruction fine-tuning.
- We build **FTibBench**, a systematic benchmark suite for Tibetan VLMs, by translating and adapting five widely used multimodal benchmarks, including BinaryVQA and MMBench, to the Tibetan setting, enabling comprehensive evaluation of Tibetan VLMs across diverse capability dimensions.

## 2 Related Work

### 2.1 Vision-Language Models

The development of vision–language models has been largely driven by two complementary lines of research: large-scale image–text pretraining and unified generative modeling. Early work typically learns cross-modal representations from web-scale image–text pairs via contrastive objectives, exemplified by CLIP (Radford et al., 2021), ALIGN (Jia et al., 2021), and BLIP (Li et al., 2022, 2023a).

As instruction following has emerged as a de facto interface for LLMs, multimodal research has increasingly shifted toward LLM-centric, generative VLMs. For example, Flamingo (Alayrac et al., 2022) introduces cross-modal connector modules between vision and language backbones, while PaLI (Chen et al., 2022) emphasizes joint scaling of vision and language. In the open-source ecosystem, LLaVA (Liu et al., 2023) and Instruct-BLIP (Dai et al., 2023) demonstrate that converting heterogeneous multimodal tasks into an “instruc-

172	tion–response” format is a key step toward building	
173	general-purpose visual assistants, while Qwen-VL	
174	(Bai et al., 2023) systematically highlights the im-	
175	portance of a strong backbone, multi-stage training,	
176	and curated multilingual multimodal corpora for	
177	general capabilities.	
178		
	<b>2.2 Data and Evaluation for VLMs</b>	
179	The advancement of VLMs has been largely en-	
180	abled by the joint maturation of high-quality in-	
181	struction data and diagnostic evaluation suites. For	
182	extremely low-resource languages such as Tibetan,	
183	the availability of a reusable data pipeline span-	
184	ning continual pretraining to instruction tuning is	
185	often a key determinant of whether an open and	
186	sustainable research ecosystem can be established.	
187	On the data side, instruction construction is in-	
188	creasingly moving beyond purely synthetic QA to-	
189	ward broader task coverage and higher annotation	
190	quality. Vision-Flan (Xu et al., 2024), for example,	
191	reformulates a wide range of academic datasets	
192	into a unified visual instruction format and demon-	
193	strates the effectiveness of a two-stage instruction-	
194	tuning recipe, first leveraging high-quality human-	
195	labeled tasks and then scaling with synthetic align-	
196	ment data. LoResMT (Xiao et al., 2025) further	
197	explores systematic pipelines that transform par-	
198	allel text corpora into multimodal training data in	
199	low-resource settings.	
200	On the evaluation side, benchmarks are shift-	
201	ing from coarse-grained leaderboards toward diag-	
202	nostic and reliability-oriented assessments. POPE	
203	(Li et al., 2023b) evaluates object hallucination in	
204	VLMs. MME (Fu et al., 2025), in contrast, of-	
205	fers a more comprehensive capability profile by	
206	covering both perception- and cognition-level sub-	
207	tasks. Beyond these, BinaryVQA (Borji, 2023)	
208	probes out-of-distribution generalization and bias,	
209	while COREVQA (Chintapatla et al., 2025) targets	
210	fine-grained observation and reasoning in crowded	
211	scenes, further revealing the brittleness of current	
212	VLMs under challenging visual conditions. More	
213	recently, UPD (Miyai et al., 2025) highlights that	
214	high multiple-choice VQA scores alone do not nec-	
215	essarily imply genuine understanding. However,	
216	despite the abundance of existing evaluations, they	
217	are predominantly English-centric, and there is no	
218	widely adopted, publicly released Tibetan counter-	
219	part of mainstream multimodal benchmarks.	
	<b>2.3 Low-Resource Language Adaptation and Resources</b>	
	Low-resource language capability is typically	
	achieved through strong backbone plus target-	
	distribution adaptation strategy, such as continual	
	pretraining (Gururangan et al., 2020). The same	
	principle applies in the multimodal setting: rather	
	than training from scratch, it is often more effective	
	and cost-efficient to start from a strong multimodal	
	backbone and perform distribution alignment and	
	stage-wise adaptation.	
	For adaptation efficiency and stability,	
	parameter-efficient fine-tuning provides a solu-	
	tion for low-resource and multi-stage training.	
	Adapters enable multi-task expansion by inserting	
	lightweight modules while keeping the backbone	
	frozen (Houlsby et al., 2019), and LoRA (Hu et al.,	
	2022) substantially reduce memory and parameter	
	overhead through low-rank updates and quantized	
	training, making iterative development feasible	
	under limited compute budgets. BranchLoRA	
	(Zhang et al., 2025) further mitigates catastrophic	
	forgetting in continual learning via structured	
	routing and freezing mechanisms. Taken together,	
	these studies suggest that effective low-resource	
	adaptation should not only acquire new capabili-	
	ties, but also preserve existing ones and support	
	controllable transfer.	
	From the perspective of resource development,	
	there has been notable progress on Chinese mi-	
	nority languages in terms of textual corpora and	
	pretrained models. MC <sup>2</sup> (Zhang et al., 2023) sys-	
	tematically constructs multilingual corpora for mi-	
	nority languages in China, while CINO (Yang et al.,	
	2022) and XLM-SWCM (Su et al., 2025) train ded-	
	icated multilingual language models for Chinese	
	minority languages. However, these efforts primar-	
	ily focus on text-only resources. Publicly available	
	multimodal alignment data, instruction-tuning data,	
	and standardized evaluation pipelines for Tibetan	
	remain absent.	
	<b>3 Constructing a Comprehensive VLM Suite for Tibetan</b>	
	This section builds and releases the first compre-	
	hensive Tibetan research resource and evaluation	
	infrastructure suite for vision–language models	
	(VLMs), <b>FTibSuite</b> , for the Tibetan community.	
	It aims to address three long-standing foundational	
	gaps in Tibetan multimodal research: (i) the lack	
	of reusable training corpora, (ii) the lack of Tibetan	

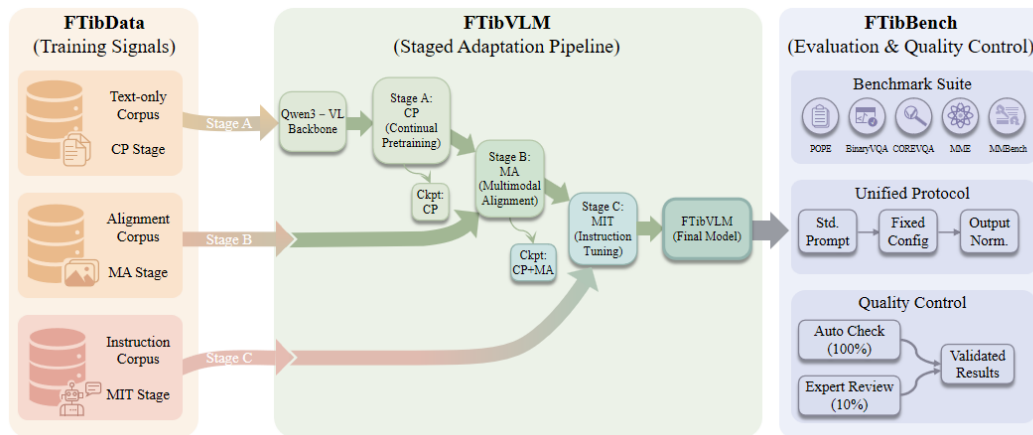


Figure 2: **FTibSuite overview**, which consists of three coupled components: **FTibData**, providing reusable multilingual and multimodal training signals; **FTibVLM**, a staged adaptation pipeline that incrementally adapts a vision-language backbone to Tibetan via continual pretraining (CP), multimodal alignment (MA), and instruction tuning (MIT) and **FTibBench**, unified evaluation framework with standardized protocols and hierarchical .

270 evaluation benchmarks that are aligned with main-  
 271 stream English benchmarks and whose quality can  
 272 be verified, and (iii) the lack of baseline models that  
 273 are reproducible and comparable under a unified  
 274 evaluation protocol.

275 To this end, we organize the resources produced  
 276 in this work following a “data–evaluation–baseline”  
 277 structure, as summarized in Figure 2. **FTibData**  
 278 provides reusable training signals for staged adap-  
 279 tation; **FTibVLM** instantiates these signals into re-  
 280 producible reference checkpoints; and **FTibBench**  
 281 standardizes evaluation with a unified protocol ac-  
 282 companied by a hierarchical quality-control work-  
 283 flow, whose feedback is used to refine translation  
 284 and parsing rules for subsequent iterations.

### 285 3.1 Training corpus

286 To introduce stable Tibetan generation capabili-  
 287 ties without modifying the model architecture, we  
 288 first conduct Tibetan-oriented continual pretrain-  
 289 ing. The goal of this stage is to explicitly shift  
 290 the backbone’s language distribution toward the  
 291 Tibetan text space, enabling more reliable Tibetan  
 292 language modeling and providing a solid linguistic  
 293 foundation for subsequent multimodal alignment  
 294 and instruction fine-tuning.

295 We use three categories of text data: a Tibetan  
 296 subset from MC<sup>2</sup> (Zhang et al., 2024), publicly  
 297 available Tibetan instruction data (e.g., tibetan-  
 298 mix-instruction-tuning-60K), and the Chinese LC-  
 299 STS corpus. After unified cleaning, the combined  
 300 dataset contains approximately 2.2 million sam-  
 301 ples, with about 70% Tibetan and 30% Chinese.  
 302 Under the low-resource setting, we retain a certain

303 proportion of Chinese data for two main reasons:  
 304 first, preserving the backbone’s original Chinese  
 305 capability is practically valuable; second, we in-  
 306 terleave source-language data during cross-lingual  
 307 continual pretraining as data replay, motivated by  
 308 discussions on mitigating catastrophic forgetting in  
 309 continual pretraining. (Zheng et al., 2024)

310 We construct the cross-modal image–text align-  
 311 ment corpus based on the Chinese captioning data  
 312 of AI Challenger (Wu et al., 2017), using a fixed  
 313 pool of 100k images. In total, we build 150k one-  
 314 image–one-caption pairs: 100k Tibetan pairs trans-  
 315 lated from primary Chinese captions as the main  
 316 grounding signal, 30k original Chinese pairs re-  
 317 tained to stabilize training, and 20k additional Chi-  
 318 nese pairs with alternative captions to enhance ex-  
 319 pression diversity.

320 We build the multimodal instruction fine-tuning  
 321 corpus based on **Vision-Flan** (Xu et al., 2024), with  
 322 fixed task-type ratios (caption 25%, VQA 40%,  
 323 classification 20%, counting 5%, others 10%). We  
 324 translate 30k sampled instances into Tibetan as the  
 325 primary instruction set, and translate another 10k  
 326 instances (with the same ratios) into Chinese to  
 327 maintain Chinese capability. We further create a  
 328 10k Tibetan–Chinese parallel subset by translating  
 329 the same image-conditioned instances into both  
 330 languages, and normalize all data into a unified  
 331 multimodal instruction–response format.

### 332 3.2 Tibetan visual-linguistic baseline

333 We build **FTibVLM** on top of the multimodal back-  
 334 bone **Qwen3-VL-8B-Instruct** (Yang et al., 2025)  
 335 and adopt a unified conversational interface of “im-

age + textual instruction” for both training and inference. The adaptation process follows a three-stage pipeline driven by the three data modules in **FTibData**.

**Stage A — Continual Pretraining (CP).** This stage adapts the backbone’s language distribution toward Tibetan through text-only continual pretraining. The goal is to establish a stable linguistic foundation for Tibetan generation and understanding, while retaining the backbone’s original Chinese capability boundary.

**Stage B — Multimodal Alignment (MA).** Given the linguistic prior obtained in CP, this stage performs caption-based image–text alignment to strengthen the correspondence between visual semantics and Tibetan expressions. This improves cross-modal grounding stability when the model is prompted in Tibetan.

**Stage C — Multimodal Instruction Tuning (MIT).** The final stage fine-tunes the model on multimodal instruction data to enhance instruction following, multi-task execution, and interactive usability. Beyond task accuracy, MIT stabilizes response format and decision behaviors under Tibetan prompts.

To support controlled analysis, we save checkpoints after each stage and evaluate three variants under the same backbone starting point: **Base** (no Tibetan adaptation), **CP+MA**, and **CP+MA+MIT** (the final FTibVLM). These variants correspond one-to-one to the three stages above and together constitute the full staged adaptation pipeline.

### 3.3 Benchmarks and Metrics

FTibBench is designed to address the lack of widely adopted and publicly released multimodal benchmarks for Tibetan. Direct translation from English benchmarks often introduces systematic noise (e.g., negation mismatch, numerical drift, entity misalignment, and option–answer mapping errors), which undermines the reliability of evaluation results. Rather than merely localizing existing datasets, FTibBench aims to provide a reproducible and auditable evaluation protocol with controlled differences across models. The full judging prompt and evaluation policy are provided in Appendix F.

**Benchmark Suite.** FTibBench covers five major multimodal benchmarks in Tibetan: **POPE** (random / popular / adversarial subsets), **BinaryVQA**, **COREVQA**, **MME**, and **MMBench-dev**. During

translation and adaptation, we preserve the original task definitions and answer spaces as much as possible, and expose all benchmarks through a unified execution interface to facilitate reuse, extension, and controlled comparison. **Details about FTibBench can be found in Appendix A.**

**Evaluation Protocol.** To ensure comparability across models, FTibBench standardizes evaluation along three components: prompt formatting, decoding configuration, and output normalization. All models are evaluated under an identical Tibetan prompt template and a fixed inference configuration. For classification-style tasks, we restrict the answer space to reduce ambiguity: multiple-choice benchmarks require outputting only the option symbol, while binary tasks are normalized to a 0/1 decision space. Outputs are subsequently processed through a unified normalization and parsing procedure, and we additionally report the proportion of unmappable outputs as a stability indicator; in our experiments, this invalid rate is 0.

**Annotation and Quality Control.** To improve the credibility of translated benchmarks, we adopt a hierarchical quality-control workflow. For each benchmark, all of instances are first screened via automated consistency checking, and an additional 10% are manually reviewed by Tibetan-language experts as a final gate. Automated verification follows a unified rubric scoring accuracy (0–2), completeness (0–2), and Tibetan linguistic naturalness (0–1), yielding a traceable quality score in 0–5. We conducted small-scale comparative tests across multiple large models together with Tibetan experts, and selected **DeepSeek-V3** (DeepSeek-AI et al., 2025) as the primary automatic verifier due to its stability in Tibetan semantic judgment. Manual inspection focuses on high-risk error categories (entity alignment, negation, numerics, option–answer mapping), and all confirmed issues are fed back to refine translations and parsing rules. The LLM-judge prompt can be found in Appendix F.

Taken together, FTibBench provides not only Tibetan counterparts of mainstream multimodal benchmarks, but also a controlled and auditable evaluation protocol, enabling fair, reproducible, and stability-aware comparison of Tibetan VLMs.

## 4 Experiments

### 4.1 Experimental Setup

**Model Setup.** We use **Qwen3-VL-8B-Instruct** (Yang et al., 2025) as the backbone model, adopt the same three-stage adaptation pipeline as mainstream open-source VLMs, and keep the model architecture unchanged. The three training stages include continual pretraining, cross-modal alignment training, and multimodal instruction fine-tuning. Our key motivation for choosing a strong backbone is that its general visual understanding and instruction-following capabilities provide a higher starting point for transferring to low-resource languages, allowing us to focus our primary efforts on completing the Tibetan data and evaluation pipeline rather than training an entire multimodal system from scratch.

**Implementation details.** All three training stages use parameter-efficient fine-tuning with LoRA, are run in bf16 precision, and are trained with DDP on  $8 \times$  RTX 4090 GPUs. We use the AdamW optimizer, adopt a cosine learning-rate schedule, and set the gradient clipping threshold to 1.0. To match the training budget with the objectives of each stage, we apply stage-specific configurations of frozen and trainable modules. In the continual pretraining stage, which focuses on language-distribution adaptation, we freeze the visual encoder and the multimodal projection layer, and inject LoRA only into the language components to enable low-cost transfer. In the cross-modal alignment and instruction fine-tuning stages, which target visual alignment and improved instruction-following ability, we freeze the visual encoder while keeping the projection layer trainable, so that we can stably preserve the backbone’s visual representations while more effectively adapting the cross-modal mapping and instruction behaviors. **Training and hyperparameter details are provided in Appendix B.**

**Benchmarks** This paper proposes **FTibBench** to evaluate models’ Tibetan multimodal capabilities, covering **POPE** (random, popular, adversarial), **BinaryVQA**, **COREVQA**, **MME**, and **MMBench-BO** (all in Tibetan; each benchmark follows the official default data split, or uses the official dev set). Chinese capability retention is reported on **MMBench-CN** (dev). **POPE**, **BinaryVQA**, and **COREVQA** report Accuracy and F1. **MME** strictly follows the official evaluation procedure and re-

ports Acc as well as the stricter Acc+. **MMBench** and **MMBench-CN** (dev) use standard multiple-choice evaluation and report the overall score.

**Evaluation protocols.** All experiments strictly adhered to a unified standard, including Tibetan prompt templates, deterministic decoding settings, a constrained output space, and unified parsing rules. Specifically, multiple-choice questions required the model to output only the letter of the option; the binary classification task normalized the answer space to 1 and 0. The invalid rate was 0 in the experiments, indicating that the output parsing is stable under this protocol, and the scoring is unaffected by parsing failures.

### 4.2 Experimental Results

Tables 1 to 3 summarize the systematic evaluation results on FTibBench in the Tibetan setting. Overall, FTibVLM achieves substantial improvements over Base on POPE, BinaryVQA, COREVQA, MME, and MMBench-BO.

As shown in Table 1, on POPE, which focuses on diagnosing object hallucination, FTibVLM consistently outperforms Base across all three subsets. Specifically, on POPE-random, accuracy increases from 47.53 to 80.56, and F1 increases from 43.62 to 80.51. On the more challenging POPE-popular and POPE-adversarial subsets, accuracy reaches 81.70 and 78.63, and F1 reaches 73.40 and 78.49. These results indicate that after Tibetan-side adaptation, the model is more robust under the image-consistency and hallucination-sensitive conditions captured by POPE, and the gains remain consistent across subsets of different difficulty levels.

On BinaryVQA, FTibVLM also delivers consistent gains. Accuracy increases from 54.46 to 76.01, and F1 increases from 53.08 to 73.25, indicating that the model’s discriminative ability is substantially strengthened in the binary VQA setting with a constrained answer space. For COREVQA, which emphasizes fine-grained observation and reasoning in crowded scenes, accuracy improves from 31.49 to 50.85.

Table 2 reports the comparison results on MME, a multi-task evaluation organized by capability dimensions. Overall, FTibVLM improves both accuracy and the stricter accuracy plus metric across multiple subtasks. For example, on basic perception tasks such as existence and color, accuracy reaches 88.33 and 78.33, while accuracy plus reaches 80.00 and 63.33. On tasks closer to scene

Model	POPE(random)		POPE(popular)		POPE(adversarial)		BinaryVQA		COREVQA	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Base	47.53	43.62	46.30	60.65	46.22	60.49	54.46	53.08	31.49	42.16
<b>FTibVLM(ours)</b>	<b>80.56</b>	<b>80.51</b>	<b>81.70</b>	<b>73.40</b>	<b>78.63</b>	<b>78.49</b>	<b>76.01</b>	<b>73.25</b>	<b>50.85</b>	35.52

Table 1: Main results on Tibetan hallucination robustness and binary VQA benchmarks.

Model	existence		color		posters		scene		count	
	Acc	Acc+	Acc	Acc+	Acc	Acc+	Acc	Acc+	Acc	Acc+
Base	50.00	33.33	55.00	10.00	63.95	34.69	60.25	31.50	50.00	0.00
<b>FTibVLM(ours)</b>	<b>88.33</b>	<b>80.00</b>	<b>78.33</b>	<b>63.33</b>	<b>77.55</b>	<b>59.18</b>	<b>75.75</b>	<b>53.00</b>	<b>66.70</b>	<b>33.33</b>

Table 2: Main results on MME for Tibetan multimodal capability profiling.

Model	Overall	LR	AR	RR	FP-S	FP-C	CP
Base	42.97	52.57	40.14	38.01	43.54	42.00	43.41
<b>FTibVLM(ours)</b>	<b>67.78</b>	<b>61.65</b>	<b>63.23</b>	<b>66.07</b>	<b>65.16</b>	<b>68.29</b>	<b>76.01</b>

Table 3: Main results on MMBench(dev) for Tibetan multimodal understanding and reasoning.

Model	Overall	LR	AR	RR	FP-S	FP-C	CP
Base	88.50	84.47	88.41	85.17	91.87	84.80	89.72
<b>FTibVLM(ours)</b>	88.15	83.50	87.12	84.04	91.42	84.05	90.80

Table 4: Chinese capability retention on MMBench-CN(dev) after Tibetan adaptation.

Model	MMBench(dev)	POPE		
		random	popular	adversarial
Base	42.97	47.53	46.38	46.22
CP + MA	60.87	71.10	73.83	69.35
FTibVLM	67.78	80.56	81.70	78.63

Table 5: Stage-wise ablation on MMBench(dev) and POPE for Tibetan adaptation (base. vs. caption alignment. vs. + instruction SFT).

531 understanding such as posters and scene, FTibVLM  
532 also achieves steady gains, with scene improving  
533 from 60.25 to 75.75. On the more challenging  
534 count task, accuracy rises from 50.00 to 66.70 and  
535 accuracy plus rises from 0.00 to 33.33. A finer-  
536 grained MME subtask analysis is reported in Ap-  
537 pendix C (Table 8), which shows that improve-  
538 ments are most pronounced on basic perception  
539 and decision-oriented dimensions, while OCR- and  
540 text-related subtasks remain comparatively chal-  
541 lenging. Given that OCR appears to be a primary  
542 bottleneck, we further conduct a targeted Tibetan  
543 OCR adaptation study; details are provided in Ap-

pendix E.

FTibVLM increases the overall score from 42.97  
544 to 67.78, demonstrating a substantial gain on the  
545 Tibetan multiple-choice comprehensive evaluation.  
546 Breaking down by dimensions, the model achieves  
547 61.65, 63.23, and 66.07 on LR, AR, and RR, and  
548 65.16, 68.29, and 76.01 on FP-S, FP-C, and CP.  
549 These results indicate that the multi-dimensional ca-  
550 pability categories covered by MMBench reliably  
551 capture the overall improvements of the model in  
552 Tibetan multimodal understanding and reasoning.  
553 To complement the coverage of the main evaluation  
554 on cross-modal semantic consistency and visual en-  
555 tailment reasoning, we conduct an additional diag-  
556 nostic study on SNLI-VE; the experimental setup  
557 and full results are reported in Appendix D.

Overall, the experimental results show that the  
560 unified evaluation protocol and hierarchical quality-  
561 control pipeline established by FTibBench can reli-  
562 ably differentiate Tibetan multimodal capabilities  
563 across models under the same setting. Within  
564 this evaluation loop, FTibVLM exhibits consistent  
565

and substantial improvements over Base across the core tasks in FTibBench, demonstrating that our stage-wise adaptation driven by FTibData effectively enhances Tibetan multimodal understanding and reasoning, and provides a reproducible and diagnosable strong baseline for future work.

### 4.3 Ablation Studies

#### 4.3.1 Stage-Wise Ablation

To quantify the marginal contributions of different data modules in FTibData and conduct an interpretable comparison under the unified evaluation protocol established by FTibBench, we evaluate three stage-wise checkpoints on FTibBench: **Base**, **CP+MA** (after continual pretraining and multimodal alignment), and the final model **FTibVLM** (further incorporating multimodal instruction tuning on top of CP+MA). As shown in Table 5, overall performance exhibits a consistent upward trend as modules are introduced, with larger gains from Base to CP+MA and additional robust improvements from CP+MA to FTibVLM.

This stage-wise improvement aligns with the functional division of the three supervision signals. CP establishes a Tibetan language-distribution foundation while reducing cross-lingual forgetting, MA strengthens the alignment between visual semantics and Tibetan expressions to improve cross-modal consistency and discriminative stability, and MIT further boosts performance in interactive and multi-task settings by shaping instruction following and overall capability. Overall, these controlled results support a resource-module-driven and interpretable improvement conclusion: under a unified evaluation protocol and a fixed implementation setup, performance gains emerge steadily as modules are added, and can be further attributed to the incremental contributions of different modules across diagnostic dimensions of the benchmarks.

#### 4.3.2 Chinese Capability Retention

Stage-wise adaptation can inject target-language capability, but it may also introduce cross-lingual degradation. To examine whether our training pipeline affects Chinese multimodal performance, we compare our Tibetan multimodal model FTibVLM with the baseline Base on MMBench-CN (dev). As shown in Table 4, FTibVLM achieves an overall score of 88.15, which is essentially on par with Base at 88.50, exhibiting only minor fluctuations and no consistent downward trend. Notably,

FTibVLM even improves on the CP dimension, increasing from 89.72 to 90.80.

These results validate that the mixed-corpus design of FTibData is both effective and necessary. We retain a certain proportion of Chinese in the text corpus and introduce a Chinese anchor subset in the image-text alignment and instruction fine-tuning stages, with the goal of providing cross-lingual stability constraints during training and mitigating cross-lingual forgetting and output degradation caused by continual pretraining and subsequent multi-stage adaptation. Overall, while substantially strengthening Tibetan multimodal capability, FTibVLM does not exhibit a noticeable loss in overall Chinese multimodal competence, providing a more stable capability boundary for cross-lingual reuse and real-world deployment in Tibetan scenarios.

## 5 Conclusion

In this paper, we introduce FTibSuite, a resource suite for Tibetan vision-language modeling that integrates FTibData, FTibBench, and the first Tibetan VLM baseline FTibVLM, together with a reproducible training and evaluation pipeline built upon Qwen3-VL-8B-Instruct. We construct FTibData and adopt a three-stage adaptation pipeline—Tibetan continual pretraining, image-text alignment, and multimodal instruction tuning—to equip FTibVLM with Tibetan generation, grounding, and instruction-following abilities in a reproducible manner. We build FTibBench by migrating five established multimodal benchmarks into the Tibetan setting, covering hallucination robustness, binary decision stability, dense-scene understanding, capability profiling, and multiple-choice reasoning. To improve benchmark reliability, we use DeepSeek-V3 for automatic verification and rubric-based scoring, and conduct Tibetan-expert review and annotation for high-risk cases, helping partially fill the infrastructure gap for Tibetan multimodal research. Experiments show consistent Tibetan gains with minimal degradation of the backbone’s Chinese capability, and staged checkpoints support controlled analysis of how different adaptation signals contribute to improvements. Future work will focus on improving Tibetan multimodal supervision quality, expanding benchmark coverage, and developing more robust multilingual adaptation strategies, especially for OCR and in-image text understanding.

663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
  
676  
  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
  
699  
  
700  
701  
702  
703  
704  
  
705  
706  
707  
708  
709  
710

## Limitations

This work focuses on data- and training-driven adaptation on top of a strong open-source backbone, rather than proposing new model architectures, so the improvements are bounded by the capabilities of the underlying design. In addition, the current training pipeline still has room to improve: some Tibetan multimodal supervision is obtained through translation and dataset repurposing, which may introduce noise and limit robustness. Future work could strengthen these aspects with higher-quality Tibetan-native multimodal data and more principled multilingual adaptation strategies.

## Ethical Considerations

This work aims to promote inclusive vision-language modeling by extending Tibetan multimodal research through a resource suite that supports more reproducible training and evaluation. The resources in FTibSuite are constructed by adapting existing datasets and benchmark designs; we make efforts to respect original licenses and document provenance and usage constraints for each component. To improve benchmark reliability, we employ a tiered quality-control workflow with large-model automatic verification and scoring followed by Tibetan-expert review and annotation for high-risk cases. While these procedures reduce common adaptation errors and evaluation noise, residual artifacts and pretrained biases may persist, and benchmark scores should not be taken as complete evidence of real-world Tibetan multimodal competence. Finally, stronger Tibetan VLM capability may be misused, such as for generating misleading content, and we therefore encourage transparent reporting of limitations and careful deployment in high-stakes settings.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Nahid Alam, Karthik Reddy Kanjula, Surya Guthikonda, Timothy Chung, Bala Krishna S Vegesna, Abhipsha Das, Anthony Susevski, Ryan Sze-Yin Chan, SM Udin, Shayekh Bin Islam, and 1 others. 2025. Behind maya: Building a multilingual vision language model. *arXiv preprint arXiv:2505.08910*.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736. 711–717

Bo An. 2023. Prompt-based for low-resource tibetan text classification. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(8):1–13. 718–721

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*. 722–725

Ali Borji. 2023. Binaryvqa: A versatile test set to evaluate the out-of-distribution generalization of vqa models. *arXiv preprint arXiv:2301.12032*. 726–728

Lifeng Chen, Ryan Lai, and Tianming Liu. 2025. Adapting large language models to low-resource tibetan: A two-stage continual and supervised fine-tuning study. *arXiv preprint arXiv:2512.03976*. 729–732

Xi Chen, Xiao Wang, Soravit Changpinyo, Anthony J Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, and 1 others. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*. 733–738

Ishant Chintapatla, Kazuma Choji, Naaisha Agarwal, Andrew Lin, Hannah You, Charles Duong, Kevin Zhu, Sean O’Brien, and Vasu Sharma. 2025. Corevqa: A crowd observation and reasoning entailment visual question answering benchmark. *arXiv preprint arXiv:2507.13405*. 739–744

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 8440–8451. 745–751

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, and 1 others. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*. 752–757

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267. 758–763

DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, 764–766



879 tasks in visual instruction tuning. *arXiv preprint*  
880 *arXiv:2402.11690*.

881 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
882 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,  
883 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-  
884 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao  
885 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41  
886 others. 2025. *Qwen3 technical report*. *Preprint*,  
887 *arXiv:2505.09388*.

888 Ziqing Yang, Zihang Xu, Yiming Cui, Baoxin Wang,  
889 Min Lin, Dayong Wu, and Zhigang Chen. 2022.  
890 Cino: A chinese minority pre-trained language model.  
891 *arXiv preprint arXiv:2202.13558*.

892 Chen Zhang, Mingxu Tao, Quzhe Huang, Jiuheg Lin,  
893 Zhibin Chen, and Yansong Feng. 2023.  $Mc^2$ : A  
894 multilingual corpus of minority languages in china.  
895 *CoRR*.

896 Chen Zhang, Mingxu Tao, Quzhe Huang, Jiuheg Lin,  
897 Zhibin Chen, and Yansong Feng. 2024.  $Mc2$ : To-  
898 wards transparent and culturally-aware nlp for mi-  
899 nority languages in china. In *Proceedings of the*  
900 *62nd Annual Meeting of the Association for Compu-*  
901 *tational Linguistics (Volume 1: Long Papers)*, pages  
902 8832–8850.

903 Duzhen Zhang, Yong Ren, Zhong-Zhi Li, Yahan Yu,  
904 Jiahua Dong, Chenxing Li, Zhilong Ji, and Jin-  
905 feng Bai. 2025. Enhancing multimodal continual  
906 instruction tuning with branchlora. *arXiv preprint*  
907 *arXiv:2506.02041*.

908 Wenzhen Zheng, Wenbo Pan, Xu Xu, Libo Qin, Li Yue,  
909 and Ming Zhou. 2024. Breaking language barriers:  
910 Cross-lingual continual pre-training at scale. *arXiv*  
911 *preprint arXiv:2407.02118*.

912	<b>A FTibBench benchmark adaptation and</b>	<b>A.3 High-Risk Error Types and Checklist</b>	957
913	<b>quality control details</b>	As is shown in Table 6, we treat the following error	958
914	<b>A.1 Benchmark Composition, Splits, and Scale</b>	types as high risk and prioritize them for automated	959
915	FTibBench currently includes Tibetan versions of	screening and manual review.	960
916	five representative multimodal evaluation bench-	<b>A.4 Automated Quality Control: Evaluation</b>	961
917	marks, covering complementary dimensions such	<b>Rubric and Field-Level</b>	962
918	as hallucination robustness, binary decision mak-	We rate each sample along three dimensions includ-	963
919	ing, dense-scene understanding, multi-dimensional	ing accuracy, completeness, and expression fluency,	964
920	capability profiling, and comprehensive multiple-	with a total score ranging from 0 to 5, and record	965
921	choice understanding.	brief diagnostic comments to support revision and	966
922	<b>POPE.</b> POPE contains three splits (adversarial,	regression analysis. The scoring dimensions are	967
923	popular, and random) and is used to evaluate hallu-	listed in the Table 7.	968
924	cination tendencies and robustness under question	<b>A.5 Score-Triggered Revision and Manual</b>	969
925	answering about target existence.	<b>Review Strategy</b>	970
926	<b>BinaryVQA.</b> BinaryVQA is a binary-	We adopt a tiered quality-control strategy of “ <i>au-</i>	971
927	classification VQA benchmark whose answer	<i>tomatic scoring + human fallback</i> ” to balance	972
928	space is strictly 0/1, and is used to evaluate	quality and cost:	973
929	decision stability and output controllability.	<ul style="list-style-type: none"> <li>• <b>Total <math>\leq 2</math>: mandatory revision and manda-</b></li> </ul>	974
930	<b>COREVQA.</b> COREVQA targets fine-grained	<b>tory human review.</b> Such samples typically	975
931	observation and reasoning in dense/complex	exhibit missing key terms, semantic drift, or	976
932	scenes, emphasizing counting, relations, and local	clearly unnatural phrasing, which may com-	977
933	entity understanding.	promise evaluation consistency and fairness.	978
934	<b>MME.</b> MME is a multi-dimensional capability	They are therefore prioritized for correction	979
935	profiling benchmark, used to diagnose a model’s ca-	and verified by human reviewers.	980
936	pability structure across multiple task dimensions.	<ul style="list-style-type: none"> <li>• <b>Total <math>\geq 3</math>: entered into a spot-check review</b></li> </ul>	981
937	It includes 14 category subsets: artwork, celebrity,	<b>pool.</b> Issues are usually minor, such as slight	982
938	code_reasoning, color, commonsense_reasoning,	verbosity, minor over-translation, or less-than-	983
939	count, existence, landmark, numerical_calculation,	natural wording. We randomly sample <b>10%</b>	984
940	OCR, position, posters, scene, and text_translation.	from this pool for human review: annotators	985
941	<b>MMBench-BO.</b> MMBench-BO is a multiple-	re-score and label the samples, and we check	986
942	choice benchmark for overall comparison of multi-	whether the human judgments are consistent	987
943	modal understanding and reasoning abilities.	with the model-generated scores.	988
944	<b>A.2 Translation and Adaptation Principles</b>	To improve the verifiability of Tibetan	989
945	To maximize fairness in cross-model comparisons,	translation/adaptation quality, we log, for	990
946	we follow the principle of “ <i>unchanged task def-</i>	sampling BinaryVQA instances, the English	991
947	<i>inition, unchanged answer space, and structure</i>	question (question_en), the Tibetan ques-	992
948	<i>aligned as much as possible</i> ” during translation	tion (question), the three dimension scores	993
949	and adaptation. We strictly maintain answer-space	(accuracy/completeness/tibetan_expression),	994
950	consistency: for binary tasks, we uniformly use	the total score (total, 0-5), and a brief diagnostic	995
951	0/1 (1 denotes “yes/present/true,” and 0 denotes	comment (comment). During human review,	996
952	“no/absent/false”); for multiple-choice tasks, we	we cross-check the automatic scoring results.	997
953	keep the original option set unchanged and restrict	Except for a small number of cases that require	998
954	outputs to A/B/C/D. Meanwhile, we ensure that	additional explanation, human judgments are	999
955	structural fields are traceably mappable, facilitat-	largely consistent with the automatic scores.	1000
956	ing subsequent alignment analyses and audits.	As is shown in Figure 3, we further present three	1001
		representative low-scoring examples to illustrate	1002
		common translation error types and the correspond-	1003
		ing reasons for score deductions.	1004



Dimension	Score	Criteria
Accuracy	2	<b>Semantically accurate:</b> Faithfully conveys the core meaning; key information remains consistent.
	1	<b>Deviations present:</b> Largely correct, but with issues such as missing/incorrect terminology, unclear correspondences, or semantic drift.
	0	<b>Incorrect:</b> Meaning is distorted; key content is wrong or there are severe mistranslation errors.
Completeness	2	<b>Complete information:</b> No obvious omissions.
	1	<b>Minor missing:</b> The main message is present, but some details are missing, which may introduce slight ambiguity.
	0	<b>Severe missing:</b> Key information is missing, changing the question intent or the answer space.
Expression Fluency	1	<b>Natural:</b> Fluent and natural expression, with no obvious traces of literal translation.
	0	<b>Awkward:</b> Stilted or unnatural phrasing; hard to read or with clear literal-translation artifacts.

Table 7: Scoring rubric for translation quality across accuracy, completeness, and expression fluency.

## C MME Subtask Evaluation Results

As shown in Table 8, compared with the base model, **FTibVLM** improves **MME Overall Acc** from **63.98%** to **75.02%** (+11.04 percentage points), and raises the stricter **Acc+** from **33.28%** to **54.51%** (+21.23 percentage points). From a task-level breakdown, the model shows particularly pronounced gains in existence, color, and count, basic perception and decision-oriented tasks with clear improvements in **Acc+** and output stability. In contrast, **OCR** and **text\_translation**, which rely more heavily on the text-recognition and cross-lingual translation pipeline, remain the main bottlenecks. Future work could further strengthen these capabilities by incorporating higher-quality Tibetan **OCR** and in-image text alignment data, as well as enforcing translation-consistency constraints.

## D Additional Diagnostic Benchmark: SNLI-VE

To complement the evaluation coverage of **FTibBench**, we additionally evaluate the model on the Tibetan three-way classification set of **SNLI-VE** to assess cross-modal logical consistency and visual entailment reasoning ability. This task requires the model to determine the relationship between an image and a textual hypothesis and output one of three labels: contradiction/neutral/entailment (encoded as 0/1/2). We adopt a unified scoring setup (*robust candidates + better gate*) and filter and aggregate predictions under a strict gating strategy (Gate: mode=strict\_entailment, min\_conf\_2=0.62, min\_margin\_2=0.1). The evaluation contains

17,901 samples, with no invalid outputs, no missing images, and no skipped samples.

As shown in Table 9, compared with Base, **FTibVLM** achieves a substantial improvement on this additional diagnostic task: overall **Accuracy** increases from **0.3715** to **0.5432**, and **Macro-F1** also rises from **0.3072** to **0.5400**. At the class level, the base model tends to over-predict contradiction (class 0), whereas **FTibVLM** produces a more balanced prediction distribution and attains higher overall F1 on the neutral and entailment classes. These results indicate that the three-stage adaptation yields clear gains in cross-modal semantic consistency and reasoning stability.

## E Tibetan OCR Adaptation

### E.1 Motivation

In our main experiments, we observe that although the three-stage training substantially improves Tibetan multimodal understanding and reasoning (e.g., on MME and VQA), the gains on tasks involving in-image text recognition (OCR) are not apparent. To examine whether *language adaptation* can directly improve Tibetan OCR recognition ability, and whether a *small amount of Tibetan OCR instruction data* can compensate for this capability, we conduct additional targeted experiments on Tibetan OCR. The results indicate that adapting a VLM to the target language, even when mixing in a small amount of OCR data into the existing instruction fine-tuning set, which does not effectively improve the model’s OCR recognition ability for that language. OCR therefore remains one of the primary bottlenecks.

SubTask	Base(ACC)	Base(ACC+)	FTibVLM(ACC)	FTibVLM(ACC+)
Existence	50.00%	3.33%	88.33%	80.00%
Color	55.00%	10.00%	78.33%	63.33%
Code Reasoning	62.50%	30.00%	80.00%	60.00%
Count	50.00%	0.00%	66.67%	33.33%
Artwork	52.00%	8.00%	68.50%	46.00%
Scene	60.25%	31.50%	75.75%	53.00%
Posters	63.95%	34.69%	77.55%	59.18%
Commonsense	52.86%	14.29%	64.29%	34.29%
Numerical Calc	52.50%	5.00%	62.50%	25.00%
Landmark	65.50%	37.50%	72.50%	51.00%
Position	50.00%	0.00%	51.67%	20.00%
Celebrity	94.12%	88.82%	93.24%	86.47%
Text Translation	52.50%	10.00%	50.00%	10.00%
OCR	90.00%	80.00%	77.50%	55.00%
OVERALL	63.98%	33.28%	75.02%	54.51%

Table 8: The performance of FTibVLM and Qwen3-VL-8B-Instruct (Base) across MME subtasks.

Model	Accuracy	Macro-P	Macro-R	Macro-F1	Prediction Distribution (0/1/2)
Base	0.3715	0.3813	0.3716	0.3716	13626/1971/2304
FTibVLM	0.5432	0.5703	0.5432	0.5400	6503/7974/3424

Table 9: The performance of FTibVLM and Qwen3-VL-8B-Instruct (Base) across SNLI-VE.

## E.2 Data and Model

**OCR training data.** We collect approximately **30k** (30,000) Tibetan OCR instruction instances from our in-house resources, and mix them into the existing multimodal instruction tuning (MIT) training data for continual training, in order to test the marginal benefit of *incremental OCR data*.

**Mixing strategy.** We train with a **3:5** mixture ratio between the OCR data and the MIT data.

**Compared model settings.** We evaluate the following three model configurations on our private OCR test set:

- **FTibVLM + OCR mixed-in training:** Starting from the existing three-stage trained model, we further train by mixing **30k** OCR instruction samples into the original instruction fine-tuning data.
- **Qwen3-VL-8B Instruct + OCR-only training on 30k:** Starting from the base instruction-tuned model, we train using **only** the **30k** OCR instruction dataset.

Model	CER	Exact Match
Base	2.1594	0.0010
Base + OCR-Only	0.3283	0.0907
CP + MA + OCR-MIT	0.2803	0.3617

Table 10: CER and Exact Match on the Tibetan OCR benchmark for different OCR training variants.

- **Base model (without the above OCR training):** Used as a lower-bound reference for OCR capability.

## E.3 Experiments Results

CER (Character Error Rate) is computed as the character-level edit distance divided by the total number of characters, where lower is better. Exact Match measures the proportion of samples whose predicted text matches the reference string exactly, where higher is better. As is shown in Figure 10, the results indicate a substantial gain in *line-level usability* after introducing OCR supervision: the base model almost never produces correct Tibetan text

on this OCR test set (Exact Match = 0.0010), while FTibVLM with mixed OCR supervision (OCR-Mix) increases Exact Match to 0.3617, i.e., an absolute improvement of +36.07 percentage points. This suggests that, after adding OCR data, the model can fully recognize entire text lines correctly for a considerable fraction of samples, leading to a tangible improvement in practical usability. In comparison, continuing training the base instruction model using only the 30k OCR dataset yields a smaller gain (Exact Match = 0.0907), and OCR-Mix achieves a clearly higher line-level accuracy.

Despite this, CER remains relatively high, implying that OCR performance is not yet stable or uniformly reliable across samples. Although OCR-Mix reduces CER to 0.2803, this value still indicates non-trivial character-level errors for many instances. Notably, the changes in CER and Exact Match are not perfectly aligned: the large jump in Exact Match resembles a shift where a subset of samples moves from “almost entirely wrong” to “entirely correct,” rather than a uniform reduction of character errors across all samples. This pattern typically suggests that training primarily benefits easier sub-distributions (e.g., clear fonts, regular layouts, higher resolution, and less background clutter), while difficult cases (low resolution, occlusion, complex backgrounds, and font/style variations) still suffer from frequent character mistakes. Moreover, the base model sometimes exhibits  $CER > 1$ , indicating extremely large character-level divergence from the target (e.g., many deletions/substitutions or irrelevant outputs), which is consistent with its near-zero Exact Match.

**E.4 Discussion and Implications**

This supplementary experiment suggests that *language-side adaptation alone* is insufficient to obtain stable Tibetan OCR capability; improving OCR still requires *dedicated supervision signals* targeting visual text recognition. Even with 30k OCR instruction instances, although line-level correctness increases markedly, the character error rate indicates that OCR remains far from “stable and reliable.” To systematically improve Tibetan OCR in future work, it may be necessary to incorporate:

- Larger-scale Tibetan OCR data that better matches real-world scene distributions.
- Higher-resolution inputs and stronger text-region alignment supervision, e.g., line-level

and box-level alignment, as well as text-region augmentation.

- OCR-oriented training strategies and model/component adaptations.

**F LLM-judge Prompt**

**F.1 LLM-judge Prompt for Translation Quality Control**

For the Tibetan translation quality-control setting, we used the prompt in Figure 4 to score an English→Tibetan translation. The evaluation was conducted using DeepSeek-V3 as the LLM-judge.

1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182

### LLM Judge Prompt for Translation Quality Control

You are a strict machine-translation quality inspector. Please score the following English→Tibetan translation, **strictly** following the rubric.

#### [Scoring Dimensions]

- 1) **Accuracy (0–2)**: 2 = no substantive semantic errors; 1 = minor deviation; 0 = clearly misleading.
- 2) **Completeness (0–2)**: 2 = no obvious omission/addition; 1 = minor omission/addition; 0 = harms the core meaning.
- 3) **Tibetan Expression (0–1)**: 1 = fluent and natural; 0 = disfluent/awkward/ambiguous.

#### [Output Requirements]

- Output **only** a single JSON object.
- Fields: accuracy, completeness, tibetan\_expression, total, comment.
- total = sum of the three scores.
- comment: a brief explanation in **English** ( ≤ 40 words).

Figure 4: Prompt for English→Tibetan translation quality scoring.