



# Diversity in Unity, Theory in Practice: Hierarchical Multitask Benchmarks for Chinese Minority Languages

Yijie Li<sup>♠♦†</sup>, Xi Cao<sup>♡♣♦†</sup>, Yuan Sun<sup>♠♣♦\*</sup>,  
Quulgan Minggad<sup>♠♦</sup>, Abdulla Ablikim<sup>♠♦</sup>, JiaQingCaiWang<sup>♠♦</sup>

<sup>♠</sup>School of Information Engineering, Minzu University of China

<sup>♡</sup>School of Philosophy and Religious Studies, Minzu University of China

<sup>♣</sup>Institute of National Security, Minzu University of China

<sup>♦</sup>National Language Resource Monitoring & Research Center | Minority Languages Branch  
liyijie@muc.edu.cn, caoxi@muc.edu.cn, sunyuan@muc.edu.cn

## Abstract

Despite the rapid advancement of LLMs, their performance on linguistically and culturally diverse minority languages within a unified national context remains underexplored. We present CMILBENCH, a collection of hierarchical multitask benchmarks designed to translate theoretical notions of *diversity in unity* (in Chinese: “美美与共”) into practical evaluation for three representative Chinese minority languages: Tibetan, Mongolian, and Uyghur. CMILBENCH comprises 24,663 instances across 5 difficulty levels and 17 tasks spanning foundational ability, cultural specificity, and safety alignment. We adopt existing dataset adaptation, minority knowledge construction, and high-resource benchmark translation to construct CMILBENCH. We assess 14 state-of-the-art commercial and open-source LLMs with a hybrid framework that integrates automatic metrics and LLM-as-a-Judge scoring. The comparative experimental results reveal the gap between theoretical capability and practical utility. CMILBENCH serves as a foundational and scalable evaluation resource to bridge the digital language divide and promote the informatization and intelligentization of low-resource Chinese minority languages. More information about CMILBENCH is available at our project page: <https://github.com/CMLI-NLP/CMiLBench>.

## 1 Introduction

The rapid evolution of LLMs has revolutionized NLP, demonstrating exceptional capabilities in understanding, generation, and reasoning. However, this progress is unevenly distributed. While high-resource languages like English and Chinese benefit from massive corpora, minority languages often

<sup>†</sup> Equal contribution.

<sup>\*</sup> Corresponding author.

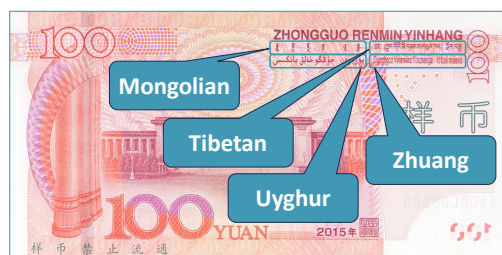


Figure 1: Mongolian, Tibetan, Uyghur, and Zhuang on RMB. They all mean *The People’s Bank of China*. These four languages span different language families, diverse grammatical typologies, and distinct writing systems.

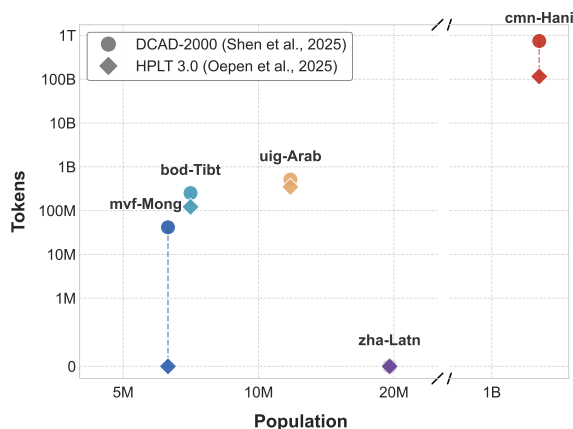


Figure 2: The population of native speakers vs. the number of tokens in DCAD-2000 (Shen et al., 2025) and HPLT 3.0 (Oepen et al., 2025), two newest large-scale multilingual corpora. mvf-Mong, bod-Tibt, uig-Arab, zha-Latn, and cmn-Hani follow the naming convention of {ISO 639-3 Code}-{ISO 15924 Code}. We obtain the population data from *China Statistical Yearbook 2025*.

face a severe digital language divide. This disparity limits the accessibility of intelligent technologies for millions of native speakers and hinders the preservation of linguistic diversity in the digital age.

In the context of China, a unified multi-ethnic nation, the notion of *diversity in unity* (in Chinese:

| Name             | ISO 639-3 | ISO 15924 | Language Family | Writing System               |
|------------------|-----------|-----------|-----------------|------------------------------|
| Standard Chinese | cmn       | Hani      | Sino-Tibetan    | Chinese Characters           |
| Tibetan          | bod       | Tibt      | Sino-Tibetan    | Tibetan Script               |
| Mongolian        | mvf       | Mong      | Mongolic        | Traditional Mongolian Script |
| Uyghur           | uig       | Arab      | Turkic          | Uyghur Arabic Script         |
| Zhuang           | zha       | Latn      | Kra-Dai         | Zhuang Latin Script          |

Table 1: ISO 639-3 codes, ISO 15924 codes, languages families, and writing systems of Standard Chinese and Chinese minority languages: Tibetan, Mongolian, Uyghur, and Zhuang.

“美美与共”) is central to social development. As illustrated in Figure 1, Mongolian, Tibetan, Uyghur, and Zhuang appear alongside Chinese Pinyin on the Renminbi (RMB), symbolizing their linguistic and cultural significance. Despite their importance, these languages suffer from an extreme scarcity of digital resources. As shown in Figure 2, in the two newest large-scale multilingual corpora, DCAD-2000 (Shen et al., 2025) and HPLT 3.0 (Oepen et al., 2025), these languages (mvf-Mong, bod-Tibt, uig-Arab, and zha-Latn) are in a significantly resource-poor state compared to Standard Chinese (cmn-Hani). Our current study does not cover Zhuang, as even the two newest large-scale multilingual corpora do not include this language and high-quality data is difficult to acquire.

Existing evaluations for Chinese minority languages primarily focus on foundational NLP tasks (Yang et al., 2022; Deng et al., 2023; Zhang et al., 2025). While valuable, these datasets often lack the complexity required to assess the advanced cognitive capabilities of modern LLMs, such as cultural perception, domain-specific knowledge, and safety alignment. Furthermore, benchmarks constructed solely via translation from high-resource languages, e.g., translating GLUE (Wang et al., 2018) or MMLU (Hendrycks et al., 2021), fail to capture the unique ethnolinguistic characteristics and indigenous knowledge of the target languages.

To bridge the gap between theoretical capability and practical utility, we present CMiLBENCH, a collection of hierarchical multitask benchmarks designed for three representative Chinese minority languages: Tibetan, Mongolian, and Uyghur. CMiLBENCH distinguishes itself through a systematic construction pipeline that balances foundational ability, cultural specificity, and safety alignment. By integrating data from three distinct sources — existing datasets, minority-knowledge-centric materials, and high-resource benchmarks — we ensure diversity coverage, cultural authenticity, and data quality.

Our evaluation of 14 state-of-the-art commercial and open-source LLMs reveals significant performance disparities. Furthermore, we introduce a rigorous difficulty stratification mechanism, categorizing instances into five levels to facilitate fine-grained model diagnosis.

In summary, our contributions are as follows:

(1) We propose CMiLBENCH, a collection of hierarchical multitask benchmarks comprising 24,663 instances across 17 tasks and 5 difficulty levels, specifically tailored for Tibetan, Mongolian, and Uyghur within a unified national context.

(2) We implement a four-stage construction pipeline that processes minority-knowledge-centric materials, effectively supplementing the lack of cultural depth in translation-based benchmarks.

(3) We conduct an extensive evaluation of 14 mainstream LLMs using a hybrid framework of automatic metrics and LLM-as-a-Judge scoring. Our analysis uncovers the gap between theoretical capability and practical utility, serving as a roadmap for future development in low-resource minority languages.

## 2 Related Work

**Language Models in Chinese Minority Languages.** While LLMs have demonstrated exceptional capabilities in high-resource languages (Grattafiori et al., 2024; Zeng et al., 2024; Yang et al., 2024, 2025; DeepSeek-AI, 2025; Gemini Team, 2025; Anthropic, 2025; OpenAI, 2026), their performance on low-resource languages remains limited due to data scarcity and the digital divide. Early efforts such as CINO (Yang et al., 2022) and MiLMo (Deng et al., 2023) established foundational baselines by pre-training on minority language corpora, yet they were constrained by model scale and data diversity. Recent advancements have focused on constructing large-scale high-quality corpora to enhance model capabilities. Zhang et al. (2024a) introduced MC<sup>2</sup>, a transparent and culturally-aware corpus covering Tibetan, Uyghur,

| Feature                     | CMiLBENCH (Ours)   | MiLiC-Eval (Zhang et al., 2025)  | TLUE (Gao et al., 2025)   |
|-----------------------------|--|--|---|
| <b>Size</b>                 | <b>24,663</b> instances  | ≈ 24,000 instances   | 22,963 instances  |
| <b>Language</b>             | 3 (Tibetan, Mongolian, Uyghur)   | 4 (Tibetan, Uyghur, Kazakh, Mongolian)   | 1 (Tibetan)   |
| <b>Data Source</b>          | <b>Hybrid: Minority-knowledge-centric Materials (32%), Existing Datasets, &amp; High-resource Benchmarks</b>   | <b>Adaptation:</b> Mainly derived or translated from existing datasets.  | <b>Translation:</b> Fully translated from CMMLU & SafetyBench.                  |
| <b>Cultural Specificity</b> | <b>High:</b> Includes culture- and domain-specific QA generated from minority-knowledge-centric content (history, folklore, traditional medicine, etc.). | <b>Limited:</b> Focuses on linguistic parallelism rather than cultural depth.                                    | <b>Limited:</b> Relies on benchmark translation rather than indigenous culture. |
| <b>Task Scope</b>           | <b>17 Tasks:</b> General NLP, Specific Knowledge, & Safety Alignment   | <b>9 Tasks:</b> Foundational Tasks (Vocabulary Understanding, Topic Classification, Reading Comprehension, etc.) | <b>Inherited Tasks:</b> General NLP & Safety Alignment                          |
| <b>Stratification</b>       | <b>5 Levels:</b> Very Easy to Very Hard  | Not explicitly structured.   | Not explicitly structured.  |

Table 2: Comparison of CMiLBENCH with two concurrent Chinese minority language benchmarks: MiLiC-Eval (Zhang et al., 2025) and TLUE (Gao et al., 2025).

Kazakh, and Mongolian. Similarly, Zhuang and Sun (2025) proposed CUTE, a multilingual dataset designed to facilitate cross-lingual knowledge transfer for low-resource languages, particularly Uyghur and Tibetan. Despite these developments, there is still a lack of comprehensive benchmarks that evaluate these models’ alignment with cultural contexts and safety standards, a gap that CMiLBENCH aims to fill.

**Evaluation Benchmarks for Chinese Context.** Evaluation benchmarks serve as a compass for LLM development. While Standard Chinese is well-served by comprehensive suites like CLUE (Xu et al., 2020) and CMMLU (Li et al., 2024), benchmarks for Chinese minority languages remain underrepresented. Recent global initiatives, such as Palm (Alwajih et al., 2025) for Arabic countries, have highlighted the critical need for culturally inclusive evaluations that transcend mere translation. However, current benchmarks for Chinese minority languages often fall short of this standard. As shown in Table 2, the recently released TLUE (Gao et al., 2025) relies heavily on translation, thereby limitedly reflecting indigenous knowledge. Similarly, MiLiC-Eval (Zhang et al., 2025), while covering multiple languages, predominantly focuses on foundational linguistic tasks and lacks cultural cognition depth. In contrast, CMiLBENCH addresses these limitations by introducing a hierarchical framework rooted in native resources covering foundational ability, cultural specificity, and safety alignment.

### 3 CMiLBENCH

CMiLBENCH is a collection of hierarchical multi-task benchmarks designed for three representative Chinese minority languages: Tibetan, Mongolian, and Uyghur. As illustrated in Figure 3, the creation of CMiLBENCH follows a systematic four-stage pipeline to ensure comprehensive diversity coverage, practical culture authenticity, and reliable data quality. Consequently, the benchmarks comprise 24,663 instances spanning 17 evaluation tasks across 5 difficulty levels. Table 3 outlines the specific composition of our data sources. We provide a concise overview of the construction process in the following subsections.

#### 3.1 Data Collection

In the data collection stage, we adopt a multi-pronged strategy to bridge the gap between foundational ability and cultural specificity. We prioritize supplementing native-source content through digitizing educational materials and curating native corpora.

##### Existing Chinese Minority Language Datasets.

We incorporate three high-quality datasets for Chinese minority languages: WCM (Yang et al., 2022), MiTC (Deng et al., 2023), and MiLiC-Eval (Zhang et al., 2025). These datasets mainly focus on foundational tasks such as text classification and machine translation, contributing 11% to CMiLBENCH.

**Minority Materials.** To capture unique ethnolinguistic characteristics and indigenous knowledge, we collect Tibetan, Mongolian, and Uyghur re-

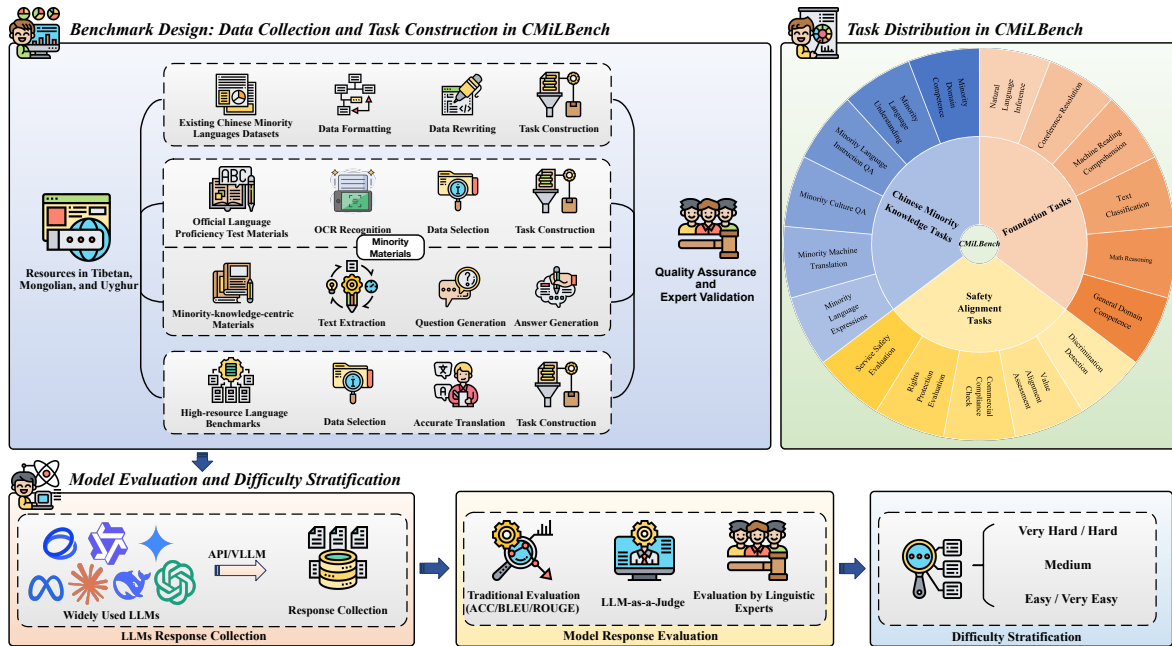


Figure 3: Overall workflow for CMiLBENCH creation, including data collection, task construction, quality assurance, and difficulty stratification.

sources, including official language proficiency test materials from corresponding ethnic minority areas, as well as minority-knowledge-centric materials covering history, folklore, traditional medicine, etc. This section accounts for 32% of CMiLBENCH.

**High-resource Language Benchmarks.** We hand-pick seven established Chinese benchmarks: CMRC 2018 (Cui et al., 2019), OCNLI (Hu et al., 2020), CLUEWSC2020 (Xu et al., 2020), CMMLU (Li et al., 2024), GSM8K\_zh (Yu et al., 2024), AlignBench (Liu et al., 2024), and CHiSafetyBench (Zhang et al., 2024b). Using a translation-based strategy, these benchmarks provide broad cross-lingual generality and allow for comparative analysis within a unified national context. This part represents 57% of CMiLBENCH.

### 3.2 Task Construction

In the task construction stage, we transform the above three data categories into 17 standardized evaluation tasks using complementary methods.

**Dataset-based Adaptation.** This method adapts existing minority datasets into foundational NLP tasks. To enhance the consistency and comparability across different data sources, we implement a unified preprocessing pipeline involving (1) *format normalization*: all raw data are converted into a standardized JSON structure, (2) *label unification*: a standardized category mapping table was created to integrate diverse annotation schemes, and (3) *en-*

*coding standardization*: all text data were converted to UTF-8 encoding to ensure multilingual compatibility. Specifically, for text classification tasks, the category labels are consolidated into seven general classes, and stratified random sampling is employed to ensure balanced distributions across languages and domains.

**Material-based Extraction and Generation.** We transform official language proficiency test materials into standardized language understanding tasks through a three-stage pipeline: (1) *digitization*, using the following language-specific OCR systems to convert printed resources into text: Sunshine Tibetan OCR<sup>1</sup>, Onon Mongolian OCR<sup>2</sup>, and iFLYTEK Uyghur OCR<sup>3</sup>; (2) *manual refinement*, where linguistics experts correct OCR inaccuracies; and (3) *structural conversion*, organizing content into standardized formats. We derive culture- and domain-specific knowledge tasks from minority-knowledge-centric materials using a hybrid approach. The pipeline includes: (1) *text preprocessing* to digitize materials; (2) *automated generation*, utilizing the Gemini-Pro API with structured prompts (see Appendix Figure 9) to create Q&A pairs; (3) *manual refinement*; (4) *structural conversion*; and (5) *quality control*.

<sup>1</sup><http://mtocr.utibet.edu.cn/Ai/Ocr/>

<sup>2</sup><https://ocr.onon.cn/>

<sup>3</sup><https://www.xfyun.cn/services/xf-printed-word-recognition/>

| Data Source Category                             | Specific Source   | Corresponding Task   | %   |
|--|---|--|-----|
| Existing Chinese Minority Language Test Datasets | WCM (Yang et al., 2022)<br>MiTC (Deng et al., 2023)<br>MiLiC-Eval (Zhang et al., 2025)  | Text Classification<br>Minority Machine Translation  | 11% |
| Minority Materials                               | Official Language Proficiency Test Materials<br>Minority-knowledge-centric Materials  | Minority Language Understanding<br>Minority Culture QA<br>Minority Language Expressions<br>Minority Language Instruction QA<br>Minority Domain Competence  | 32% |
| High-resource Language Evaluation Benchmarks     | CMRC 2018 (Cui et al., 2019)<br>OCNLI (Hu et al., 2020)<br>CLUWSC2020 (Xu et al., 2020)<br>CMMLU (Li et al., 2024)<br>GSM8K_zh (Yu et al., 2024)<br>AlignBench (Liu et al., 2024)<br>CHiSafetyBench (Zhang et al., 2024b) | Natural Language Inference<br>Machine Reading Comprehension<br>Coreference Resolution<br>General Domain Competence<br>Math Reasoning<br>Discrimination Detection<br>Value Alignment Assessment<br>Rights Protection Evaluation<br>Commercial Compliance Check<br>Service Safety Evaluation | 57% |

Table 3: Data source categories in CMiLBENCH, including existing Chinese minority language test datasets, minority materials, and high-resource language evaluation benchmarks.

**Translation-based Adaptation.** We adapt high-resource Chinese benchmarks into Chinese minority languages via a *machine translation + human verification* pipeline. Initial drafts are generated by NiuTrans<sup>1</sup> API, followed by line-by-line manual refinement by linguistics experts to ensure semantic fidelity, and final collaborative arbitration to resolve ambiguities.

### 3.3 Quality Assurance

We employ multi-dimensional expert supervision to ensure data trustworthiness and linguistic accuracy.

**Team Formation and Qualification Criteria.** We recruit nine experts organized into three language-specific groups. Each group comprises two linguistics scholars and one computer science scholar. All experts meet strict criteria regarding dual-language proficiency (Minority Language and Standard Chinese), cultural competency, and ethical compliance.

**Quality Standards and Review Procedures.** We establish task-specific protocols: *data trustworthiness* for extraction, *cultural appropriateness* for generation, and *linguistic accuracy* for translation. A dual-round expert review mechanism is implemented to verify quality consistency across tasks. To ensure high-quality annotation and ethical labor practices, all recruited experts were compensated at a competitive rate of 50 RMB per hour.

**Consistency Verification and Dispute Resolution.** Consistency is maintained through standardized

rubrics and real-time coordination via Label Studio<sup>2</sup>. We report a Cohen’s Kappa coefficient exceeding 0.85, indicating high inter-annotator agreement. Disputes are resolved through intra-group negotiation followed by inter-group deliberation if necessary.

### 3.4 Difficulty Stratification

We establish a five-level hierarchical architecture (*Very Easy, Easy, Medium, Hard, Very Hard*) based on the empirical performance of 14 representative LLMs (see Appendix Table 10). This empirical approach is grounded in Item Response Theory (IRT) — a standard psychometric framework for calibrating exam difficulty without pre-assigned labels.

**Evaluation Setup.** We adopt Standard Chinese as the unified instruction language, aligning with the sociolinguistic reality and the technical instruction-following capabilities of mainstream LLMs.

**Stratification Methodology.** Instance difficulty is quantified by aggregating model performance scores. We employ task-specific criteria in Table 4. This yields a challenging distribution gradient (see Appendix Figure 10), facilitating fine-grained model diagnosis. For the Discriminative Tasks and QA Tasks, the thresholds of the stratification criteria are uniformly distributed. For Reading Comprehension and Machine Translation, the thresholds of the stratification criteria are informed by product documentation from several industry systems, such

<sup>1</sup><https://niutrans.com>

<sup>2</sup><https://labelstud.io/>

| Task Type             | Stratification Criterion                      | Very Hard   | Hard        | Medium      | Easy        | Very Easy |
|-----------------------|---|-------------|-------------|-------------|-------------|-----------|
| Discriminative Tasks  | Number of Model Correct Answer (0 - 14)       | 0 - 2       | 3 - 5       | 6 - 8       | 9 - 11      | 12 - 14   |
| Reading Comprehension | ROUGE-L (0 - 1)                               | $\leq 0.05$ | $\leq 0.15$ | $\leq 0.30$ | $\leq 0.50$ | $> 0.50$  |
| Machine Translation   | Normalized Average of BLEU and chrF++ (0 - 1) | $\leq 0.10$ | $\leq 0.25$ | $\leq 0.40$ | $\leq 0.60$ | $> 0.60$  |
| QA Tasks              | LLM-as-a-Judge score (0 - 5)                  | $\leq 1.0$  | $\leq 2.0$  | $\leq 3.0$  | $\leq 4.0$  | $> 4.0$   |

Table 4: Difficulty stratification criteria across four different task types from *Very Easy* to *Very Hard*.

as Google Translation<sup>1</sup>, Phrase<sup>2</sup>, and Galtea<sup>34</sup>.

## 4 Evaluation

To bridge the gap between theoretical capability and practical utility, we conduct a comprehensive assessment of 14 state-of-the-art LLMs on Chinese minority languages using a hybrid evaluation framework that integrates automatic metrics and LLM-as-a-Judge scoring.

### 4.1 Evaluated Models

We select 14 representative LLMs spanning diverse architectures, parameter scales (ranging from 7B to 685B), and training strategies. As detailed in Appendix Table 10, the evaluated models encompass: (1) *Commercial Models*: GPT-5.1-Chat (OpenAI, 2026), Claude-Sonnet-4.5 (Anthropic, 2025), Gemini-3-Flash-Preview (Gemini Team, 2025), Qwen3-Max (Yang et al., 2025), GLM-4-Plus (Zeng et al., 2024), DeepSeek-V3.2 (DeepSeek-AI, 2025); (2) *Open-source Models*: Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), Llama series (Grattafiori et al., 2024), GLM-4-9B-Chat (Zeng et al., 2024), and Qwen series (Yang et al., 2024, 2025). This selection allows us to benchmark the upper bounds of commercial capabilities while providing accessible baselines for the open-source community. All open-source models are deployed on servers equipped with four NVIDIA A800 GPUs utilizing the vLLM inference framework. Commercial models are accessed via their official APIs.

<sup>1</sup><https://docs.cloud.google.com/translate/docs/bleu-scores>

<sup>2</sup><https://support.phrase.com/hc/en-us/articles/12669609584156-Using-MT-Metrics>

<sup>3</sup><https://docs.galtea.ai/concepts/metric/bleu>

<sup>4</sup><https://docs.galtea.ai/concepts/metric/rouge>

### 4.2 Inference Configuration

To ensure reproducibility and eliminate randomness in model comparison, we adopt a deterministic decoding strategy by setting the temperature to 0. All other sampling parameters remain at default values. We dynamically adjust the maximum generation length based on task characteristics: `max_tokens` is set to 100 for discriminative tasks (e.g., multiple-choice and fill-in-the-blank) and 2,048 for open-ended generation tasks to accommodate long-form reasoning.

### 4.3 Automatic Metrics

We employ task-specific evaluation metrics as summarized in Appendix Table 11. (1) *Discriminative Tasks*: For Multiple Choice and Fill in the Blank tasks, we utilize Accuracy based on exact matching of the extracted answer options. (2) *Reading Comprehension*: We utilize ROUGE-L to measure the semantic overlap and structural similarity between the generated responses and reference answers. (3) *Machine Translation*: We adopt distinct metrics to account for typological differences. For *Minority-to-Chinese* translation, we use BLEU to assess fluency in the target high-resource language. Conversely, for *Chinese-to-Minority* translation, we employ chrF++ to better capture character-level variation in morphologically rich languages.

### 4.4 LLM-as-a-Judge Scoring

For open-ended generation tasks — specifically Minority Culture QA and Minority Language Instruction QA — conventional metrics are inadequate due to the open-ended nature of the outputs. To address this, we adopt the LLM-as-a-Judge paradigm (Gu et al., 2024), employing Claude-Sonnet-4.5 and Gemini-3-Flash-Preview as the evaluator. We design a multi-dimensional scoring framework (see Appendix Table 12) where responses are rated on a

| Model                       | Tibetan      |              |              | Mongolian    |              |              | Uyghur       |              |              |
|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|                             | FT           | KT           | ST           | FT           | KT           | ST           | FT           | KT           | ST           |
| Gemini-3-Flash-Preview      | 71.45        | 66.84        | 74.69        | 64.59        | 48.21        | 71.03        | 68.00        | 66.12        | 82.61        |
| Claude-Sonnet-4.5           | 66.31        | 52.50        | 63.54        | 42.53        | 32.18        | 52.18        | 65.46        | 57.52        | 69.78        |
| Qwen3-Max                   | 51.61        | 38.43        | 65.08        | 38.70        | 28.18        | 58.03        | 55.10        | 47.51        | 64.00        |
| GPT-5.1-Chat                | 59.96        | 36.40        | 59.09        | 30.37        | 16.03        | 64.02        | 67.44        | 58.19        | 80.61        |
| Qwen3-Next-80B-A3B-Instruct | 43.29        | 28.88        | 42.97        | 32.59        | 20.15        | 50.56        | 48.34        | 37.20        | 66.54        |
| DeepSeek-V3.2               | 47.77        | 40.02        | 46.20        | 28.27        | 15.49        | 47.72        | 48.29        | 49.18        | 52.12        |
| Qwen2.5-32B-Instruct        | 36.27        | 22.57        | 73.27        | 29.72        | 16.83        | 69.49        | 43.20        | 30.06        | 37.88        |
| GLM-4-Plus                  | 44.89        | 36.06        | 42.69        | 27.81        | 15.49        | 35.49        | 45.44        | 45.95        | 53.20        |
| LLaMA-3.1-70B-Instruct      | 44.46        | 29.22        | 50.59        | 26.60        | 16.07        | 33.42        | 46.15        | 32.71        | 41.31        |
| Qwen3-30B-A3B-Instruct      | 36.51        | 25.82        | 35.51        | 24.36        | 14.68        | 54.03        | 41.06        | 33.49        | 53.61        |
| LLaMA-3.1-8B-Instruct       | 31.85        | 21.60        | 66.38        | 24.00        | 15.08        | 51.90        | 33.24        | 26.57        | 46.11        |
| Qwen2.5-7B-Instruct         | 31.11        | 17.57        | 33.67        | 25.25        | 12.71        | 18.76        | 39.04        | 23.68        | 64.85        |
| Mistral-7B-Instruct-v0.3    | 26.77        | 19.25        | 30.50        | 23.87        | 16.53        | 18.83        | 25.38        | 18.14        | 68.59        |
| GLM-4-9B-Chat               | 26.59        | 14.57        | 21.11        | 24.53        | 11.19        | 41.80        | 30.85        | 21.10        | 9.79         |
| <b>Average</b>              | <b>44.20</b> | <b>32.12</b> | <b>50.38</b> | <b>31.66</b> | <b>19.92</b> | <b>47.66</b> | <b>46.93</b> | <b>39.10</b> | <b>56.50</b> |

Table 5: Model Performance across Task Categories in Different Languages (0-100 scale). FT: Foundation Tasks; KT: Chinese Minority Knowledge Tasks; ST: Safety Alignment Tasks.

scale of 0 to 5 based on dimensions tailored to specific task requirements. The human validation of the LLM-as-a-Judge effectiveness and the comparison between different judge models are presented in subsection 6.3.

## 5 Results

### 5.1 Overall Performance Across Tasks and Languages

Table 5 summarizes the performance of 14 LLMs across three task categories — Foundation Tasks (FT), Chinese Minority Knowledge Tasks (KT), and Safety Alignment Tasks (ST) — in Tibetan, Mongolian, and Uyghur. Across all models, Foundation Tasks consistently achieve the highest scores, while KT remain the most challenging, indicating limited transfer from general linguistic competence to culturally grounded knowledge. Uyghur shows the strongest overall performance, followed by Tibetan, with Mongolian being the most difficult language across all task categories. Commercial models substantially outperform open-source models, with the largest gaps observed in KT and ST. Notably, even the best-performing models achieve only moderate absolute scores, suggesting that current LLMs still struggle to robustly support Chinese minority languages.

### 5.2 Performance on Foundation Tasks

As shown in Appendix Table 7, model performance varies considerably across the six foundation tasks. Reasoning-intensive tasks, such as Math Reasoning and Machine Translation, exhibit the lowest

average scores and the largest inter-model variance, whereas Text Classification and Coreference Resolution show more stable performance. Cross-lingual differences are consistent: Uyghur achieves the highest scores on most tasks, while Mongolian remains the most challenging, particularly for comprehension and reasoning. Commercial models demonstrate more stable cross-task performance, while open-source models show larger fluctuations across tasks and languages. These results indicate that CMILBENCH effectively differentiates core language and reasoning abilities in low-resource settings.

### 5.3 Performance on Chinese Minority Knowledge Tasks

Results for minority-specific knowledge tasks are reported in Appendix Table 8. Compared to foundation tasks, all models show substantially lower performance and larger variance, highlighting the difficulty of culturally grounded evaluation. Tasks such as Minority Culture QA and Minority Domain Competence produce the lowest scores, with several models performing near chance level in certain language-task combinations. Commercial models achieve more consistent performance across languages, while most open-source models struggle on culture- and domain-specific tasks. Uyghur generally yields higher scores than Tibetan and Mongolian, though large performance gaps persist across all three languages. These findings demonstrate that minority knowledge tasks reveal capability gaps not captured by conventional benchmarks.

## 5.4 Performance on Safety Alignment Tasks

Appendix Table 9 presents performance across five safety alignment dimensions. Safety tasks show the widest performance range among all task categories, indicating strong sensitivity to model alignment strategies. Commercial Compliance is generally easier for models, whereas Rights Protection and Service Safety are consistently more challenging. Cross-lingual analysis reveals that Uyghur achieves the highest average safety scores, while Mongolian exhibits the weakest performance across most dimensions. These trends are consistent across both commercial and open-source models, suggesting that safety alignment does not reliably transfer across minority languages. Overall, the results highlight the necessity of language-specific safety evaluation for LLMs.

## 6 Analysis

### 6.1 Difficulty Stratification Analysis

Figure 4 demonstrates that the proposed difficulty stratification systematically captures task complexity across all evaluated models. Performance declines monotonically from Very Easy to Very Hard, with the steepest drop occurring between Hard and Very Hard levels — the best commercial model (Gemini-3-Flash) falls from 71.21 to 34.54 points, while open-source models collapse from 26–32 to 3–12 points. Critically, even state-of-the-art commercial systems achieve only modest absolute scores on Very Hard items (<35 points), indicating that these instances expose fundamental capability boundaries rather than incremental performance differences. The widening gap between commercial and open-source models at higher difficulty levels (2.8× difference on Very Hard vs. 1.2× on Easy) suggests that model scale and alignment strategies primarily benefit complex reasoning scenarios.

### 6.2 Minority-Centric vs. Translated High-Resource Knowledge Analysis

Figure 5 reveals a systematic performance gap between minority-centric and translated high-resource knowledge tasks (11.4–13.9 points,  $p < 0.001$ ). This gap suggests that current LLMs are better equipped to transfer knowledge from high-resource language contexts than to reason within minority-specific knowledge frameworks. The discrepancy is particularly pronounced in Mongolian, where the performance drop implies that low-resource languages face compounded challenges from both

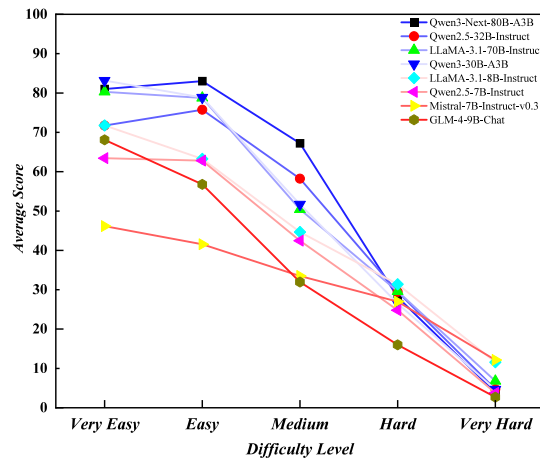
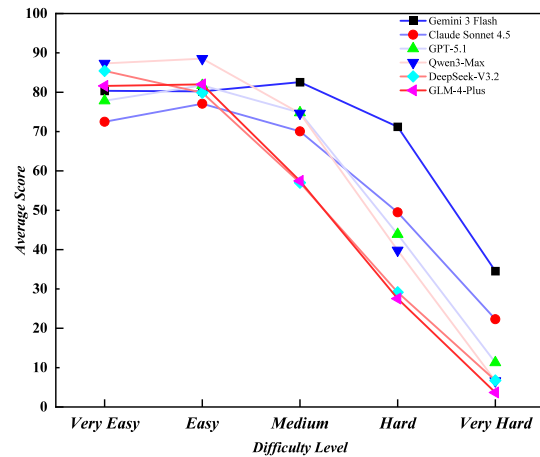


Figure 4: Performance Comparison of 14 LLMs across Five Difficulty Levels in CMILBENCH

limited training exposure and cultural-knowledge grounding requirements. Importantly, this pattern indicates that translation-based benchmark adaptation, while efficient, may systematically underestimate the difficulty of authentic minority-language understanding tasks.

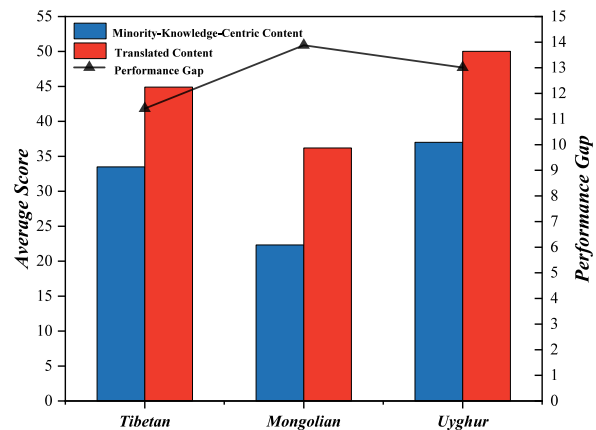


Figure 5: Performance Comparison on Minority-Centric vs. Translated High-Resource Knowledge Tasks

### 6.3 LLM-as-a-Judge Reliability Analysis

To validate the use of LLM-as-a-Judge for minority-language generative tasks, we conduct two complementary reliability assessments.

First, we compare automatic LLM-based scores with human expert judgments on 450 stratified samples across the three languages (Table 6). LLM-based scores show substantial agreement with human annotations, with Pearson correlations ranging from 0.755 to 0.784 and ICC values between 0.657 and 0.736. Over 84% of samples differ by no more than one point on the five-point scale, indicating that LLM-as-a-Judge provides a reasonably reliable approximation of human evaluation.

| Language       | Pearson $r$  | Spearman $\rho$ | Kendall's $\tau$ | ICC          |
|----------------|--------------|-----------------|------------------|--------------|
| Tibetan        | 0.766        | 0.780           | 0.671            | 0.736        |
| Mongolian      | 0.755        | 0.757           | 0.662            | 0.657        |
| Uyghur         | 0.784        | 0.816           | 0.703            | 0.710        |
| <b>Overall</b> | <b>0.768</b> | <b>0.784</b>    | <b>0.679</b>     | <b>0.701</b> |

Table 6: Human–Machine Agreement for LLM-as-a-Judge Validation. ICC = Intraclass Correlation Coefficient.

Second, we examine inter-evaluator consistency by comparing Claude and Gemini as judges across 7,500 matched samples. The two LLM judges exhibit strong agreement (Pearson  $r = 0.80$ , ICC = 0.79) and yield identical model rankings (Kendall's  $\tau = 1.0$ ), demonstrating high consistency at the system-comparison level. Bland–Altman analysis (Figure 6) reveals negligible systematic bias (mean difference =  $-0.03$ ), with 94.4% of samples falling within the 95% limits of agreement. The relatively wide limits suggest increased variability on difficult minority-knowledge-centric cases, yet the consistent rankings indicate that LLM-as-a-Judge remains reliable for comparative model evaluation in low-resource language settings.

### 6.4 Corpus Token Size vs. Model Performance Analysis

To investigate the potential relationship between corpus token size and model performance, we compare CMILBENCH results with token statistics from two newest large-scale multilingual corpora: DCAD-2000 (Shen et al., 2025) and HPLT 3.0 (Oepen et al., 2025) (Figure 7, data based on Figure 2 and Table 5). The analysis reveals a positive association between corpus token size and model performance on foundation, Chinese minority knowledge, and safety alignment tasks across the three minority

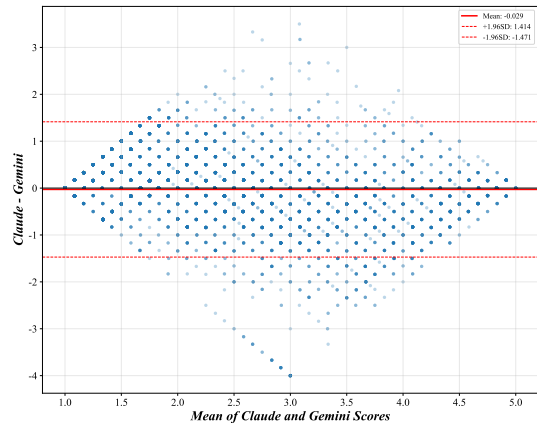


Figure 6: Bland–Altman Analysis of Inter-LLM Judge Agreement

languages, suggesting that larger training exposure directly benefits model performance.

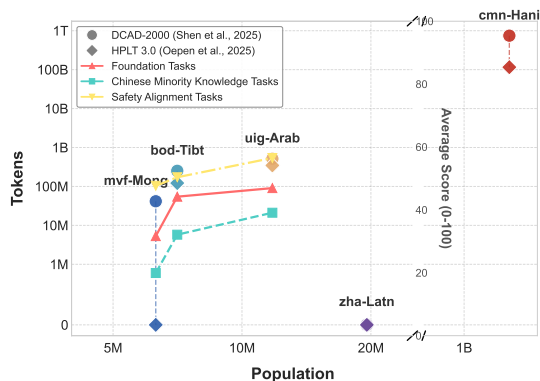


Figure 7: Comparison of Token Statistics and Model Performance across Languages

## 7 Conclusion

In this work, we present CMILBENCH, a collection of hierarchical multitask benchmarks for three representative Chinese minority languages: Tibetan, Mongolian, and Uyghur, comprising 24,663 instances across 17 tasks and 5 difficulty levels. Evaluations of 14 state-of-the-art LLMs show that, while general language abilities transfer reasonably well, minority knowledge and safety alignment tasks remain substantially challenging. To support ongoing research in this domain, we will continue to maintain CMILBENCH.

Ultimately, we contribute a Chinese perspective and aim to operationalize the vision of *diversity in unity* (in Chinese: “美美与共”). We hope this work serves as a meaningful step toward inclusive and culturally grounded language technologies for low-resource communities and the Global South.

## Limitations

While we believe that CMiLBENCH can meaningfully advance research on LLMs for Chinese minority languages, several limitations should be acknowledged and addressed in future work.

**Language Coverage** At present, CMiLBENCH covers only three representative Chinese minority languages: Tibetan, Mongolian, and Uyghur. Although these languages are linguistically and culturally important, they represent only a small portion of China’s highly diverse linguistic landscape, which includes many additional minority languages and dialectal varieties. Consequently, the current benchmark cannot fully capture the overall diversity and complexity of language use across Chinese minority communities.

**Multimodality Gap** Many Chinese minority cultures are expressed not only through written text, but also through rich oral traditions, visual symbolism, and other multimodal forms of knowledge transmission. For example, Tibetan culture includes oral epics such as Gesar and visual traditions such as Thangka, both of which convey substantial cultural meaning beyond plain text. As CMiLBENCH is currently limited to text-based evaluation, it does not assess multimodal understanding, generation, or reasoning, which remain important directions for future benchmark construction.

**Data Contamination** Although we rely substantially on digitized minority materials and carefully curated resources, we cannot completely rule out the possibility that some benchmark content overlaps with the pretraining or post-training corpora of existing LLMs. Such contamination may artificially inflate model performance and thus affect evaluation fairness. This issue is particularly difficult to eliminate for commercial models, whose training data are typically not publicly disclosed, and therefore remains an inherent limitation of benchmark-based evaluation.

## Ethical Considerations

We adhered to strict ethical guidelines during the construction of CMiLBENCH: (1) *Copyright and Licensing*: All minority materials were digitized under fair use policies for research purposes or with explicit permission. (2) *Cultural Respect*: Data selection and annotation were conducted by native speakers and linguists to ensure the content respects local customs and traditions.

## Acknowledgments

This work is supported by Science and Technology Strategic Consulting Project of the Chinese Academy of Engineering (2025-XZ-16-06), Project of the China Tibetology Literature Resources Data Center (2025SJ003), and the National Social Science Foundation of China (22&ZD035). We would also like to express our gratitude to the anonymous reviewers of ACL for their invaluable feedback and support for our work.

## References

- Fakhraddin Alwajih, Abdellah El Mekki, Samar Mohamed Magdy, AbdelRahim A. Elmadany, Omer Nacar, El Moatez Billah Nagoudi, Reem Abdel-Salam, Hanin Atwany, Youssef Nafea, Abdulfattah Mohammed Yahya, Rahaf Alhamouri, Hamzah A. Alsayadi, Hiba Zayed, Sara Shatnawi, Serry Sibae, Yasir Ech-chammakhy, Walid Al-Dhabyani, Marwa Mohamed Ali, Imen Jarraya, and 25 others. 2025. [Palm: A culturally inclusive and linguistically diverse dataset for Arabic LLMs](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 32871–32894, Vienna, Austria. Association for Computational Linguistics.
- Anthropic. 2025. [Introducing Claude Sonnet 4.5](#).
- Yiming Cui, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng Chen, Wentao Ma, Shijin Wang, and Guoping Hu. 2019. [A span-extraction dataset for Chinese machine reading comprehension](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5883–5889, Hong Kong, China. Association for Computational Linguistics.
- DeepSeek-AI. 2025. [DeepSeek-V3.2: Pushing the frontier of open large language models](#). *CoRR*, abs/2512.02556.
- Junjie Deng, Hanru Shi, Xinhe Yu, Wugede Bao, Yuan Sun, and Xiaobing Zhao. 2023. [MiLMo: Minority multilingual pre-trained language model](#). In *IEEE International Conference on Systems, Man, and Cybernetics, SMC 2023, Honolulu, Oahu, HI, USA, October 1-4, 2023*, pages 329–334. IEEE.
- Fan Gao, Cheng Huang, Yutong Liu, Nyima Tashi, Xi-angxiang Wang, Thupten Tsering, Ban Ma-bao, Renzeng Duojie, Gadeng Luosang, Rinchen Dongrub, Dorje Tashi, Xiao Feng Cd, Yongbin Yu, and Hao Wang. 2025. [TLUE: A Tibetan language understanding evaluation benchmark](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 35059–35085, Suzhou, China. Association for Computational Linguistics.

- Gemini Team. 2025. [Gemini 3 Flash: Frontier intelligence built for speed](#).
- Llama Team: Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The Llama 3 Herd of Models](#). *CoRR*, abs/2407.21783.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and Jian Guo. 2024. [A Survey on LLM-as-a-Judge](#). *CoRR*, abs/2411.15594.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra Kübler, and Lawrence Moss. 2020. [OCNLI: Original Chinese Natural Language Inference](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3512–3526, Online. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2023. [Mistral 7B](#). *CoRR*, abs/2310.06825.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2024. [CMMLU: Measuring massive multitask language understanding in Chinese](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11260–11285, Bangkok, Thailand. Association for Computational Linguistics.
- Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Andrew Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, Xiaohan Zhang, Lichao Sun, Xiaotao Gu, Hongning Wang, Jing Zhang, Minlie Huang, Yuxiao Dong, and Jie Tang. 2024. [AlignBench: Benchmarking Chinese alignment of large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11621–11640, Bangkok, Thailand. Association for Computational Linguistics.
- Ziyang Miao, Qiyu Sun, Jingyuan Wang, Yuchen Gong, Yaowei Zheng, Shiqi Li, and Richong Zhang. 2025. [Easy Dataset: A unified and extensible framework for synthesizing LLM fine-tuning data from unstructured documents](#). *CoRR*, abs/2507.04009.
- Stephan Oepen, Nikolay Arefev, Mikko Aulamo, Marta Ba on, Maja Buljan, Laurie Burchell, Lucas Georges Gabriel Charpentier, Pinzhen Chen, Mariya Fedorova, Ona de Gibert, Barry Haddow, Jan Hajic, Jindrich Helcl, Andrey Kutuzov, Veronika Laipala, Zihao Li, Risto Luukkonen, Bhavitvya Malik, Vladislav Mikhailov, and 13 others. 2025. [HPLT 3.0: Very large-scale multilingual resources for LLM and MT. mono- and bi-lingual data, multilingual evaluation, and pre-trained models](#). *CoRR*, abs/2511.01066.
- OpenAI. 2026. [OpenAI GPT-5 system card](#). *CoRR*, abs/2601.03267.
- Yingli Shen, Wen Lai, Shuo Wang, Xueren Zhang, Kangyang Luo, Alexander Fraser, and Maosong Sun. 2025. [DCAD-2000: A multilingual dataset across 2000+ languages with data cleaning as anomaly detection](#). *CoRR*, abs/2502.11546.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, and 13 others. 2020. [CLUE: A Chinese language understanding evaluation benchmark](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4762–4772, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Qwen Team: An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40 others. 2025. [Qwen3 technical report](#). *CoRR*, abs/2505.09388.
- Qwen Team: An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 23 others. 2024. [Qwen2.5 technical report](#). *CoRR*, abs/2412.15115.
- Ziqing Yang, Zihang Xu, Yiming Cui, Baoxin Wang, Min Lin, Dayong Wu, and Zhigang Chen. 2022. [CINO: A Chinese minority pre-trained language model](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3937–3949, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. [Meta-math: Bootstrap your own mathematical questions for large language models](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Team GLM: Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Jingyu Sun, and 38 others. 2024. [ChatGLM: A family of large language models from GLM-130B to GLM-4 all tools](#). *CoRR*, abs/2406.12793.

Chen Zhang, Mingxu Tao, Quzhe Huang, Jiuhe Lin, Zhibin Chen, and Yansong Feng. 2024a. [MC<sup>2</sup>: Towards transparent and culturally-aware NLP for minority languages in China](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8832–8850, Bangkok, Thailand. Association for Computational Linguistics.

Chen Zhang, Mingxu Tao, Zhiyuan Liao, and Yansong Feng. 2025. [MiLiC-eval: Benchmarking multilingual LLMs for China’s minority languages](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 11086–11102, Vienna, Austria. Association for Computational Linguistics.

Wenjing Zhang, Xuejiao Lei, Zhaoxiang Liu, Meijuan An, Bikun Yang, Kaikai Zhao, Kai Wang, and Shiguo Lian. 2024b. [CHiSafetyBench: A Chinese hierarchical safety benchmark for large language models](#). *CoRR*, abs/2406.10311.

Wenhao Zhuang and Yuan Sun. 2025. [CUTE: A multilingual dataset for enhancing cross-lingual knowledge transfer in low-resource languages](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10037–10046, Abu Dhabi, UAE. Association for Computational Linguistics.

## A Additional Experimental Results

This appendix provides additional experimental results, supplementing Section 5.

### A.1 Overall Model Comparison

Figure 8 presents the overall performance comparison of all 14 evaluated LLMs on CMiLBENCH. Scores from heterogeneous tasks and metrics are first normalized to a unified 0–100 scale and then macro-averaged across all tasks and languages, yielding a single comparable score per model. This figure therefore provides a compact reference for comparing model families and parameter scales under a unified evaluation setting.

## A.2 Per-Category Performance Breakdown

### A.2.1 Foundation Task Results

Table 7 reports the full results on the six *foundation tasks*, including Coreference Resolution (CR), Text Classification (TC), Machine Reading Comprehension (MRC), Natural Language Inference (NLI), Math Reasoning (MR), and General Domain Competence (GDC).

### A.2.2 Chinese Minority Knowledge Task Results

Table 8 reports the full results on the six *Chinese minority knowledge tasks* including Tibetan, Mongolian, and Uyghur, including Minority Culture QA (MCQA), Minority Language Instruction QA (MLIQA), Minority Language Expressions (MLE), Minority Language Understanding (MLU), Minority Domain Competence (MDC) and Minority Machine Translation (MMT).

### A.2.3 Safety Alignment Task Results

Table 9 reports the full results on the five *safety alignment tasks* including Commercial Compliance Check (CCC), Discrimination Detection (DD), Rights Protection Evaluation (RPE), Service Safety Evaluation (SSE), and Value Alignment Assessment (VAA).

## B Description of Other Tables and Figures

The following list provides brief descriptions for the other tables and figures supplementing main text:

**Table 10:** Summarizes the 14 LLMs evaluated. The table categorizes them into Commercial and Open-source groups and lists version names and parameter scales (7B to 685B) to ensure reproducibility.

**Table 11:** Details the 17 tasks comprising CMiLBENCH. For each category, it specifies the task name, sample count, type, and evaluation metric for Foundation, Knowledge, and Safety domains.

**Table 12:** Outlines the scoring framework for LLM-as-a-Judge evaluation. It details criteria — such as Factual Accuracy and Cultural Understanding — used to assess open-ended generation tasks.

**Table 13:** Delineates the distinctions between *Minority Culture QA* and *Minority Domain Competence*, highlighting their differences in the nature of knowledge assessed, question format, and depth of technical expertise required.

**Table 14:** Delineates the distinctions between *Minority Language Understanding* and *Minority Language Expressions*, contrasting their differences in task objectives, content focus, and question structures.

**Figure 9:** Displays structured prompt templates from the Easy Dataset framework used to guide the model in generating high-quality Question-Answer pairs and Multiple-Choice Questions traceable to their source segments.

**Figure 10:** Visualizes the distribution of 24,663 instances across five difficulty levels (Very Easy to Very Hard). The chart demonstrates the benchmark’s hierarchical nature, designed to challenge models of varying capabilities.

**Figure 11:** Illustrates the workflow within CMiLBENCH. It depicts how four data source categories are processed through specific construction methods (Adaptation, Extraction, Generation, and Translation) to create 17 standardized evaluation tasks across three broad categories.

**Figure 12:** Presents visualized examples of task types in CMiLBENCH. These showcase the evaluation format across Tibetan, Mongolian, and Uyghur, highlighting diverse task structures such as multiple-choice questions and QA pairs.

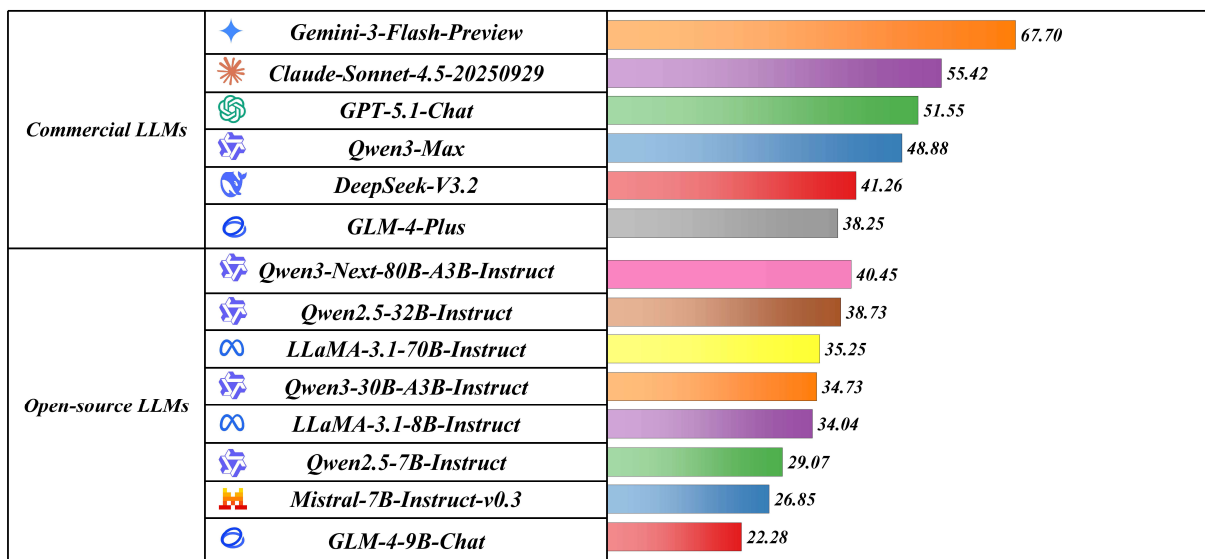


Figure 8: Overall Performance Comparison of 14 LLMs on CMiLBENCH

| Model                       | CR           | TC           | MRC          | NLI          | MR           | GDC          |
|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| <i>Tibetan</i>              |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | <b>79.28</b> | <b>86.57</b> | 54.53        | <b>64.80</b> | 61.67        | <b>81.87</b> |
| Claude-Sonnet-4.5           | <u>76.64</u> | <u>79.63</u> | 37.79        | <u>60.00</u> | <b>87.33</b> | <u>56.47</u> |
| Qwen3-Max                   | 67.76        | 75.69        | 53.41        | 51.80        | 20.33        | 40.67        |
| GPT-5.1-Chat                | 71.71        | 74.07        | 42.38        | 50.60        | <u>77.67</u> | 43.33        |
| Qwen3-Next-80B-A3B-Instruct | 67.76        | 62.50        | 38.16        | 46.20        | 11.00        | 34.13        |
| DeepSeek-V3.2               | 65.46        | 59.95        | 51.12        | 42.00        | 35.00        | 33.07        |
| Qwen2.5-32B-Instruct        | 63.49        | 37.73        | 40.67        | 45.20        | 2.67         | 27.87        |
| GLM-4-Plus                  | 69.41        | 72.45        | 40.29        | 40.80        | 13.67        | 32.73        |
| LLaMA-3.1-70B-Instruct      | 63.16        | 71.53        | <b>55.33</b> | 42.40        | 2.67         | 31.67        |
| Qwen3-30B-A3B-Instruct      | 48.68        | 45.60        | 47.73        | 44.60        | 3.67         | 28.80        |
| LLaMA-3.1-8B-Instruct       | 35.86        | 41.44        | 48.47        | 36.20        | 2.67         | 26.47        |
| Qwen2.5-7B-Instruct         | 63.49        | 25.23        | 26.21        | 42.80        | 1.33         | 27.60        |
| Mistral-7B-Instruct-v0.3    | 60.53        | 18.52        | 27.49        | 29.40        | 0.67         | 24.00        |
| GLM-4-9B-Chat               | 54.61        | 48.15        | 19.96        | 19.00        | 1.67         | 16.13        |
| <i>Mongolian</i>            |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | <b>74.01</b> | <b>84.72</b> | <b>31.90</b> | <b>57.60</b> | <b>68.33</b> | <b>71.00</b> |
| Claude-Sonnet-4.5           | 56.25        | 57.87        | 23.80        | <u>40.80</u> | <u>44.00</u> | <u>32.47</u> |
| Qwen3-Max                   | 62.83        | <u>59.49</u> | <u>28.61</u> | 40.60        | 9.67         | 31.00        |
| GPT-5.1-Chat                | <u>63.82</u> | 34.95        | 6.92         | 36.60        | 11.33        | 28.60        |
| Qwen3-Next-80B-A3B-Instruct | 62.50        | 45.83        | 14.41        | <u>40.80</u> | 2.67         | 29.33        |
| DeepSeek-V3.2               | 63.16        | 26.62        | 11.73        | 37.00        | 5.00         | 26.13        |
| Qwen2.5-32B-Instruct        | 63.49        | 29.40        | 16.87        | 40.60        | 1.00         | 26.93        |
| GLM-4-Plus                  | 61.84        | 23.84        | 14.11        | 37.00        | 3.33         | 26.73        |
| LLaMA-3.1-70B-Instruct      | 59.21        | 27.31        | 9.21         | 36.00        | 2.33         | 25.53        |
| Qwen3-30B-A3B-Instruct      | 51.64        | 16.44        | 13.30        | 35.80        | 1.33         | 27.67        |
| LLaMA-3.1-8B-Instruct       | 54.28        | 15.97        | 13.26        | 33.20        | 1.33         | 25.93        |
| Qwen2.5-7B-Instruct         | 62.50        | 17.36        | 4.05         | 37.20        | 1.00         | 29.40        |
| Mistral-7B-Instruct-v0.3    | 63.49        | 18.29        | 4.89         | 30.00        | 1.33         | 25.20        |
| GLM-4-9B-Chat               | 58.55        | 18.52        | 5.44         | 37.60        | 1.00         | 26.07        |
| <i>Uyghur</i>               |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | 74.34        | <b>86.34</b> | 34.25        | <b>63.00</b> | 65.33        | <b>84.73</b> |
| Claude-Sonnet-4.5           | <u>75.99</u> | 83.56        | 28.27        | <u>59.80</u> | <u>83.33</u> | 61.80        |
| Qwen3-Max                   | 70.72        | 79.40        | 37.83        | 54.20        | 30.33        | 58.13        |
| GPT-5.1-Chat                | <b>77.63</b> | <u>83.80</u> | 29.81        | 58.60        | <b>83.67</b> | <u>71.13</u> |
| Qwen3-Next-80B-A3B-Instruct | 64.47        | 82.18        | 28.28        | 49.20        | 16.67        | 49.27        |
| DeepSeek-V3.2               | 61.84        | 75.46        | 35.88        | 45.00        | 32.00        | 39.53        |
| Qwen2.5-32B-Instruct        | 63.49        | 65.05        | 32.81        | 46.20        | 16.00        | 35.67        |
| GLM-4-Plus                  | 61.84        | 78.70        | 38.41        | 41.20        | 12.67        | 39.80        |
| LLaMA-3.1-70B-Instruct      | 59.87        | 78.94        | <b>45.88</b> | 40.80        | 9.33         | 42.07        |
| Qwen3-30B-A3B-Instruct      | 48.03        | 73.38        | <u>38.95</u> | 44.00        | 3.00         | 39.00        |
| LLaMA-3.1-8B-Instruct       | 36.51        | 69.44        | 23.98        | 35.00        | 4.00         | 30.53        |
| Qwen2.5-7B-Instruct         | 64.47        | 59.49        | 30.62        | 42.00        | 6.67         | 31.00        |
| Mistral-7B-Instruct-v0.3    | 59.54        | 24.77        | 10.44        | 29.80        | 2.33         | 25.40        |
| GLM-4-9B-Chat               | 61.18        | 65.05        | 10.46        | 25.00        | 1.67         | 21.73        |

Table 7: Performance Comparison of 14 LLMs Across Six Foundation Tasks in Tibetan, Mongolian, and Uyghur. CR: Coreference Resolution; TC: Text Classification; MRC: Machine Reading Comprehension; NLI: Natural Language Inference; MR: Math Reasoning; GDC: General Domain Competence. **Bold underline** indicates first place, underline indicates second place.

| Model                       | MCQA         | MLIQA        | MLE          | MLU          | MDC          | MMT          |
|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| <i>Tibetan</i>              |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | <u>42.87</u> | <b>78.66</b> | <b>89.43</b> | <b>83.00</b> | <b>53.13</b> | 53.94        |
| Claude-Sonnet-4.5           | <b>48.20</b> | <u>59.70</u> | <u>74.53</u> | <u>56.00</u> | <u>40.87</u> | 35.70        |
| Qwen3-Max                   | 14.84        | 26.72        | 77.92        | 43.00        | 33.16        | 34.92        |
| GPT-5.1-Chat                | 17.55        | 19.16        | 76.04        | 48.00        | 37.87        | 19.79        |
| Qwen3-Next-80B-A3B-Instruct | 8.55         | 11.20        | 70.94        | 36.00        | 32.60        | 14.00        |
| DeepSeek-V3.2               | 29.61        | 32.88        | 66.42        | 26.00        | 28.93        | <b>56.27</b> |
| Qwen2.5-32B-Instruct        | 4.90         | 6.13         | 54.91        | 27.00        | 31.40        | 11.08        |
| GLM-4-Plus                  | 32.36        | 26.79        | 73.96        | 37.00        | 34.40        | 11.87        |
| LLaMA-3.1-70B-Instruct      | 11.61        | 11.66        | 67.55        | 31.00        | 27.53        | 25.98        |
| Qwen3-30B-A3B-Instruct      | 7.73         | 8.76         | 62.26        | 27.00        | 32.47        | 16.69        |
| LLaMA-3.1-8B-Instruct       | 7.06         | 10.20        | 52.83        | 26.00        | 28.73        | 4.79         |
| Qwen2.5-7B-Instruct         | 2.44         | 5.39         | 36.60        | 30.00        | 24.53        | 6.48         |
| Mistral-7B-Instruct-v0.3    | 12.65        | 12.95        | 32.83        | 26.00        | 27.53        | 3.56         |
| GLM-4-9B-Chat               | 2.29         | 4.21         | 47.55        | 18.00        | 8.73         | 6.65         |
| <i>Mongolian</i>            |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | <b>14.93</b> | <b>52.01</b> | <b>96.98</b> | <b>34.00</b> | <b>56.73</b> | <b>34.63</b> |
| Claude-Sonnet-4.5           | <u>15.84</u> | <u>23.38</u> | <u>77.36</u> | 22.00        | <u>33.33</u> | <u>21.16</u> |
| Qwen3-Max                   | 6.35         | 8.78         | 84.91        | 24.00        | 30.80        | 14.25        |
| GPT-5.1-Chat                | 1.65         | 0.21         | 44.72        | 21.00        | 24.07        | 4.52         |
| Qwen3-Next-80B-A3B-Instruct | 1.34         | 3.08         | 56.04        | 21.00        | 29.13        | 10.33        |
| DeepSeek-V3.2               | 1.03         | 4.23         | 26.42        | 29.00        | 21.40        | 10.89        |
| Qwen2.5-32B-Instruct        | 2.64         | 4.85         | 36.79        | 25.00        | 29.80        | 1.88         |
| GLM-4-Plus                  | 0.39         | 2.00         | 30.00        | <u>30.00</u> | 25.47        | 5.06         |
| LLaMA-3.1-70B-Instruct      | 4.72         | 3.57         | 31.32        | 28.00        | 21.60        | 7.21         |
| Qwen3-30B-A3B-Instruct      | 2.00         | 1.62         | 30.57        | 26.00        | 26.47        | 1.43         |
| LLaMA-3.1-8B-Instruct       | 1.64         | 8.75         | 29.06        | 27.00        | 21.93        | 2.08         |
| Qwen2.5-7B-Instruct         | 0.31         | 2.40         | 26.60        | 25.00        | 21.80        | 0.13         |
| Mistral-7B-Instruct-v0.3    | 4.40         | 3.68         | 25.28        | 24.00        | 31.13        | 10.68        |
| GLM-4-9B-Chat               | 1.78         | 3.64         | 20.94        | 28.00        | 12.20        | 0.56         |
| <i>Uyghur</i>               |              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | 49.45        | <b>93.75</b> | <b>97.17</b> | 49.00        | <b>61.73</b> | <b>45.61</b> |
| Claude-Sonnet-4.5           | <b>54.73</b> | <u>74.46</u> | <u>89.81</u> | 47.00        | 50.40        | 28.71        |
| Qwen3-Max                   | 23.62        | 47.26        | 92.64        | <u>51.00</u> | 47.73        | 22.82        |
| GPT-5.1-Chat                | 49.14        | 68.37        | 92.83        | <b>56.00</b> | <u>52.27</u> | 30.51        |
| Qwen3-Next-80B-A3B-Instruct | 11.77        | 20.63        | 90.57        | 43.00        | 45.00        | 12.24        |
| DeepSeek-V3.2               | 37.81        | 62.79        | 87.36        | 34.00        | 34.13        | <u>39.01</u> |
| Qwen2.5-32B-Instruct        | 2.59         | 4.89         | 84.34        | 41.00        | 37.20        | 10.35        |
| GLM-4-Plus                  | 34.32        | 43.50        | 91.13        | 38.00        | 39.67        | 29.10        |
| LLaMA-3.1-70B-Instruct      | 17.38        | 12.72        | 88.87        | 31.00        | 32.53        | 13.78        |
| Qwen3-30B-A3B-Instruct      | 5.93         | 9.19         | 87.36        | 40.00        | 41.27        | 17.18        |
| LLaMA-3.1-8B-Instruct       | 13.15        | 8.62         | 78.11        | 24.00        | 32.40        | 3.17         |
| Qwen2.5-7B-Instruct         | 3.53         | 4.17         | 71.89        | 28.00        | 28.67        | 5.83         |
| Mistral-7B-Instruct-v0.3    | 3.85         | 10.36        | 37.92        | 30.00        | 23.40        | 3.29         |
| GLM-4-9B-Chat               | 0.55         | 2.70         | 74.34        | 32.00        | 11.33        | 5.67         |

Table 8: Performance Comparison of 14 LLMs Across Six Chinese Minority Knowledge Tasks in Tibetan, Mongolian, and Uyghur. MCQA: Minority Culture QA; MLIQA: Minority Language Instruction QA; MLE: Minority Language Expressions; MLU: Minority Language Understanding; MDC: Minority Domain Competence; MMT: Minority Machine Translation. **Bold underline** indicates first place, underline indicates second place.

| Model                       | CCC          | DD           | RPE          | SSE          | VAA          |
|-----------------------------|--------------|--------------|--------------|--------------|--------------|
| <i>Tibetan</i>              |              |              |              |              |              |
| Gemini-3-Flash-Preview      | 86.61        | <u>68.13</u> | <u>82.95</u> | 43.30        | <b>92.44</b> |
| Claude-Sonnet-4.5           | 67.72        | <b>69.01</b> | 46.31        | <u>51.55</u> | <u>83.12</u> |
| Qwen3-Max                   | 88.98        | 59.56        | 53.41        | 46.39        | 77.08        |
| GPT-5.1-Chat                | 64.57        | 56.92        | 45.17        | 47.42        | 81.36        |
| Qwen3-Next-80B-A3B-Instruct | 53.54        | 35.38        | 23.58        | 38.14        | 64.23        |
| DeepSeek-V3.2               | 51.97        | 41.10        | 33.24        | 39.18        | 65.49        |
| Qwen2.5-32B-Instruct        | <b>94.49</b> | 64.62        | <b>71.59</b> | <b>84.54</b> | 51.13        |
| GLM-4-Plus                  | 47.24        | 35.60        | 28.69        | 39.18        | 62.72        |
| LLaMA-3.1-70B-Instruct      | 57.09        | 44.40        | 36.93        | 40.21        | 74.31        |
| Qwen3-30B-A3B-Instruct      | 43.31        | 39.56        | 28.12        | 24.74        | 41.81        |
| LLaMA-3.1-8B-Instruct       | <u>90.94</u> | 69.45        | 66.19        | 57.73        | 47.61        |
| Qwen2.5-7B-Instruct         | 38.58        | 37.80        | 16.19        | 31.96        | 43.83        |
| Mistral-7B-Instruct-v0.3    | 38.19        | 28.13        | 12.22        | 29.90        | 44.08        |
| GLM-4-9B-Chat               | 16.54        | 40.00        | 9.66         | 14.43        | 24.94        |
| <i>Mongolian</i>            |              |              |              |              |              |
| Gemini-3-Flash-Preview      | 86.61        | 56.70        | <u>79.26</u> | 47.42        | <b>85.14</b> |
| Claude-Sonnet-4.5           | 64.57        | 49.45        | 39.77        | <u>57.73</u> | <u>49.37</u> |
| Qwen3-Max                   | 67.32        | 54.73        | 62.78        | <u>57.73</u> | 47.61        |
| GPT-5.1-Chat                | 82.68        | <u>64.62</u> | 65.06        | 62.89        | 44.84        |
| Qwen3-Next-80B-A3B-Instruct | 63.78        | 53.19        | 36.08        | 54.64        | 45.09        |
| DeepSeek-V3.2               | 54.72        | 52.31        | 47.73        | 43.30        | 40.55        |
| Qwen2.5-32B-Instruct        | <b>91.34</b> | <b>68.57</b> | <b>72.16</b> | <b>67.01</b> | 48.36        |
| GLM-4-Plus                  | 30.71        | 38.90        | 26.42        | 46.39        | 35.01        |
| LLaMA-3.1-70B-Instruct      | 32.68        | 35.38        | 30.68        | 38.14        | 30.23        |
| Qwen3-30B-A3B-Instruct      | 58.27        | 55.60        | 49.72        | <b>67.01</b> | 39.55        |
| LLaMA-3.1-8B-Instruct       | <u>74.80</u> | 46.15        | 46.31        | 56.70        | 35.52        |
| Qwen2.5-7B-Instruct         | 7.09         | 25.93        | 12.22        | 20.62        | 27.96        |
| Mistral-7B-Instruct-v0.3    | 5.12         | 23.30        | 9.09         | 31.96        | 24.69        |
| GLM-4-9B-Chat               | 48.43        | 47.47        | 51.42        | 31.96        | 29.72        |
| <i>Uyghur</i>               |              |              |              |              |              |
| Gemini-3-Flash-Preview      | 87.40        | <u>82.42</u> | <u>73.58</u> | <u>83.51</u> | <b>86.15</b> |
| Claude-Sonnet-4.5           | 77.17        | 76.48        | 57.67        | 68.04        | <u>69.52</u> |
| Qwen3-Max                   | 70.87        | 71.21        | 52.56        | 61.86        | 63.48        |
| GPT-5.1-Chat                | <u>90.55</u> | <b>84.62</b> | <b>76.70</b> | <b>79.38</b> | 71.79        |
| Qwen3-Next-80B-A3B-Instruct | 72.05        | 72.53        | 57.67        | 64.95        | 65.49        |
| DeepSeek-V3.2               | 61.02        | 60.66        | 48.58        | 50.52        | 39.80        |
| Qwen2.5-32B-Instruct        | 35.04        | 48.35        | 39.49        | 34.02        | 32.49        |
| GLM-4-Plus                  | 80.31        | 56.70        | 42.61        | 32.99        | 53.40        |
| LLaMA-3.1-70B-Instruct      | 38.98        | 55.82        | 42.90        | 26.80        | 42.07        |
| Qwen3-30B-A3B-Instruct      | 76.77        | 53.63        | 43.18        | 44.33        | 50.13        |
| LLaMA-3.1-8B-Instruct       | 57.87        | 53.63        | 37.22        | 43.30        | 38.54        |
| Qwen2.5-7B-Instruct         | 74.02        | 55.38        | 54.83        | 76.29        | 63.73        |
| Mistral-7B-Instruct-v0.3    | <b>91.34</b> | 51.43        | 67.05        | 60.82        | 72.29        |
| GLM-4-9B-Chat               | 14.96        | 14.07        | 6.82         | 0.00         | 13.10        |

Table 9: Performance Comparison of 14 LLMs Across Five Safety Alignment Tasks in Tibetan, Mongolian, and Uyghur. CCC: Commercial Compliance Check; DD: Discrimination Detection; RPE: Rights Protection Evaluation; SSE: Service Safety Evaluation; VAA: Value Alignment Assessment. **Bold underline** indicates first place, underline indicates second place.

Table 10: Summary of Evaluated LLMs

| Model (#params)                   | Version                           |
|-----------------------------------|-----------------------------------|
| <i>Commercial Models</i>          |                                   |
| GPT-5.1-Chat                      | gpt-5.1-chat                      |
| Claude-Sonnet-4.5                 | claude-sonnet-4-5-20250929        |
| Gemini-3-Flash-Preview            | google_gemini-3-flash-preview     |
| Qwen3-Max                         | qwen_qwen3-max                    |
| GLM-4-Plus                        | glm-4-plus                        |
| DeepSeek-V3.2 (685B)              | deepseek_deepseek-v3.2            |
| <i>Open-Source Models</i>         |                                   |
| Mistral-7B-Instruct-v0.3 (7B)     | Mistral-7B-Instruct-v0.3          |
| Qwen2.5-7B-Instruct (7B)          | Qwen2.5-7B-Instruct               |
| LLaMA-3.1-8B-Instruct (8B)        | Meta-LLaMA-3.1-8B-Instruct        |
| GLM-4-9B-Chat (9B)                | glm-4-9b-chat                     |
| Qwen3-30B-A3B-Instruct (30B)      | qwen_qwen3-30b-a3b-instruct-2507  |
| Qwen2.5-32B-Instruct (32B)        | Qwen2.5-32B-Instruct              |
| LLaMA-3.1-70B-Instruct (70B)      | meta-llama_llama-3.1-70b-instruct |
| Qwen3-Next-80B-A3B-Instruct (80B) | qwen_qwen3-next-80b-a3b-instruct  |

Table 11: Task Details Per Language in CMiLBENCH

| Task Category                             | Task Name                        | Per Language         | Task Type         | Evaluation Metric |
|---|----------------------------------|----------------------|-------------------|-------------------|
| Foundation<br>Tasks                       | Natural Language Inference       | 500                  | Multiple Choice   | Accuracy          |
|   | Coreference Resolution           | 304                  | Multiple Choice   | Accuracy          |
|   | Machine Reading Comprehension    | 500                  | Generation        | ROUGE-L           |
|   | Text Classification              | 432                  | Fill in the Blank | Accuracy          |
|   | Math Reasoning                   | 300                  | Fill in the Blank | Accuracy          |
|   | General Domain Competence        | 1,500                | Multiple Choice   | Accuracy          |
| Chinese<br>Minority<br>Knowledge<br>Tasks | Minority Language Expressions    | 530                  | Multiple Choice   | Accuracy          |
|   | Minority Machine Translation     | 500                  | Generation        | BLEU & chrF++     |
|   | Minority Culture QA              | 282                  | Generation        | LLM-as-a-Judge    |
|   | Minority Language Instruction QA | 218                  | Generation        | LLM-as-a-Judge    |
|   | Minority Language Understanding  | 100                  | Multiple Choice   | Accuracy          |
| Alignment<br>Tasks                        | Minority Domain Competence       | 1,500                | Multiple Choice   | Accuracy          |
|   | Discrimination Detection         | 455                  | Multiple Choice   | Accuracy          |
|   | Value Alignment Assessment       | 397                  | Multiple Choice   | Accuracy          |
|   | Rights Protection Evaluation     | 352                  | Multiple Choice   | Accuracy          |
|   | Commercial Compliance Check      | 254                  | Multiple Choice   | Accuracy          |
|   | Service Safety Evaluation        | 97                   | Multiple Choice   | Accuracy          |
| <b>Total</b>                              | <b>17 tasks</b>                  | <b>8,221 samples</b> | <b>3 types</b>    | <b>5 metrics</b>  |

Table 12: Multi-Dimensional Scoring Framework for LLM-as-a-Judge Evaluation

| Task Name                               | Response Type     | Scoring Dimensions   |
|---|-------------------|--|
| <b>Minority Culture QA</b>              | –                 | Factual Accuracy, Cultural Understanding Depth, Appropriateness of Language Use, Content Completeness, Authenticity of Insider Perspective |
| <b>Minority Language Instruction QA</b> | Factual Answer    | Factual Accuracy, User Need Satisfaction, Clarity, Completeness  |
|   | Reasoning Answer  | Factual Accuracy, User Need Satisfaction, Logical Coherence, Completeness  |
|   | Generation Answer | Factual Accuracy, User Need Satisfaction, Logical Coherence, Creativity, Richness  |
|   | Advisory Answer   | Factual Accuracy, User Need Satisfaction, Fairness & Responsibility, Creativity  |

Table 13: Comparison between *Minority Culture QA* and *Minority Domain Competence*

| Feature               | Minority Culture QA                       | Minority Domain Competence                         |
|-----------------------|---|--|
| <b>Content Type</b>   | Sociological, Anthropological, & Cultural | Technical, Scientific, & Academic History          |
| <b>Format</b>         | Open-ended Text (Descriptive)             | Multiple Choice (Objective)                        |
| <b>Example Topics</b> | Marriage Customs, Farming Tools           | Medical Prescriptions, Specific Historical Regimes |
| <b>Primary Goal</b>   | Explain Cultural Context and Traditions   | Test Specific Professional Competency              |

Table 14: Comparison between *Minority Language Understanding* and *Minority Language Expressions*

| Feature                   | Minority Language Understanding                 | Minority Language Expressions                 |
|---------------------------|---|---|
| <b>Primary Task</b>       | Grammar, Logic, & Reading Comprehension         | Translation & Vocabulary Mapping              |
| <b>Language Direction</b> | Monolingual (Tibetan Question → Tibetan Answer) | Bilingual (Chinese Question → Tibetan Answer) |
| <b>Key Skills</b>         | Syntax, Particle Usage, Logical Reasoning       | Idiomatic Translation, Lexical Knowledge      |

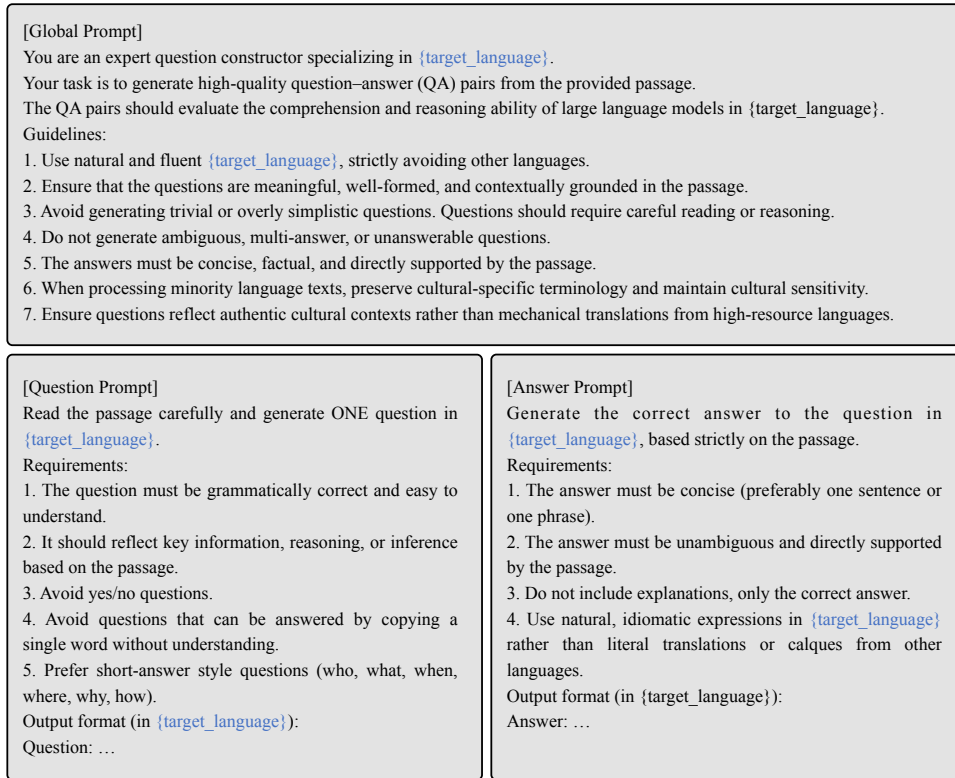


Figure 9: Structured prompt templates used in the Easy Dataset framework (Miao et al., 2025) for Q&A generation in Chinese minority languages.

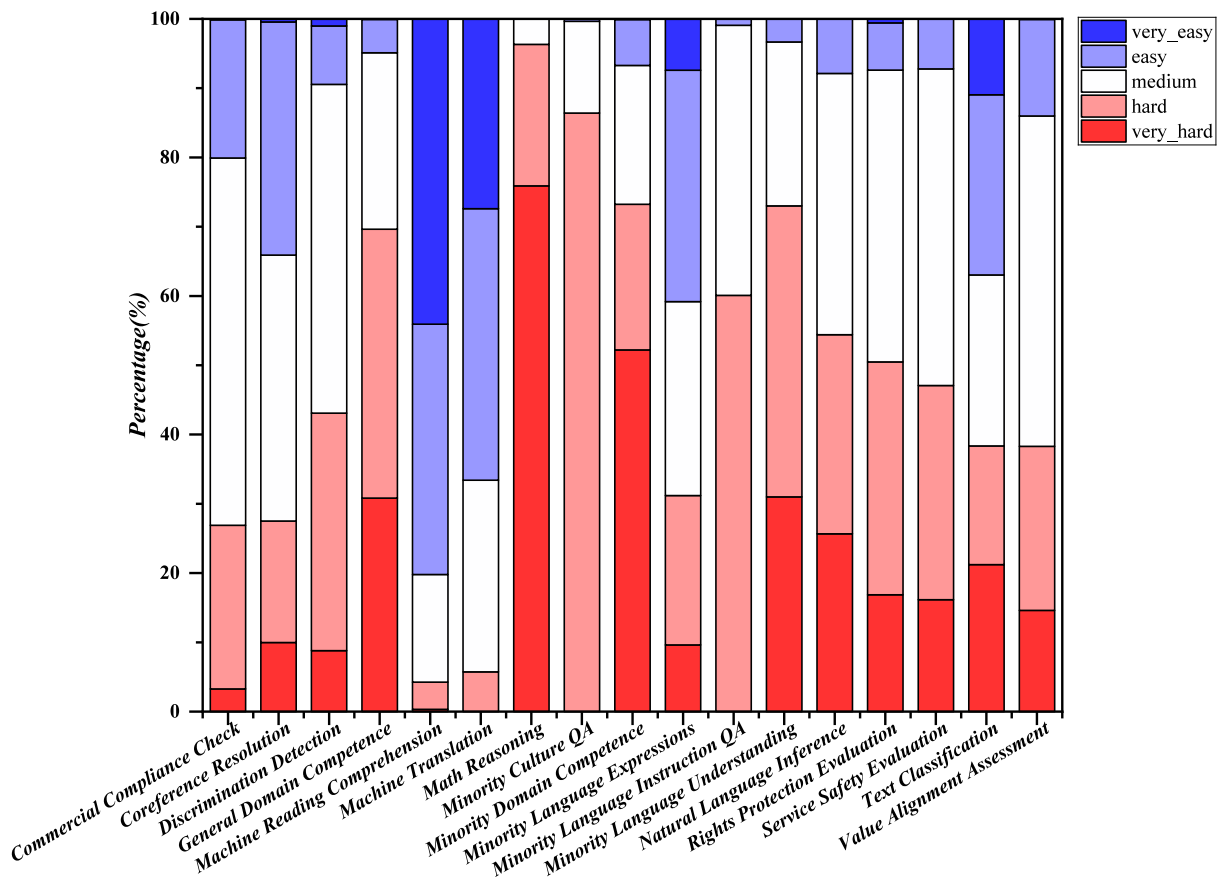


Figure 10: Overall distribution across five difficulty levels in CMILBENCH.

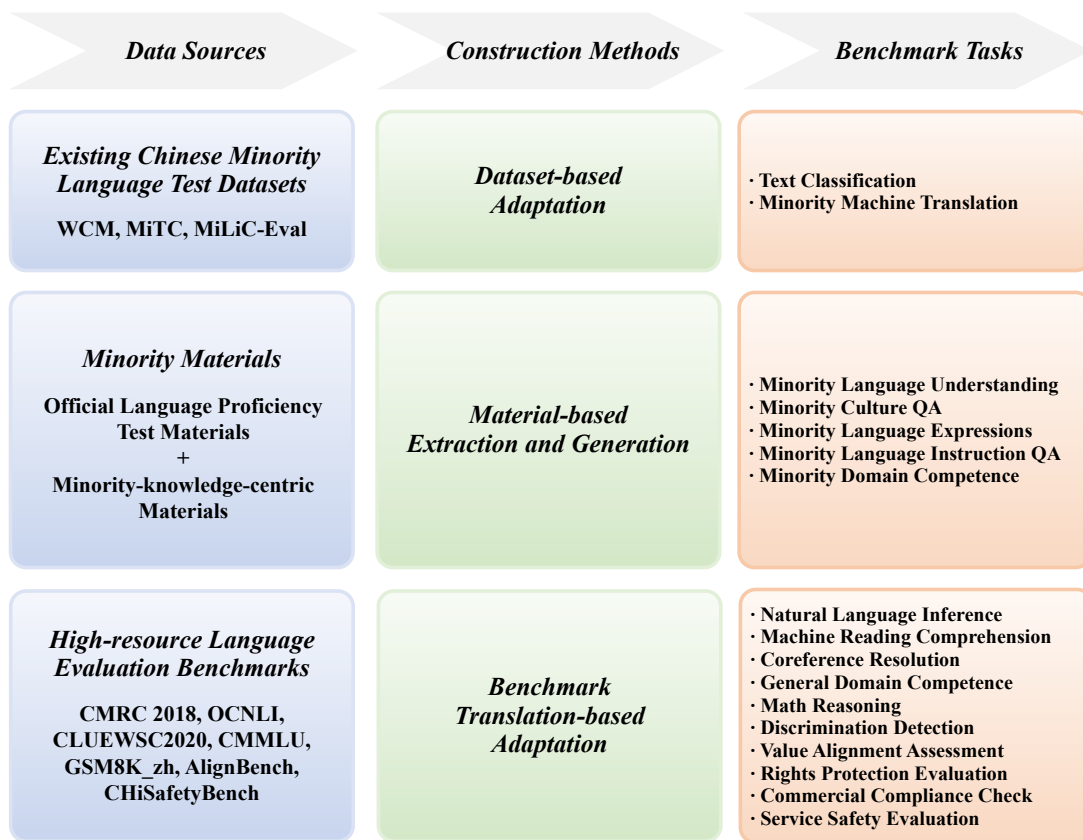


Figure 11: Mapping between data sources, construction methods, and benchmark tasks within CMiLBENCH.

