

---

# TCM-Ladder: A Benchmark for Multimodal Question Answering on Traditional Chinese Medicine

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Traditional Chinese Medicine (TCM), as an effective alternative medicine, has  
2 been receiving increasing attention. In recent years, the rapid development of  
3 large language models (LLMs) tailored for TCM has underscored the need for an  
4 objective and comprehensive evaluation framework to assess their performance  
5 on real-world tasks. However, existing evaluation datasets are limited in scope  
6 and primarily text-based, lacking a unified and standardized multimodal question-  
7 answering (QA) benchmark. To address this issue, we introduce **TCM-Ladder**,  
8 the first multimodal QA dataset specifically designed for evaluating large TCM  
9 language models. The dataset spans multiple core disciplines of TCM, includ-  
10 ing fundamental theory, diagnostics, herbal formulas, internal medicine, surgery,  
11 pharmacognosy, and pediatrics. In addition to textual content, TCM-Ladder in-  
12 corporates various modalities such as images and videos. The datasets were  
13 constructed using a combination of automated and manual filtering processes  
14 and comprise **52,000+** questions in total. These questions include single-choice,  
15 multiple-choice, fill-in-the-blank, diagnostic dialogue, and visual comprehension  
16 tasks. We trained a reasoning model on TCM-Ladder and conducted comparative  
17 experiments against 9 state-of-the-art general domain and 5 leading TCM-specific  
18 LLMs to evaluate their performance on the datasets. Moreover, we propose **Ladder-**  
19 **Score**, an evaluation method specifically designed for TCM question answering  
20 that effectively assesses answer quality regarding terminology usage and semantic  
21 expression. To our knowledge, this is the first work to evaluate mainstream gen-  
22 eral domain and TCM-specific LLMs on a unified multimodal benchmark. The  
23 datasets and leaderboard are publicly available at <https://tcmladder.com> or  
24 <https://54.211.107.106> and will be continuously updated. The source code  
25 is available at <https://github.com/orangeshushu/TCM-Ladder>.

26 

## 1 Introduction

27 The development of large language models tailored to the field of Traditional Chinese Medicine  
28 (TCM) [1, 2] has emerged as a significant research direction. Given the unique and intricate nature  
29 of the TCM knowledge system, the construction of intelligent tools specifically designed for this  
30 domain can substantially enhance the efficiency of medical students, clinicians, and researchers.  
31 Such models have the potential to facilitate accurate and timely access to specialized information for  
32 clinical decision-making, knowledge retrieval, and academic inquiry, thereby supporting effective  
33 reasoning and practical application within the TCM framework.

34 TCM diagnostic methods including inspection, auscultation and olfaction, inquiry, and palpation  
35 embody a representative process of multimodal information acquisition, integration, and reasoning  
36 [3]. Fundamentally, this diagnostic paradigm reflects the nature of multimodal fusion in clinical

37 decision-making. However, existing large language models (LLMs) tailored for TCM still face  
38 notable limitations in real-world applications. These limitations are primarily manifested in their  
39 relatively small model scales, insufficient reasoning capacity, and the lack of deep integration of  
40 multimodal information. The acquisition of high-quality TCM data poses significant challenges, as  
41 it requires deep expertise in traditional medicine, sustained clinical data collection, and extensive  
42 manual annotation. Currently, most mainstream medical benchmark datasets [4, 5, 6, 7, 8] are  
43 predominantly focused on Western medicine and have yet to systematically address the core tasks  
44 unique to TCM, including syndrome differentiation, symptom-based diagnosis, and formula-herb  
45 matching. Furthermore, the training and evaluation of existing TCM large language models remain  
46 heavily reliant on unimodal textual data, neglecting other essential modalities that are widely utilized  
47 in clinical practice. These include diagnostic images (e.g., tongue and pulse), medicinal herb atlases,  
48 and structured case records. Such an overdependence on textual data severely constrains the models'  
49 ability to capture the holistic and multimodal nature of TCM knowledge, thereby impeding their  
50 performance in complex and realistic clinical scenarios.

51 Therefore, the construction of a standardized evaluation dataset for TCM that integrates text, images,  
52 audio, and structured data is of great importance. On one hand, such a dataset would enable a  
53 comprehensive and accurate assessment of existing LLMs in handling complex multimodal tasks,  
54 thereby providing a realistic reflection of their overall performance in clinical applications. On  
55 the other hand, a unified and standardized evaluation framework would facilitate fair and objective  
56 comparisons across different TCM-specific models, supporting continuous optimization and iterative  
57 improvement of model capabilities.

58 To address the aforementioned gaps, we propose **TCM-Ladder**, which, to the best of our knowledge,  
59 is the first large-scale multimodal dataset specifically designed for the training and evaluation of large  
60 language models in TCM. TCM-Ladder encompasses a wide spectrum of domain-specific knowledge,  
61 including fundamental TCM theories, diagnostics, formulae, pharmacology, clinical medicine, as  
62 well as visual modalities such as tongue images, herbal medicine illustrations, acupuncture, and tuina  
63 (therapeutic massage), thereby offering a comprehensive foundation for developing and benchmarking  
64 TCM-specific LLMs.

65 As illustrated in Figure 1, based on the TCM-Ladder dataset, we design a series of evaluation tasks  
66 to comprehensively assess the capabilities of TCM-specific large language models across multiple  
67 dimensions. We constructed a total of **21,326 high-quality questions** and **25,163 diagnostic long-text**  
68 **dialogues** based on domain-specific literature and publicly available databases across various subfields  
69 of TCM. In addition, we release a **visual dataset** comprising **6,061 images of medicinal herbs**,  
70 **1,394 tongue images**, **6,420 audio clips**, and **49 videos**, forming a comprehensive multimodal  
71 foundation to support diverse evaluation tasks. All textual and visual data were independently  
72 reviewed and validated by certified TCM practitioners to ensure accuracy, clinical relevance, and  
73 authoritative quality. Subsequently, we benchmarked the performance of 9 state-of-the-art general  
74 domain LLMs[9, 10, 11, 12, 13, 14, 15, 16, 17] and 5 TCM-specific models[18, 19, 20] using the  
75 TCM-Ladder dataset. Additionally, we fine-tuned a GPT-4-based model, Bencao[21], and trained a  
76 Qwen2.5-7B [22] based reasoning model, which uses a training subset constructed from TCM-Ladder  
77 to support TCM-specific reasoning tasks.

78 Our contributions can be summarized as follows:

- 79 • We construct **TCM-Ladder**, a multimodal dataset designed for both training and evaluating  
80 TCM-specific and general domain LLMs. The dataset encompasses multiple TCM sub-  
81 disciplines and a variety of data modalities.
- 82 • We design a comprehensive set of tasks including single-choice questions, multiple-choice  
83 questions, fill-in-the-blank tasks, visual understanding tasks, and long-form question an-  
84 swering to evaluate models' reasoning abilities across different tasks.
- 85 • We introduce **Ladder-Score**, an evaluation metric that integrates TCM-specific terminology  
86 and LLM-assisted semantic scoring to assess term accuracy and reasoning quality in TCM  
87 question answering.
- 88 • We systematically evaluate the performance of several general domain and TCM-specific  
89 LLMs on TCM-Ladder. To our knowledge, this is the first work to conduct a comparative  
90 evaluation of diverse LLMs on a unified multimodal TCM dataset.



Figure 1: Overview of the architectural composition of TCM-Ladder. TCM-Ladder encompasses six task types aimed at evaluating the comprehensive capabilities of large language models in Traditional Chinese Medicine. These include: (1) single-choice questions, which assess basic knowledge recognition; (2) multiple-choice questions, designed to test the model’s ability to integrate and reason over complex concepts; (3) long-form diagnostic question answering, which evaluates clinical reasoning based on detailed symptom descriptions and patient inquiries; (4) fill-in-the-blank tasks, which measure generative accuracy and contextual understanding without the aid of answer options; (5) image-based comprehension tasks, involving the interpretation of medicinal herb and tongue images to assess multimodal reasoning across visual and textual inputs; and (6) additional audio and video resources, such as diagnostic sounds, pulse recordings, and tuina (massage) videos, which support the development and evaluation of multimodal TCM models incorporating auditory and dynamic visual data.

- 91 • We develop an interactive data visualization website that not only presents evaluation results,  
 92 but also allows researchers to explore existing data and contribute new entries, thereby  
 93 providing a standardized, extensible, and multimodal infrastructure for future benchmarking  
 94 of TCM-specific LLMs.

## 95 2 Related Works

96 In recent years, the expanding application of LLMs in medicine and the sciences has driven the  
 97 progressive development of evaluation datasets tailored for TCM, evolving from modern medical  
 98 domains to TCM-specific tasks, and from classification-based to generation-based paradigms.  
 99 **Huatuo-26M**[23], released in 2020, remains the largest Chinese medical QA dataset, comprising  
 100 over 26 million question-answer pairs sourced from online encyclopedias, medical knowledge bases,  
 101 and telemedicine transcripts. Despite its scale, the dataset suffers from noisy labels, informal expressions,  
 102 redundancy, and a lack of TCM-specific annotations, limiting its utility for TCM applications.  
 103 **CBLUE**[24] introduced a standardized multi-task evaluation suite for Chinese biomedical NLP,  
 104 covering named entity recognition, relation extraction, etc. **PromptCBLUE**[25] extended this frame-  
 105 work via instruction tuning and prompt reformulation to facilitate few-shot and zero-shot evaluation.  
 106 However, both benchmarks were designed around modern medical reasoning and do not reflect the  
 107 unique logic or semantic structure of TCM diagnosis.

108 To address these gaps, **TCMBench**[26] compiled 5,473 structured questions from national TCM li-  
 109 censing examinations, providing a focused benchmark for foundational knowledge assessment. Never-  
 110 theless, it lacks multimodal input (e.g., tongue and pulse images) and real-world diagnostic reasoning  
 111 tasks. **TCMEval-SDT**[27] introduced syndrome differentiation based on 300 clinical cases, eval-

ating the model’s reasoning over symptom–pathomechanism–syndrome chains. While it improved interpretability, its scale and disease diversity remain limited. Subsequently, **TCM-3CEval**[28] proposed a cognitive three-axis framework—basic knowledge, classical text comprehension, and clinical decision-making—enabling fine-grained cognitive evaluation. However, tasks were still text-only and often reduced the complexity of classical TCM literature to overly simplistic answers. **TCMD**[29] presented a human-annotated open-ended QA benchmark emphasizing reasoning and generation, though annotation costs limited its scale and case diversity. **ShenNong\_TCM\_Dataset**[30] adopted a novel approach, combining knowledge graphs with ChatGPT-based generation to create 110,000+ instruction–response pairs on herbal medicine and treatment plans. While valuable for instruction tuning, the absence of expert validation raises concerns over factual accuracy and stylistic fidelity. **CHBench**[31] introduced a safety-focused benchmark with 9,492 community-sourced questions, highlighting deficiencies in LLM reliability under ethically sensitive conditions. However, its scope remains narrow. **MedBench**[32] represents the most comprehensive Chinese medical LLM evaluation to date, integrating 20 datasets and over 300,000 questions across diverse tasks, including QA, clinical case analysis, diagnostic reasoning, and summarization. The platform supports dynamic sampling and randomized option ordering to prevent overfitting. However, access to API use is restricted due to data privacy concerns. Benchmarks like **CMB**[33] and **CMExam**[34] further extend to structured exam QA, offering high coverage but lacking realistic patient–physician interaction.

TCM-Ladder distinguishes itself from existing datasets in several key aspects. First, it establishes a large-scale, open-ended QA dataset that spans a wide range of TCM subfields, including basic theory, diagnostics, internal medicine, surgery, pediatrics, and pharmacology. This breadth enables more thorough and representative evaluation of TCM-specific LLMs across multiple knowledge domains. Second, TCM-Ladder integrates visual elements such as herbal medicine images and tongue diagnostics. This multimodal design reflects traditional TCM diagnostic practices, requiring LLMs to demonstrate both textual reasoning and visual understanding capabilities. Third, TCM-Ladder incorporates a variety of task formats. This comprehensive task structure facilitates an in-depth evaluation of the strengths and limitations of LLMs, providing guidance for the future development of TCM-specific models.

Table 1: Overview of TCM and Medical QA Datasets. En:English, Zh: Chinese.

Dataset	Format	TCM Coverage	Size	Source	Domain	Task	Verified	Language
<b>Huatuo-26M</b>	Text	✗	26M+	Online QA platforms and physician records	Medicine	QA, Dialogue	✗	Zh
<b>CBLUE</b>	Text	✗	13 subtasks	Clinical trials, EHRs, logs, textbooks	Biomedical	Classification, NER, RE, NLI	Partial	Zh
<b>PromptCBLUE</b>	Text	✗	11 prompt datasets	Prompt-formatted CBLUE	Biomedical	Same as CBLUE	✗	Zh
<b>TCMD</b>	Text	✓	1,500+	Professional TCM practitioners	TCM	NER, Term Normalization	✓	Zh
<b>TCM-3CEval</b>	Text	✓	4,000+	Expert-annotated multi-rater QA	TCM	QA	✓	Zh
<b>ShenNong_TCM_Dataset</b>	Text	✓	113,000	TCM knowledge graph GPT-3.5 assisted	TCM	Dialogue	✗	Zh
<b>CMB</b>	Text	Partial	280,839 MCQ, 74 consults	Textbooks, forums, exams	TCM	MCQ, Dialogue	✓	Zh
<b>CMExam</b>	Text	Partial	60,000+	TCM licensing exam	Medicine	MCQ, QA	Partial	Zh
<b>CHBench</b>	Text	Partial	9,492	Community health Q&A	Health	QA	✓	Zh
<b>MedBench</b>	Text	Partial	40,041	Clinical exam questions	Medicine	MCQ, QA	Partial	Zh
<b>TCMBench</b>	Text	✓	5,473	TCM licensing exam	TCM	QA	✗	Zh
<b>TCM-Ladder (Ours)</b>	Text, images, audio, video	✓	52,000+	Research, books, QA platforms	TCM	MCQ, FIB, QA, Dialogue, Image Understanding	✓	Zh & En

### 3 TCM-Ladder Dataset

#### 3.1 Data Collection

- We collected a question-answering dataset covering various domains of TCM, including several publicly available datasets previously published in academic literature under permissive licenses. For the textual data, we identified seven subfields: fundamental theory, diagnostics, herbal formulas, internal medicine, surgery, pharmacognosy, and pediatrics.
- For the Chinese herbal medicine image data, we collected over 6,061 images of medicinal herbs based on the herb names referenced in *The Pharmacology of Chinese Herbs*[35]. The dataset comprises images sourced from publicly available online resources, as well as photographs we captured at

149 traditional Chinese medicine manufacturing facilities. Sample images and collection process are  
150 provided in **Appendix G**.

151 The clinical tongue images were collected by a tongue imaging device [14] at Shanghai University  
152 of Traditional Chinese Medicine. This device is designed for tongue diagnosis and provides stable  
153 and consistent lighting conditions during image acquisition. Another subset of the proprietary data  
154 was obtained from our previous work, the *iTongue* [36, 37] diagnostic software. All data collection  
155 procedures were approved by the institutional ethics review board. To protect the privacy of tongue  
156 image contributors, only a subset of tongue image patches and corresponding labels have been  
157 released.

158 The video data was recorded by faculty members from the Department of Acupuncture and Tuina  
159 at Shanghai University of Traditional Chinese Medicine. These instructional videos cover essential  
160 techniques, procedural explanations, and key operational steps. Audio and pulse diagnosis data were  
161 sourced from publicly available datasets referenced in academic publications[38, 39, 40, 41]. A  
162 detailed list is available in the supplementary materials. We manually filtered and removed samples  
163 with poor quality or missing information from the collected data.

### 164 3.2 Construction of the Dataset

165 The textual question-answering (QA) data consist of two parts. The first part comprises 5,000  
166 TCM-related QA pairs manually written by licensed TCM practitioners under a standardized question  
167 design protocol (see **Appendix I**). To ensure answer accuracy, each question was independently  
168 reviewed and verified by two additional TCM physicians. The second part of the textual QA data was  
169 collected from publicly available sources, including the National Physician Qualification Examination  
170 of China and various open-access online resources. Detailed data sources and construction guidelines  
171 are provided in the **Appendix B**.

172 The visual question-answering (VQA) tasks were constructed through both manual annotation  
173 and automated generation based on existing knowledge bases. For the manually created subset,  
174 domain experts selected high-quality images from the Chinese herbal medicine image repository  
175 and generated corresponding questions based on each herb's name and medicinal properties. The  
176 automatically generated subset was produced through a procedural pipeline. For example, an image  
177 labeled as *Astragalus membranaceus* (*Huangqi*) was selected as the correct answer, while three  
178 distractor images were randomly sampled from the knowledge base. A question was then constructed  
179 using a predefined template library, such as "Which of the following images shows *Huangqi*?" The  
180 design of tongue image understanding tasks followed a similar approach. Details of the construction  
181 process and implementation code can be found in **Appendix G**.

### 182 3.3 Deduplication and Preprocessing

183 Detecting duplication and semantic similarity in the data is critical for both model evaluation and  
184 training, as it helps prevent evaluation failures and reduces the risk of overfitting caused by redundant  
185 content. Given the diverse sources of the original data, we conducted a comprehensive similarity  
186 detection process on the aggregated dataset and removed highly similar questions to enhance overall  
187 data quality. The methods employed included string edit distance[42], TF-IDF [43, 44] with cosine  
188 similarity, and BERT-based[45, 46] semantic encoding. Subsequently, all questions and answers  
189 were manually reviewed by two licensed physicians. The selection criteria and detailed experimental  
190 procedures are provided in **Appendix I**. Subsequently, we divided the dataset into three subsets: 10%  
191 for evaluation, 10% for validation, and 80% for training. To ensure balanced representation, each  
192 subset contains question-answer pairs spanning all subfields.

### 193 3.4 Dataset Statistics

194 Table 2 presents the statistics of all constructed question-answer pairs across different categories.  
195 The TCM-Ladder dataset comprises a total of 52,169 TCM-related QA instances, including 6,061  
196 Chinese herbal medicine images and 1,394 annotated tongue image patches. The distribution of each  
197 data type is illustrated in Figure 2.

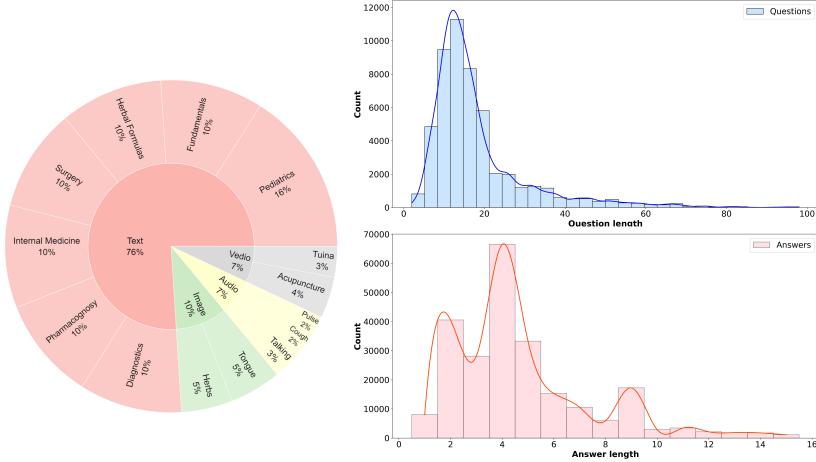


Figure 2: Data distribution and length statistics in TCM-Ladder. The left illustrates the dataset composition across text, image, and audio modalities, along with TCM subfields. The right plots show the distribution of question and answer lengths.

## 198 4 Ladder-Score

199 Evaluating free-form question answering  
200 presents notable challenges, as the responses are  
201 often descriptive and lack a predefined standard  
202 format. This issue is further exacerbated in  
203 the context of TCM diagnostic tasks, where  
204 large language models are capable of generating  
205 diverse and nuanced answers. Even when the  
206 expressions differ, the underlying responses  
207 may still be factually correct. Traditional  
208 evaluation metrics such as BLEU[47] and  
209 ROUGE[48] often fail to capture this semantic  
210 equivalence adequately. Recently proposed  
211 methods[49, 50, 51] employ instruction-tuned  
212 models to score candidate answers on a rubric-  
213 based scale. We propose a novel evaluation  
214 metric for TCM question answering, named  
215 **Ladder-Score**. This score comprises two  
216 components: *TermScore*, which assesses the  
217 accuracy and completeness of TCM terminology usage, and *SemanticScore*, derived from large  
218 language models to evaluate multiple aspects including logical consistency, semantic accuracy,  
219 comprehensiveness of knowledge, and fluency of expression. As shown in Equation 1, the  
220 Ladder-Score is a weighted combination of these two components:

$$\text{Ladder-Score} = \alpha \cdot \text{TermScore} + \beta \cdot \text{SemanticScore} \quad (1)$$

221 where  $\alpha = 0.4$  and  $\beta = 0.6$ , which can be adjusted based on practical needs.

222 For details on the scoring criteria, terminology dictionary, and calculation examples, please refer to  
223 Appendix H.

## 224 5 Experiments

### 225 5.1 Experiment Setup

226 We evaluated 9 state-of-the-art general domain LLMs and 5 TCM-specific models on the TCM-  
227 Ladder dataset across five task settings: single-choice questions, multiple-choice questions, fill-in-

228 the-blank questions, image-based understanding, and long-form dialogue tasks. Evaluations were  
229 conducted under zero-shot settings, and models received only the task instructions as input. For  
230 single-choice and image understanding tasks, we used the Top-1 prediction accuracy[52] as the  
231 primary evaluation metric. For multiple-choice tasks, we adopted exact match accuracy to assess  
232 performance comprehensively. For fill-in-the-blank and long-form dialogue tasks, we evaluated  
233 models using metrics such as accuracy, BLEU[47], ROUGE[48], METEOR[53] and BERTScore[54].

234 **5.2 Model Training**

235 We trained two models using the TCM-Ladder dataset. The first is *Bencao*[21], an online model  
236 fine-tuned from ChatGPT, and the second is Ladder-base, which is built upon the pretrained *Qwen2.5-  
237 7B-Instruct*[55] model and enhanced with Group Relative Policy Optimization (GRPO)[56] to  
238 improve its reasoning capabilities. The Bencao model was trained on knowledge extracted from  
239 over 700 classical Chinese medicine books, none of which contained any question-answer pairs.  
240 Additionally, the training subset of TCM-Ladder was used as its knowledge base.  
241 The GRPO stage for Ladder-base was conducted on two NVIDIA A100 PCIe GPUs (80GB each).  
242 The temperature and top-p sampling of Ladder-base were 0.7 and 0.8. Training was performed for 2  
243 epochs with a group size of 6 and a batch size of 12, resulting in a total training time of approximately  
244 60 hours. Model training and inference were implemented using HuggingFace Transformers, while  
245 the GRPO process was carried out using the TRL (Transformer Reinforcement Learning) library[57].  
246 Details of the training process can be found in **Appendix C**.

247 **5.3 Human Evaluation**

248 We conducted a human evaluation on 20% of the TCM-Ladder test set. Due to the coverage of  
249 multiple subfields, establishing a reliable human upper bound poses a significant challenge, as  
250 accurately answering questions across all domains requires extensive interdisciplinary expertise. To  
251 investigate this issue, we recruited two licensed clinical TCM physicians, who were not involved  
252 in the original data annotation. Human evaluators were asked to select the correct answers based  
253 on the question stems and to identify the correct herbal medicine and tongue images. In terms of  
254 top-1 accuracy for answer retrieval, the human evaluators achieved a performance of 64%, which  
255 was approximately 4% lower than that of the best-performing model (*Bencao*). This suggests that  
256 LLMs may already possess strong comprehension capabilities in the domains of herbal medicine and  
257 tongue image recognition.

258 **5.4 Main Results**

259 **5.4.1 Text-Based Single-Choice Question Answering**

260 As shown in Figure 3, Ladder-base (ours) consistently outperforms other models across all subject  
261 areas, achieving the highest overall accuracy. Notably, its performance is especially strong in  
262 Pharmacognosy, Herbal Formulas, and Pediatrics, where exact match scores exceed 0.85. Our other  
263 model, Bencao (ours), also demonstrates robust performance, particularly in Diagnostics and Internal  
264 Medicine. Among the general domain LLMs, Gemini 2.5 Pro, Deepseek, and Tongyi Qwen show  
265 relatively stable accuracy across domains, with scores ranging from 0.65 to 0.75, though they still  
266 fall short compared to domain-specialized models. In contrast, Claude 3, GPT-4o mini, and Bentsao  
267 underperform, especially in the more clinically nuanced domains such as Surgery and Pediatrics,  
268 suggesting limited capability in handling complex, multi-faceted TCM tasks. These findings highlight  
269 the advantage of domain-specific fine-tuning and multi-source integration, as utilized in Ladder-base,  
270 for enhancing the accuracy and generalization of LLMs on structured TCM knowledge assessments.

271 **5.4.2 Visual Question Answering**

272 To further assess the models’ capability in visual understanding tasks within Traditional Chinese  
273 Medicine (TCM), we evaluated 10 large language models (LLMs) on two image-based benchmarks:  
274 Herbs classification and Tongue image diagnosis. As illustrated in Figure 4, performance varies  
275 considerably across models. Among the evaluated models, Bencao (ours) achieves the highest  
276 accuracy in both tasks, with over 80% on herb recognition and above 65% on tongue classification,  
277 demonstrating strong multimodal understanding grounded in TCM-specific training. General domain

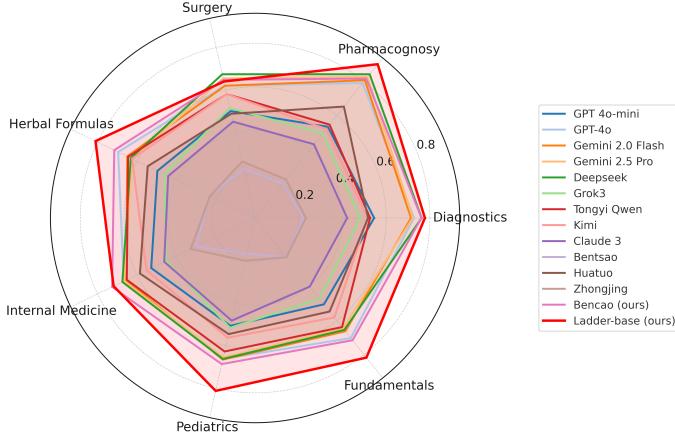


Figure 3: Performance of general-domain and TCM-specific language models on single and multiple-choice question answering tasks.

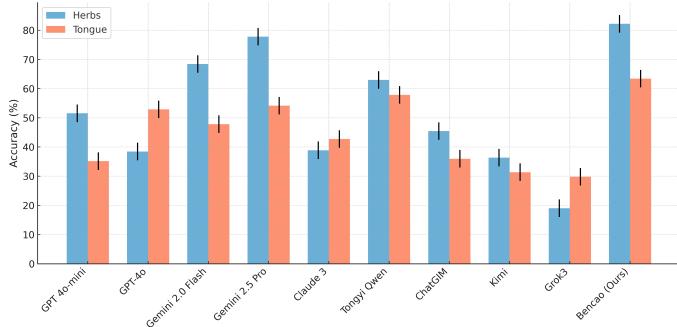


Figure 4: The performance of large language models on questions regarding Chinese herbal medicine and tongue images.

278 LLMs such as Gemini 2.5 Pro, Gemini 2.0 Flash, and Tongyi Qwen exhibit moderate performance,  
 279 with herb classification accuracy around 65–75%, but show a relative drop in tongue image tasks  
 280 (around 50–60%), likely due to the greater complexity and domain specificity of tongue diagnosis.

281 In contrast, models like GPT-4o, Claude 3, Kimi, and Grok3 demonstrate limited performance,  
 282 particularly in the **tongue classification task**, where accuracies often fall below 40%. This reveals  
 283 their insufficient visual comprehension of TCM-related imagery. It is worth noting that models such  
 284 as Ladder-base and Zhongjing are not included in this figure, as they are not equipped with image  
 285 understanding capabilities at this stage. Their current design focuses on structured text-based TCM  
 286 evaluation and does not support visual input.

#### 287 5.4.3 Diagnostic dialogue and Fill-in-the blank Questions

288 As shown in Table 3, in the diagnostic dialogue task, our model Ladder-base achieved the highest  
 289 scores in BLEU-4 (0.0249), ROUGE-L (0.2431), and METEOR (0.2268), while also maintaining  
 290 a strong Ladder-score (0.803). This indicates that Ladder-base generates answers with high lexical  
 291 similarity, semantic accuracy, and alignment with TCM diagnostic logic. Notably, Tongyi Qwen  
 292 achieved the best **Ladder-score (0.861)** and the highest METEOR (0.2328), showcasing its strength  
 293 in generating fluently worded responses. Bencao (ours) achieved the best BERTScore (0.9663),  
 294 reflecting its semantic closeness to gold references.

295 In the fill-in-the-blank task, Bencao significantly outperformed all other models, achieving the highest  
 296 Exact Match Accuracy of 0.9034, followed by Tongyi Qwen (0.8786) and Deepseek (0.874). Our  
 297 Ladder-base model also performed competitively with 0.8623 accuracy, further demonstrating its  
 298 generalizability beyond free-form dialogue. Overall, the results demonstrate that Ladder-base excels

299 in structured diagnostic dialogue tasks, generating semantically accurate and logically coherent  
300 responses, while *Bencao* shows outstanding performance in fill-in-the-blank tasks, reflecting strong  
301 factual recall and precise terminology usage. Domain-specific models consistently outperform general  
302 domain LLMs, particularly in tasks that require accurate retrieval of structured TCM knowledge and  
303 professional terms.

Table 3: Performance Comparison on Diagnostic Dialogue and Fill-in-the-blank Tasks

Model	Diagnostic dialogue					Fill-in-the-blank
	BLEU-4	ROUGE-L	METEOR	BERTScore	Ladder-score	
GPT 4o-mini	0.0034	0.1125	0.119	0.9433	0.718	0.432
GPT-4o	0.0040	0.1447	0.2073	0.9620	0.828	0.514
Gemini 2.0 Flash	0.0067	0.1518	0.2155	0.9633	0.836	0.436
Gemini 2.5 Pro	0.0180	0.1353	0.2393	0.9605	0.859	0.7143
Deepseek	0.0047	0.1533	0.1293	0.9455	0.825	0.874
Grok3	0.0063	0.1751	0.1691	0.9526	0.686	0.6389
Tongyi Qwen	0.0225	0.1818	<b>0.2328</b>	0.9642	<b>0.861</b>	0.8786
Kimi	0.0100	0.1878	0.1586	0.9559	0.708	0.8378
Claude 3	0.0068	0.2267	0.2203	0.9561	0.756	0.489
Bentsao	0.0024	0.1135	0.1725	0.9531	0.613	0.162
Huatuo	0.0086	0.1375	0.1742	0.9635	0.855	0.2347
Zhongjing	0.0044	0.1951	0.1134	0.9539	0.573	0.2167
Bencao (ours)	0.0073	0.2156	0.2013	<b>0.9663</b>	0.791	<b>0.9034</b>
Ladder-base (ours)	<b>0.0249</b>	<b>0.2431</b>	0.2268	0.9549	0.803	0.8623

## 304 6 Application Website

305 In addition to releasing the raw dataset, we provide access to all TCM-Ladder data and leaderboard re-  
306 sults through an interactive website (<https://tcmladder.com/>). This platform enables researchers  
307 to explore, verify, and contribute to the open-access data. We encourage the research community to  
308 submit additional data through the platform, and we intend to expand the dataset continuously as part  
309 of our ongoing efforts. Our objective is to establish a long-term and reliable data foundation for the  
310 training and evaluation TCM-specific LLMs.

## 311 7 Limitations and Societal Impact

312 Although TCM-Ladder encompasses question-answer pairs from multiple disciplines within TCM,  
313 its current scale remains insufficient to cover the full breadth of TCM knowledge. TCM diagnosis is  
314 inherently a multimodal process—textual information represents only one component. At present,  
315 the utilization of data related to tongue diagnosis, pulse diagnosis, and olfactory inspection remains  
316 limited, and these modalities require further supplementation and enrichment. Expanding and  
317 continuously updating the scope and scale of data included in TCM-Ladder will be a critical direction  
318 for future research.

## 319 8 Conclusion

320 We introduce **TCM-Ladder**, the first multimodal benchmark dataset designed explicitly for evaluating  
321 LLMs in the context of TCM. In addition, we propose a novel evaluation metric, Ladder-Score,  
322 which enables more precise analysis of the semantic alignment between candidate and reference  
323 answers. We conduct comprehensive experiments involving 9 state-of-the-art general domain and 5  
324 TCM-specific LLMs, marking the first systematic comparison on a unified benchmark. Furthermore,  
325 we fine-tune two TCM-specific models using a subset of TCM-Ladder, and observe significant  
326 performance improvements over zero-shot baselines. Our work establishes a reproducible and  
327 extensible benchmark for TCM-specific, providing a foundation for future development and evaluation  
328 in this emerging research area.

329 **References**

- 330 [1] Jin-Ling Tang, Bao-Yan Liu, and Kan-Wen Ma. Traditional chinese medicine. *The Lancet*, 372(9654):1938–  
331 1940, 2008.
- 332 [2] Dennis Normile. The new face of traditional chinese medicine. *Science*, 299(5604):188–190, 2003.
- 333 [3] Tianhan Xue and Rustum Roy. Studying traditional chinese medicine. *Science*, 300(5620):740–741, 2003.
- 334 [4] Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. Pathvqa: 30000+ questions for  
335 medical visual question answering. *arXiv preprint arXiv:2003.10286*, 2020.
- 336 [5] Di Jin, Eileen Pan, Nassim Oufattolle, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease  
337 does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied  
338 Sciences*, 11(14):6421, 2021.
- 339 [6] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren,  
340 Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database  
341 of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019.
- 342 [7] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales,  
343 Ajay Tanwani, Heather Cole-Lewis, Stephen Pfahl, et al. Large language models encode clinical knowledge.  
344 *Nature*, 620(7972):172–180, 2023.
- 345 [8] Boya Zhang, Alban Bornet, Anthony Yazdani, Philipp Khlebnikov, Marija Milutinovic, Hossein  
346 Rouhizadeh, Poorya Amini, and Douglas Teodoro. A dataset for evaluating clinical research claims  
347 in large language models. *Scientific Data*, 12(1):86, 2025.
- 348 [9] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,  
349 Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*,  
350 2024.
- 351 [10] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
352 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv  
353 preprint arXiv:2303.08774*, 2023.
- 354 [11] Google DeepMind. Gemini 2.0: Deepmind’s multimodal llm. <https://deepmind.google/technologies/gemini/>, 2023. Accessed: 2025-05-16.
- 356 [12] Google DeepMind. Gemini 2.5: Next-generation multimodal reasoning model. <https://deepmind.google/technologies/gemini/>, 2024. Accessed: 2025-05-16.
- 358 [13] Xiao Bi, Deli Chen, Guanting Chen, Shanhua Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai  
359 Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism.  
360 *arXiv preprint arXiv:2401.02954*, 2024.
- 361 [14] Jun Li, Pei Yuan, Xiaojuan Hu, Jingbin Huang, Longtao Cui, Ji Cui, Xuxiang Ma, Tao Jiang, Xinghua Yao,  
362 Jiacai Li, et al. A tongue features fusion approach to predicting prediabetes and diabetes with machine  
363 learning. *Journal of biomedical informatics*, 115:103693, 2021.
- 364 [15] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han,  
365 Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 366 [16] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao,  
367 Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv  
368 preprint arXiv:2501.12599*, 2025.
- 369 [17] Anthropic. Claude 3. <https://www.anthropic.com/index/clause-3>, 2024. Accessed: 2025-05-16.
- 370 [18] Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: Tuning  
371 llama model with chinese medical knowledge. *arXiv preprint arXiv:2304.06975*, 2023.
- 372 [19] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou,  
373 and Benyou Wang. Huatuogpt-o1, towards medical complex reasoning with llms. *arXiv preprint  
374 arXiv:2412.18925*, 2024.
- 375 [20] Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan.  
376 Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback  
377 and real-world multi-turn dialogue. In *Proceedings of the AAAI conference on artificial intelligence*,  
378 volume 38, pages 19368–19376, 2024.

- 379 [21] Bencao. <https://chatgpt.com/g/g-6750c5262fb48191a08ef4d899a3dd1f-bencao>, 2025. Accessed: May 16, 2025.
- 380 [22] Kai Gao, Liu Liu, Shuangshuang Lei, Zhinong Li, Peipei Huo, Zhihao Wang, Lei Dong, Wenxin Deng, Dechao Bu, Xiaoxi Zeng, et al. Herb 2.0: an updated database integrating clinical and experimental evidence for traditional chinese medicine. *Nucleic Acids Research*, 53(D1):D1404–D1414, 2025.
- 381 [23] Jianquan Li, Xidong Wang, Xiangbo Wu, Zhiyi Zhang, Xiaolong Xu, Jie Fu, Prayag Tiwari, Xiang Wan, and Benyou Wang. Huatuo-26m, a large-scale chinese medical qa dataset. *arXiv preprint arXiv:2305.01526*, 2023.
- 382 [24] Ningyu Zhang, Mosha Chen, Zhen Bi, Xiaozhuan Liang, Lei Li, Xin Shang, Kangping Yin, Chuanqi Tan, Jian Xu, Fei Huang, et al. Cblue: A chinese biomedical language understanding evaluation benchmark. *arXiv preprint arXiv:2106.08087*, 2021.
- 383 [25] Wei Zhu, Xiaoling Wang, Huanran Zheng, Mosha Chen, and Buzhou Tang. Promptcblue: A chinese prompt tuning benchmark for the medical domain. *arXiv preprint arXiv:2310.14151*, 2023.
- 384 [26] Wenjing Yue, Xiaoling Wang, Wei Zhu, Ming Guan, Huanran Zheng, Pengfei Wang, Changzhi Sun, and Xin Ma. Tcmbench: A comprehensive benchmark for evaluating large language models in traditional chinese medicine. *arXiv preprint arXiv:2406.01126*, 2024.
- 385 [27] Zhe Wang, Meng Hao, Suyuan Peng, Yuyan Huang, Yiwei Lu, Keyu Yao, Xiaolin Yang, and Yan Zhu. Tcmeval-sdt: a benchmark dataset for syndrome differentiation thought of traditional chinese medicine. *Scientific Data*, 12(1):437, 2025.
- 386 [28] Tianai Huang, Lu Lu, Jiayuan Chen, Lihao Liu, Junjun He, Yuping Zhao, Wenchao Tang, and Jie Xu. Tcm-3ceval: A triaxial benchmark for assessing responses from large language models in traditional chinese medicine. *arXiv preprint arXiv:2503.07041*, 2025.
- 387 [29] Ping Yu, Kaitao Song, Fengchen He, Ming Chen, and Jianfeng Lu. Tcmd: A traditional chinese medicine qa dataset for evaluating large language models. *arXiv preprint arXiv:2406.04941*, 2024.
- 388 [30] Jinyang Zhu, Qingyue Gong, Chunfang Zhou, and Huidan Luan. Zhongjing: A locally deployed large language model for traditional chinese medicine and corresponding evaluation methodology: A large language model for data fine-tuning in the field of traditional chinese medicine, and a new evaluation method called tcmeval are proposed. In *Proceedings of the 2023 4th International Symposium on Artificial Intelligence for Medicine Science*, pages 1036–1042, 2023.
- 389 [31] Chenlu Guo, Nuo Xu, Yi Chang, and Yuan Wu. Chbench: A chinese dataset for evaluating health in large language models. *arXiv preprint arXiv:2409.15766*, 2024.
- 390 [32] Mianxin Liu, Weiguo Hu, Jinru Ding, Jie Xu, Xiaoyang Li, Lifeng Zhu, Zhian Bai, Xiaoming Shi, Benyou Wang, Haitao Song, et al. Medbench: A comprehensive, standardized, and reliable benchmarking system for evaluating chinese medical large language models. *Big Data Mining and Analytics*, 7(4):1116–1128, 2024.
- 391 [33] Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, et al. Cmb: A comprehensive medical benchmark in chinese. *arXiv preprint arXiv:2308.08833*, 2023.
- 392 [34] Junling Liu, Peilin Zhou, Yining Hua, Dading Chong, Zhongyu Tian, Andrew Liu, Helin Wang, Chenyu You, Zhenhua Guo, Lei Zhu, et al. Benchmarking large language models on cmexam-a comprehensive chinese medical exam dataset. *Advances in Neural Information Processing Systems*, 36:52430–52452, 2023.
- 393 [35] Kee C Huang. *The pharmacology of Chinese herbs*. CRC press, 1998.
- 394 [36] Jiacheng Xie, Congcong Jing, Ziyang Zhang, Jiatuo Xu, Ye Duan, and Dong Xu. Digital tongue image analyses for health assessment. *Medical Review*, 1(2):172–198, 2021.
- 395 [37] Ye Duan and Dong Xu. Itongue: an iphone app for personal health monitoring based on tongue image. 2014.
- 396 [38] Ivandro Sanches, Victor V Gomes, Carlos Caetano, Lizeth SB Cabrera, Vinicius H Cene, Thomas Beltrame, Wonkyu Lee, Sanghyun Baek, and Otávio AB Penatti. Mimic-bp: A curated dataset for blood pressure estimation. *Scientific Data*, 11(1):1233, 2024.
- 397 [39] zhidong zhang. pulse dataset, 2023.

- 430 [40] Ahsan Mehmood, Asma Sarouji, M Mahboob Ur Rahman, and Tareq Y Al-Naffouri. Your smartphone  
 431 could act as a pulse-oximeter and as a single-lead ecg. *Scientific Reports*, 13(1):19277, 2023.
- 432 [41] Andrea Nemcova, Enikö Vargova, Radovan Smisek, Lucie Marsanova, Lukas Smital, and Martin Vitek.  
 433 Brno university of technology smartphone ppg database (but ppg): Annotated dataset for ppg quality  
 434 assessment and heart rate estimation. *BioMed Research International*, 2021(1):3453007, 2021.
- 435 [42] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals. In  
 436 *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union, 1966.
- 437 [43] Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information retrieval*,  
 438 volume 39. Cambridge University Press Cambridge, 2008.
- 439 [44] Gerard Salton, Anita Wong, and Chung-Shu Yang. A vector space model for automatic indexing. *Communi-  
 440 cations of the ACM*, 18(11):613–620, 1975.
- 441 [45] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidi-  
 442 rectional transformers for language understanding. In *Proceedings of the 2019 conference of the North  
 443 American chapter of the association for computational linguistics: human language technologies, volume  
 444 1 (long and short papers)*, pages 4171–4186, 2019.
- 445 [46] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks.  
 446 *arXiv preprint arXiv:1908.10084*, 2019.
- 447 [47] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation  
 448 of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational  
 449 Linguistics*, pages 311–318, 2002.
- 450 [48] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches  
 451 out*, pages 74–81, 2004.
- 452 [49] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat,  
 453 Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing  
 454 Systems*, 36:55006–55021, 2023.
- 455 [50] Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin  
 456 Shin, Sungdong Kim, James Thorne, et al. Prometheus: Inducing fine-grained evaluation capability in  
 457 language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- 458 [51] Seongyun Lee, Seungone Kim, Sue Park, Geewook Kim, and Minjoon Seo. Prometheus-vision: Vision-  
 459 language model as a judge for fine-grained evaluation. In *Findings of the Association for Computational  
 460 Linguistics ACL 2024*, pages 11286–11315, 2024.
- 461 [52] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional  
 462 neural networks. *Advances in neural information processing systems*, 25, 2012.
- 463 [53] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved  
 464 correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evalua-  
 465 tion measures for machine translation and/or summarization*, pages 65–72, 2005.
- 466 [54] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text  
 467 generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- 468 [55] An Yang, Baosong Yang, Beichen Zhang, Bin yuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng  
 469 Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- 470 [56] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan  
 471 Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language  
 472 models. *arXiv preprint arXiv:2402.03300*, 2024.
- 473 [57] Shengchao Hu, Li Shen, Ya Zhang, Yixin Chen, and Dacheng Tao. On transforming reinforcement learning  
 474 with transformers: The development trajectory. *IEEE Transactions on Pattern Analysis and Machine  
 475 Intelligence*, 2024.

476 **NeurIPS Paper Checklist**

477 **1. Claims**

478 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's  
479 contributions and scope?

480 Answer: [Yes]

481 Justification: As stated in the abstract and introduction, to address the current scarcity of multimodal  
482 datasets in Traditional Chinese Medicine (TCM), we propose a multimodal TCM question-answering  
483 dataset. We evaluate it using 9 general-purpose and 5 TCM-specific large language models, and  
484 present the dataset and leaderboard through an online platform.

485 Guidelines:

- 486 • The answer NA means that the abstract and introduction do not include the claims made in the  
487 paper.  
488 • The abstract and/or introduction should clearly state the claims made, including the contributions  
489 made in the paper and important assumptions and limitations. A No or NA answer to this  
490 question will not be perceived well by the reviewers.  
491 • The claims made should match theoretical and experimental results, and reflect how much the  
492 results can be expected to generalize to other settings.  
493 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
494 attained by the paper.

495 **2. Limitations**

496 Question: Does the paper discuss the limitations of the work performed by the authors?

497 Answer: [Yes]

498 Justification: See **Section 7. Limitations and Societal Impact**.

499 Guidelines:

- 500 • The answer NA means that the paper has no limitation while the answer No means that the paper  
501 has limitations, but those are not discussed in the paper.  
502 • The authors are encouraged to create a separate "Limitations" section in their paper.  
503 • The paper should point out any strong assumptions and how robust the results are to violations of  
504 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
505 asymptotic approximations only holding locally). The authors should reflect on how these  
506 assumptions might be violated in practice and what the implications would be.  
507 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
508 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
509 assumptions, which should be articulated.  
510 • The authors should reflect on the factors that influence the performance of the approach. For  
511 example, a facial recognition algorithm may perform poorly when image resolution is low or  
512 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide  
513 closed captions for online lectures because it fails to handle technical jargon.  
514 • The authors should discuss the computational efficiency of the proposed algorithms and how  
515 they scale with dataset size.  
516 • If applicable, the authors should discuss possible limitations of their approach to address problems  
517 of privacy and fairness.  
518 • While the authors might fear that complete honesty about limitations might be used by reviewers  
519 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
520 aren't acknowledged in the paper. The authors should use their best judgment and recognize  
521 that individual actions in favor of transparency play an important role in developing norms that  
522 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
523 honesty concerning limitations.

524 **3. Theory assumptions and proofs**

525 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete  
526 (and correct) proof?

527 Answer: [Yes]

528 Justification: See **Appendix H**

529 Guidelines:

- 530 • The answer NA means that the paper does not include theoretical results.

- 531           • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.  
 532           • All assumptions should be clearly stated or referenced in the statement of any theorems.  
 533           • The proofs can either appear in the main paper or the supplemental material, but if they appear in  
 534           the supplemental material, the authors are encouraged to provide a short proof sketch to provide  
 535           intuition.  
 536           • Inversely, any informal proof provided in the core of the paper should be complemented by  
 537           formal proofs provided in appendix or supplemental material.  
 538           • Theorems and Lemmas that the proof relies upon should be properly referenced.

539          **4. Experimental result reproducibility**

540          Question: Does the paper fully disclose all the information needed to reproduce the main experimental  
 541          results of the paper to the extent that it affects the main claims and/or conclusions of the paper  
 542          (regardless of whether the code and data are provided or not)?

543          Answer: [Yes]

544          Justification: To ensure the reproducibility, we have publicly released all datasets, as well as the  
 545          code and access links used for models evaluation. The training process of Ladder-base is also made  
 546          available on GitHub. Please see Appendix A for details.

547          Guidelines:

- 548           • The answer NA means that the paper does not include experiments.  
 549           • If the paper includes experiments, a No answer to this question will not be perceived well by the  
 550           reviewers: Making the paper reproducible is important, regardless of whether the code and data  
 551           are provided or not.  
 552           • If the contribution is a dataset and/or model, the authors should describe the steps taken to make  
 553           their results reproducible or verifiable.  
 554           • Depending on the contribution, reproducibility can be accomplished in various ways. For  
 555           example, if the contribution is a novel architecture, describing the architecture fully might suffice,  
 556           or if the contribution is a specific model and empirical evaluation, it may be necessary to either  
 557           make it possible for others to replicate the model with the same dataset, or provide access to  
 558           the model. In general, releasing code and data is often one good way to accomplish this, but  
 559           reproducibility can also be provided via detailed instructions for how to replicate the results,  
 560           access to a hosted model (e.g., in the case of a large language model), releasing of a model  
 561           checkpoint, or other means that are appropriate to the research performed.  
 562           • While NeurIPS does not require releasing code, the conference does require all submissions  
 563           to provide some reasonable avenue for reproducibility, which may depend on the nature of the  
 564           contribution. For example  
 565           (a) If the contribution is primarily a new algorithm, the paper should make it clear how to  
 566           reproduce that algorithm.  
 567           (b) If the contribution is primarily a new model architecture, the paper should describe the  
 568           architecture clearly and fully.  
 569           (c) If the contribution is a new model (e.g., a large language model), then there should either be  
 570           a way to access this model for reproducing the results or a way to reproduce the model (e.g.,  
 571           with an open-source dataset or instructions for how to construct the dataset).  
 572           (d) We recognize that reproducibility may be tricky in some cases, in which case authors are  
 573           welcome to describe the particular way they provide for reproducibility. In the case of  
 574           closed-source models, it may be that access to the model is limited in some way (e.g.,  
 575           to registered users), but it should be possible for other researchers to have some path to  
 576           reproducing or verifying the results.

577          **5. Open access to data and code**

578          Question: Does the paper provide open access to the data and code, with sufficient instructions to  
 579          faithfully reproduce the main experimental results, as described in supplemental material?

580          Answer: [Yes]

581          Justification: We have released all datasets and the code used for evaluating the models, along with the  
 582          training process of Ladder-base, which is publicly available on GitHub. The data and code resources  
 583          can be found in the **Abstract** and **Appendix A**.

584          Guidelines:

- 585           • The answer NA means that paper does not include experiments requiring code.  
 586           • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.

- 588           • While we encourage the release of code and data, we understand that this might not be possible,  
 589           so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless  
 590           this is central to the contribution (e.g., for a new open-source benchmark).
- 591           • The instructions should contain the exact command and environment needed to run to reproduce  
 592           the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 593           • The authors should provide instructions on data access and preparation, including how to access  
 594           the raw data, preprocessed data, intermediate data, and generated data, etc.
- 595           • The authors should provide scripts to reproduce all experimental results for the new proposed  
 596           method and baselines. If only a subset of experiments are reproducible, they should state which  
 597           ones are omitted from the script and why.
- 598           • At submission time, to preserve anonymity, the authors should release anonymized versions (if  
 599           applicable).
- 600           • Providing as much information as possible in supplemental material (appended to the paper) is  
 601           recommended, but including URLs to data and code is permitted.

603       **6. Experimental setting/details**

604       Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
 605           how they were chosen, type of optimizer, etc.) necessary to understand the results?

606       Answer: [Yes]

607       Justification: See **Section 5**, **Section 3.3** and **Appendix C**.

608       Guidelines:

- 609           • The answer NA means that the paper does not include experiments.
- 610           • The experimental setting should be presented in the core of the paper to a level of detail that is  
 611           necessary to appreciate the results and make sense of them.
- 612           • The full details can be provided either with the code, in appendix, or as supplemental material.

613       **7. Experiment statistical significance**

614       Question: Does the paper report error bars suitably and correctly defined or other appropriate information  
 615           about the statistical significance of the experiments?

616       Answer: [Yes]

617       Justification: In Figure 4, we include error curves based on a 3% error margin. However, due to the  
 618           high cost associated with repeated API calls, we conducted only a single run of the experiment. As  
 619           such, no statistically derived errors are provided.

620       Guidelines:

- 621           • The answer NA means that the paper does not include experiments.
- 622           • The authors should answer "Yes" if the results are accompanied by error bars, confidence  
 623           intervals, or statistical significance tests, at least for the experiments that support the main claims  
 624           of the paper.
- 625           • The factors of variability that the error bars are capturing should be clearly stated (for example,  
 626           train/test split, initialization, random drawing of some parameter, or overall run with given  
 627           experimental conditions).
- 628           • The method for calculating the error bars should be explained (closed form formula, call to a  
 629           library function, bootstrap, etc.)
- 630           • The assumptions made should be given (e.g., Normally distributed errors).
- 631           • It should be clear whether the error bar is the standard deviation or the standard error of the  
 632           mean.
- 633           • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report  
 634           a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is  
 635           not verified.
- 636           • For asymmetric distributions, the authors should be careful not to show in tables or figures  
 637           symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 638           • If error bars are reported in tables or plots, The authors should explain in the text how they were  
 639           calculated and reference the corresponding figures or tables in the text.

640       **8. Experiments compute resources**

641       Question: For each experiment, does the paper provide sufficient information on the computer  
 642           resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

643       Answer: [Yes]

644 Justification: See **Section 5.2, Appendix C** and **Appendix D**.

645 Guidelines:

- 646 • The answer NA means that the paper does not include experiments.
- 647 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 648 • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 649 • The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

650 **9. Code of ethics**

651 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code  
652 of Ethics <https://neurips.cc/public/EthicsGuidelines>?

653 Answer: [Yes]

654 Justification: The research complies with the NeurIPS Code of Ethics. The tongue image data used in  
655 our dataset were approved by the institutional review board. All personally identifiable information  
656 has been thoroughly anonymized or removed to ensure the privacy and protection of the individuals  
657 involved.

658 Guidelines:

- 659 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 660 • If the authors answer No, they should explain the special circumstances that require a deviation  
661 from the Code of Ethics.
- 662 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due  
663 to laws or regulations in their jurisdiction).

664 **10. Broader impacts**

665 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts  
666 of the work performed?

667 Answer: [Yes]

668 Justification: See **Section 7**.

669 Guidelines:

- 670 • The answer NA means that there is no societal impact of the work performed.
- 671 • If the authors answer NA or No, they should explain why their work has no societal impact or  
672 why the paper does not address societal impact.
- 673 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,  
674 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-  
675 ment of technologies that could make decisions that unfairly impact specific groups), privacy  
676 considerations, and security considerations.
- 677 • The conference expects that many papers will be foundational research and not tied to particular  
678 applications, let alone deployments. However, if there is a direct path to any negative applications,  
679 the authors should point it out. For example, it is legitimate to point out that an improvement in  
680 the quality of generative models could be used to generate deepfakes for disinformation. On the  
681 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks  
682 could enable people to train models that generate Deepfakes faster.
- 683 • The authors should consider possible harms that could arise when the technology is being used  
684 as intended and functioning correctly, harms that could arise when the technology is being used  
685 as intended but gives incorrect results, and harms following from (intentional or unintentional)  
686 misuse of the technology.
- 687 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies  
688 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitor-  
689 ing misuse, mechanisms to monitor how a system learns from feedback over time, improving the  
690 efficiency and accessibility of ML).

691 **11. Safeguards**

692 Question: Does the paper describe safeguards that have been put in place for responsible release of  
693 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or  
694 scraped datasets)?

695 Answer: [Yes]

700 Justification: See **Appendix E**.

701 Guidelines:

- 702 • The answer NA means that the paper poses no such risks.
- 703 • Released models that have a high risk for misuse or dual-use should be released with necessary
- 704 safeguards to allow for controlled use of the model, for example by requiring that users adhere to
- 705 usage guidelines or restrictions to access the model or implementing safety filters.
- 706 • Datasets that have been scraped from the Internet could pose safety risks. The authors should
- 707 describe how they avoided releasing unsafe images.
- 708 • We recognize that providing effective safeguards is challenging, and many papers do not require
- 709 this, but we encourage authors to take this into account and make a best faith effort.

## 710 12. Licenses for existing assets

711 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,  
712 properly credited and are the license and terms of use explicitly mentioned and properly respected?

713 Answer: [Yes]

714 Justification: See **Appendix B**.

715 Guidelines:

- 716 • The answer NA means that the paper does not use existing assets.
- 717 • The authors should cite the original paper that produced the code package or dataset.
- 718 • The authors should state which version of the asset is used and, if possible, include a URL.
- 719 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 720 • For scraped data from a particular source (e.g., website), the copyright and terms of service of
- 721 that source should be provided.
- 722 • If assets are released, the license, copyright information, and terms of use in the package should
- 723 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for
- 724 some datasets. Their licensing guide can help determine the license of a dataset.
- 725 • For existing datasets that are re-packaged, both the original license and the license of the derived
- 726 asset (if it has changed) should be provided.
- 727 • If this information is not available online, the authors are encouraged to reach out to the asset's
- 728 creators.

## 729 13. New assets

730 Question: Are new assets introduced in the paper well documented and is the documentation provided  
731 alongside the assets?

732 Answer: [Yes]

733 Justification: Please see **Appendix G**. We provide a detailed description of the image acquisition  
734 procedures for Chinese herbal medicine samples and tongue images used in our study.

735 Guidelines:

- 736 • The answer NA means that the paper does not release new assets.
- 737 • Researchers should communicate the details of the dataset/code/model as part of their sub-
- 738 missions via structured templates. This includes details about training, license, limitations,
- 739 etc.
- 740 • The paper should discuss whether and how consent was obtained from people whose asset is
- 741 used.
- 742 • At submission time, remember to anonymize your assets (if applicable). You can either create an
- 743 anonymized URL or include an anonymized zip file.

## 744 14. Crowdsourcing and research with human subjects

745 Question: For crowdsourcing experiments and research with human subjects, does the paper include  
746 the full text of instructions given to participants and screenshots, if applicable, as well as details about  
747 compensation (if any)?

748 Answer: [Yes]

749 Justification: We describe the tongue image acquisition process in **Appendix G**.

750 Guidelines:

- 751 • The answer NA means that the paper does not involve crowdsourcing nor research with human
- 752 subjects.

- 753           • Including this information in the supplemental material is fine, but if the main contribution of the  
754           paper involves human subjects, then as much detail as possible should be included in the main  
755           paper.  
756           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
757           labor should be paid at least the minimum wage in the country of the data collector.

758           **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

759           Question: Does the paper describe potential risks incurred by study participants, whether such  
760           risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an  
761           equivalent approval/review based on the requirements of your country or institution) were obtained?

762           Answer: [Yes]

763           Justification: The tongue image collection process was approved by the Institutional Review Board  
764           (IRB).

765           Guidelines:

- 766           • The answer NA means that the paper does not involve crowdsourcing nor research with human  
767           subjects.  
768           • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
769           required for any human subjects research. If you obtained IRB approval, you should clearly state  
770           this in the paper.  
771           • We recognize that the procedures for this may vary significantly between institutions and  
772           locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for  
773           their institution.  
774           • For initial submissions, do not include any information that would break anonymity (if applica-  
775           ble), such as the institution conducting the review.

776           **16. Declaration of LLM usage**

777           Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard  
778           component of the core methods in this research? Note that if the LLM is used only for writing,  
779           editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or  
780           originality of the research, declaration is not required.

781           Answer: [Yes]

782           Justification: The detailed evaluation procedure of the large language models is described in **Section 5**  
783           and **Appendix F**.

784           Guidelines:

- 785           • The answer NA means that the core method development in this research does not involve LLMs  
786           as any important, original, or non-standard components.  
787           • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what  
788           should or should not be described.