

1 A Availability

- 2 You can access TCM-Ladder website at <http://tcmladder.com>
- 3 The GitHub repository with evaluation code and prompts is available at <https://github.com/orangeshushu/TCM-Ladder>
- 5 Data can be downloaded from <https://huggingface.co/datasets/timzzyus/TCM-Ladder>

7 B Data resources

8 As shown in Table 1 and Table 2, we incorporated several publicly available resources from existing
9 publications to extend and supplement the voice and pulse diagnosis datasets. For each dataset, we
10 provide the corresponding original publication reference and licensing documentation. In addition,
11 we partially incorporated several publicly available long-form dialogue datasets to supplement our
12 corpus. The sources and licensing statements of these datasets are provided in Table 3.

Table 1: Summary of publicly available voice and disease audio datasets

Dataset	Modality	Language	Task Type	Number of Samples	Year	License
COUGHVID[1] Coswara[2]	Cough audio Cough, breath, speech	English English	COVID-19 detection COVID-19 detection	20,000 5,000	2020 2022	CC-BY 4.0 CC-BY 4.0
UK COVID-19 Vocal Audio Dataset[3]	Cough, breath, speech	English	COVID-19 detection	70,000	2023	OGL v3.0
Respiratory Sound Database[4]	Lung auscultation sounds	English	Respiratory disease classification	920	2017	CC-BY 4.0
smarty4covid[5] Bridge2AI-Voice[6]	Cough, breath, voice Voice recordings	English English	COVID-19 detection Voice biomarker research	4,600 Not specified	2023 2025	CC-BY 4.0 Apache-2.0
VOICED[7]	Voice recordings	English	Pathological voice analysis	208	2018	ODC-BY 1.0
Perceptual Voice Qualities Dataset[8]	Voice recordings	English	Perpetual voice quality	360+	2020	CC-BY 4.0
COVID-19 Voice Dataset[9]	Voice recordings	English	COVID-19 detection	Not specified	2023	CC-BY 4.0
ALS IAC Speech Corpus[10]	Speech	English	ALS	Not specified	2024	CC-BY 4.0
PMC COVID-19 Voice Dataset[11]	Voice recordings	English	COVID-19 detection	Not specified	2022	OGL v3.0

Table 2: Summary of publicly available pulse datasets

Dataset	Modality	Language	Task Type	Number of Samples	Year	License
PulseDB[12]	ECG, PPG, ABP waveforms	English	Cuff-less blood pressure estimation	5,245,454	2023	ODbL
MIMIC-BP[13]	ECG, PPG, ABP waveforms	English	Blood pressure estimation	12,000	2024	ODC-BY 1.0
Pulse-ECG[14]	ECG images	English	ECG interpretation	1,160,000	2023	Apache-2.0
MTHS Dataset[15]	Video-PPG, ECG signals	English	Heart rate and SpO2 estimation	65	2023	CC BY-NC-ND 4.0
Weltory Dataset[16]	Video-PPG, ECG signals	English	Heart rate variability analysis	21	2023	CC BY-NC-ND 4.0
BUT-PPG Dataset[17]	Video-PPG signals	English	Heart rate estimation	65	2023	CC-BY 4.0

Table 3: Summary of available long-form TCM dialogue dataset

Dataset	Modality	Language	Task Type	Number of Samples	Year	License
Huatuo-26M[18]	Text	Chinese	QA, Dialogue	26M+	2023	CC-BY 4.0
TCMD[19]	Text	Chinese	Syndrome-Finding Mapping	100,000+	2024	CC-BY 4.0
CMD[20]	Text	Chinese	Medical Dialogue	25,000+ dialogues	2020	MIT

13 C Training details

14 **For Bencao:** *Bencao* is a model fine-tuned on the basis of a large language model, GPT (Generative
15 Pre-trained Transformer)[21]. The fine-tuning process leveraged a curated corpus of over 700 TCM
16 classical books as its knowledge base. Reinforcement Learning from Human Feedback (RLHF)[22]
17 was employed to iteratively improve the model’s output quality, with human annotators evaluating
18 the helpfulness and accuracy of its responses.

19 The model was initially guided by role-defining prompts. For example, one representative instruction
20 states: “*You are an experienced TCM expert. Please respond to users’ questions based on your*
21 *expertise and the content in your knowledge base. Try to combine traditional TCM terminology*

22 *with modern language to ensure both professionalism and clarity. Do not provide any diagnostic
23 conclusions, and explicitly mention the limitations of your responses."*

24 To refine the model’s output, we conducted more than 200 rounds of multi-turn instructional fine-
25 tuning. This iterative process aimed to ensure that the model’s responses were as professional,
26 objective, and informative as possible.

27 It is important to note that the training corpus consisted primarily of authoritative classical TCM
28 literature and educational textbooks. Some associated books also included tongue images and herbal
29 illustrations. However, no question-answer pairs from the TCM-Ladder benchmark were used in any
30 stage of the training or fine-tuning process.

31 **For Ladder-base:** The Ladder-base model is fine-tuned on the *Qwen2.5-7B* model using Group Relative
32 Policy Optimization(GRPO)[23], with the train set in TCM-Ladder. A consistent system prompt
33 is used for all questions: "*You are a helpful AI Assistant that provides well-reasoned and detailed
34 responses. You first think about the reasoning process as an internal monologue and then provide
35 the user with the answer. The response cannot be more than 400 words. Each question has only one
36 answer, A-E. Respond in the following format: <think>\n...</think>\n<answer>\n...Answer:A-
37 E...\n</answer>. Do not respond with another tag. Do not repeat the response.*". The final response
38 is parsed and evaluated using a rule-based reward system that assigns one point each for a correct
39 answer, proper format, and correct tags. These rewards are weighted in a ratio of 5:1:1 to encourage
40 the model to prioritize answering questions correctly.

41 For a question–answer pair (q, a) , the policy model $\pi_{\theta_{\text{old}}}$ samples a group of G responses $\{o_i\}_{i=1}^G$.
42 The advantage of the i -th response is calculated by normalizing the group-level reward $\{R_i\}_{i=1}^G$

$$\hat{A}_{l,t} = \frac{r_i - \text{mean}(\{R_i\}_{i=1}^G)}{\text{std}(\{R_i\}_{i=1}^G)} \quad (1)$$

43 GRPO first computes the mean loss across each generated sequence and then averages these losses
44 over all sampled sequences. It employs a clipped objective function combined with an explicitly
45 applied KL divergence penalty term between policy model and reference model:

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\theta) &= \mathbb{E} \left[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O | q) \right] \\ &\cdot \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[\frac{\pi_\theta(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})} \hat{A}_{i,t}, \right. \right. \\ &\quad \left. \left. \text{clip} \left(\frac{\pi_\theta(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta \| \pi_{\text{ref}}] \right\} \quad (2) \end{aligned}$$

46 D Evaluation environment

47 For general domain large language models such as *GPT-4*, *GPT-4o mini*, *Gemini 2.0 Flash*, *Gemini
48 2.5 Pro*, and *Deepseek*, we conducted evaluations using their respective APIs. The corresponding
49 code is available in our GitHub repository. For TCM-specific models such as *HuatuoGPT*[24],
50 *Zhongjing*[25], and *Bentsao*[26], we performed local evaluations, and the testing code for these
51 models is also provided in the GitHub repository to facilitate reproducibility.

52 All experiments were conducted on the Hellbender[27] computing cluster at the University of
53 Missouri. Detailed configurations of the computational resources are provided in Table 4.

54 E Safeguards

55 All released images underwent thorough manual review, and tongue images were subjected to strict
56 de-identification procedures to ensure that no personally identifiable information is included. Use of
57 our dataset requires agreement to the terms of use provided.

58 When using our released models, users are provided with appropriate risk disclaimers. For example:
59 the responses generated by this model are for informational purposes only and should not be con-

Table 4: Detailed configuration of computational resources

Model	Dell R7 40xa
Nodes	17
Cores/node	64
System Memory	238 G
GPU	A100
GPU Memory	80 GB
Numbers of GPUs	4
Local Scratch	1.6 TB

60 sidered a substitute for professional medical diagnosis. If you experience any discomfort or health
 61 concerns, please seek medical attention promptly.

62 **F Declaration of LLM usage**

63 We employed LLMs as a core component of our evaluation framework. Specifically, both general-
 64 domain and TCM-specific LLMs were systematically evaluated on *TCM-Ladder*, a multimodal
 65 benchmark designed for TCM. For multiple-choice questions, we used standardized zero-shot
 66 prompts, for example: “*Please answer the following multiple-choice question and clearly indicate*
 67 *the letter (A, B, C, D, or E) of the option you believe is correct.*” We provide implementation code
 68 for each evaluated LLM to ensure transparency and reproducibility (see GitHub links in **Appendix**
 69 **A**). Furthermore, LLMs are integral to our proposed textual evaluation metric, *Ladder-Score*, where
 70 LLMs assess model-generated answers based on well-defined rubrics. These rubrics span multiple
 71 dimensions, including logical consistency, semantic accuracy, knowledge coverage, and fluency of
 72 expression. Detailed evaluation protocols and rubric definitions are described in **Appendix H**.

73 **G Herb and tongue image data collection**

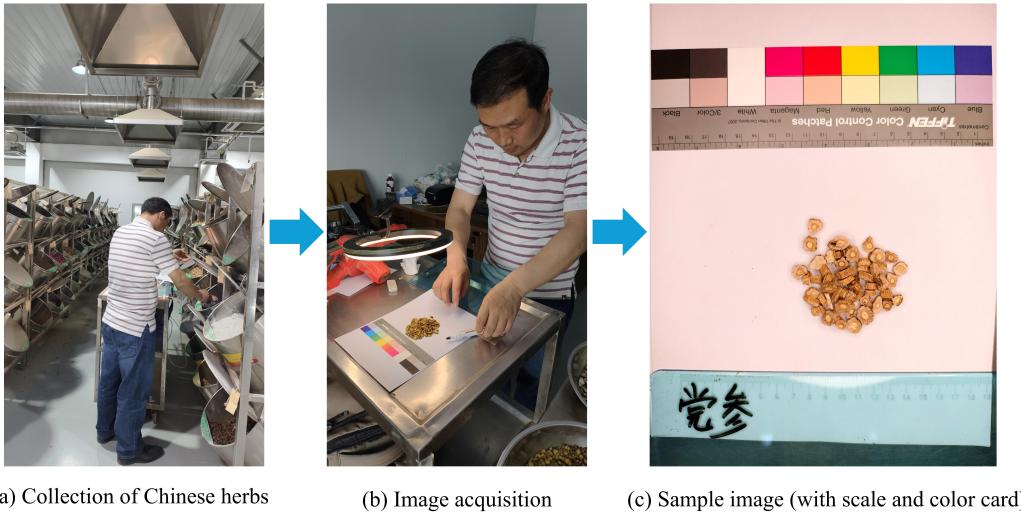
74 The images of Chinese medicinal herbs were sourced from two primary channels: a portion was
 75 collected from publicly accessible online resources, while the remainder was obtained through on-site
 76 photography conducted at *Shanghai Kangqiao Pharmaceutical Co., Ltd*[28]. As illustrated in Figure
 77 1, we first identified each herb within the pharmaceutical facility, followed by standardized image
 78 capture under controlled lighting conditions. Additionally, a measurement scale and a standard color
 79 chart were placed alongside each herb specimen to facilitate subsequent analytical processing and
 80 color correction.

81 Tongue image acquisition was conducted through two primary methods. As shown in Figure 2, a
 82 portion of the images was collected using a specialized tongue imaging device[29], which provides
 83 a stable lighting environment to ensure image consistency. Another portion was obtained via our
 84 mobile application, *iTongue*[30] (as shown in Figure 3), which enables analysis and prediction of
 85 TCM body constitution. To maximize the protection of participants’ privacy, only the tongue surface
 86 region of the images has been released publicly.

87 We use the existing image labels as ground truth answers and select three incorrect options as
 88 distractors. Based on specific question templates such as “*Which of the following images shows the*
 89 *herb [Herb Name]?*”, we generate corresponding visual questions. The code for generating these
 90 visual questions is publicly available on GitHub.

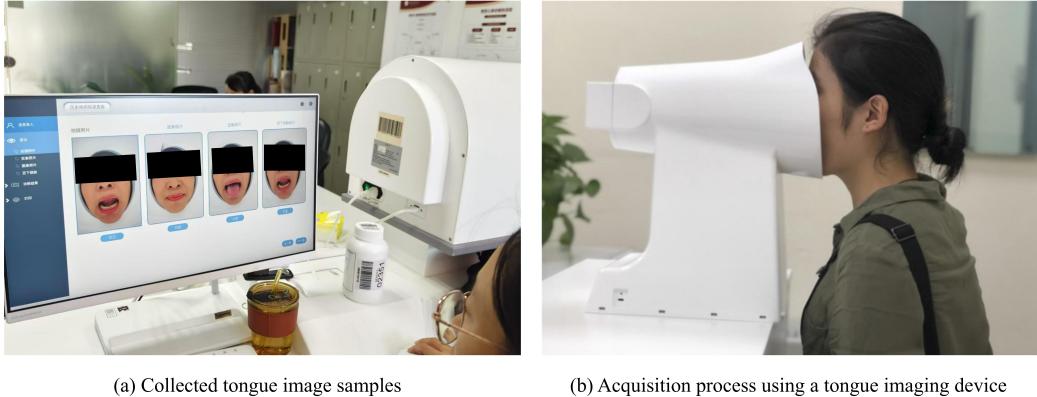
91 **H Ladder-Score: A hybrid evaluation metric for long-form TCM question 92 Answering**

93 To better evaluate the quality of long-form answers generated by TCM-specific LLMs in TCM
 94 settings, we introduce *Ladder-Score*, a hybrid evaluation metric that combines terminology matching
 95 with rubric-based semantic assessment.



(a) Collection of Chinese herbs (b) Image acquisition (c) Sample image (with scale and color card)

Figure 1: The image acquisition process of Chinese medicinal herbs. (a) Collection of physical herb specimens from a pharmaceutical factory. (b) Image capture under controlled lighting conditions. Photographs of the herbs were taken to ensure consistency and quality. (c) A sample herb image, including a scale bar and a color calibration card, which facilitates subsequent processing and color correction.



(a) Collected tongue image samples (b) Acquisition process using a tongue imaging device

Figure 2: Tongue image acquisition process using a tongue imaging device. (a) Sample tongue images collected from participants. To protect subject privacy, the eye regions have been masked. (b) Demonstration of the tongue image capture procedure using the imaging device.

96 H.1 Terminology-based score (*TermScore*)

97 We construct a domain-specific term list using *TCMBank*[31], international TCM terminology
 98 standards[32], and manually curated clinical terms (e.g., syndromes, symptoms, treatment principles,
 99 herbal formulas). Given a reference answer R and a generated candidate answer C , we extract their
 100 term sets T_R and T_C via Jieba-based tokenization and lexicon matching.

101 We compute:

- 102 • **Term Coverage:** the proportion of reference terms present in the candidate:

$$TC = \frac{|T_C \cap T_R|}{|T_R|} \quad (3)$$

- 103 • **Term Consistency:** the semantic similarity between term sets using sentence-level BERT
 104 embeddings and cosine similarity.

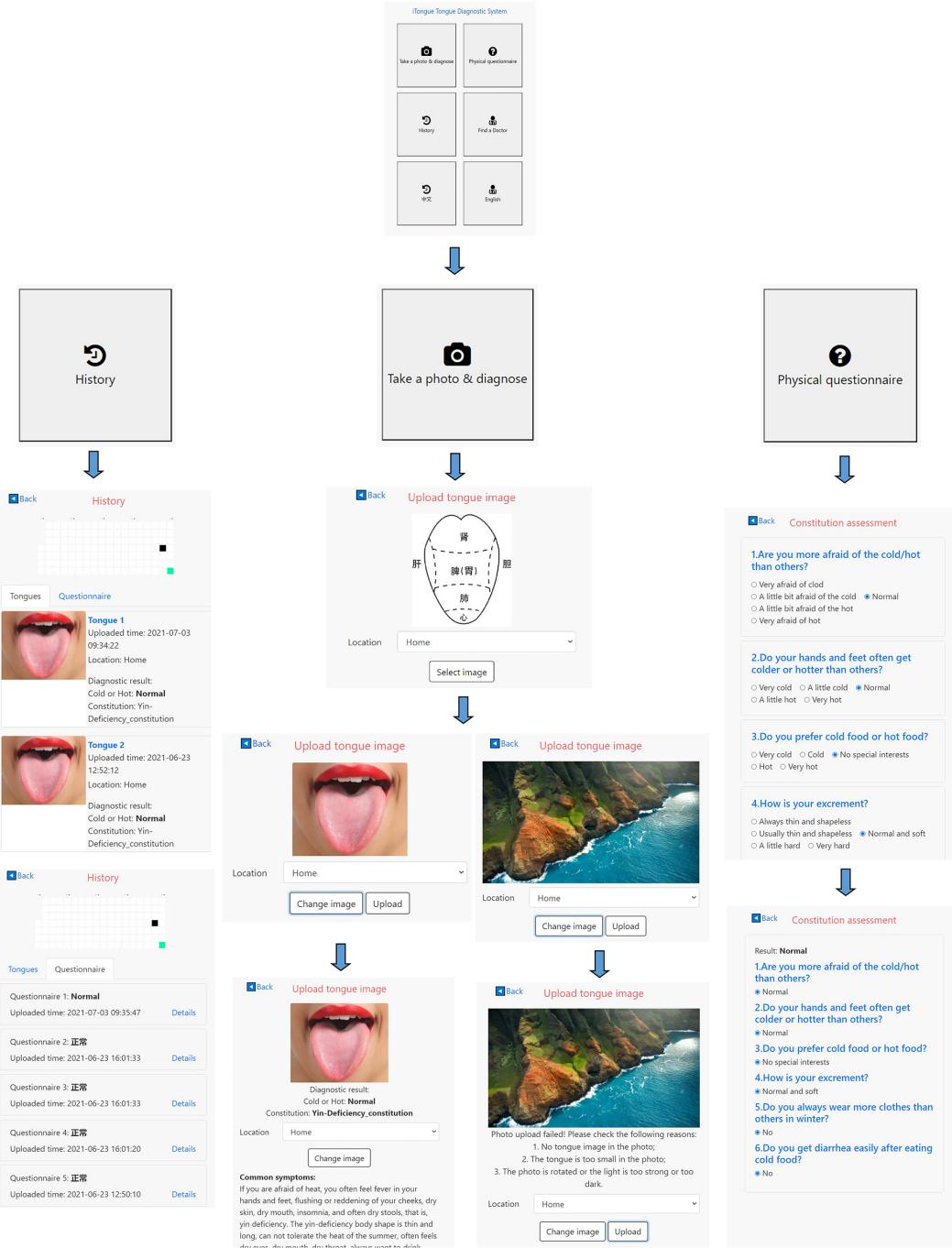


Figure 3: User interaction with *iTongue* app. Users can upload their own tongue images to receive TCM-based assessments, including cold–heat syndrome differentiation and body constitution classification. If a non-tongue image is uploaded, the application prompts the user to re-upload a valid tongue image.

105 The overall *TermScore* is:

$$TermScore = \lambda \cdot TC + (1 - \lambda) \cdot TS \quad (\lambda = 0.5) \quad (4)$$

106 **H.2 Semantic rubric score (*SemanticScore*)**

107 To assess semantic quality, we prompt an external LLM (e.g., GPT-4 or Claude) using a rubric
108 covering four dimensions as shown in Table 5.

Table 5: The four dimensions of the LLM scoring rubric.

Dimension	Description
Diagnostic Accuracy	Whether a correct syndrome or pathology is inferred
Treatment Appropriateness	Whether the recommendation aligns with clinical reasoning
Logical Coherence	Whether the answer is structured and internally consistent
Medical Language Quality	Clarity, fluency, and use of professional language

109 Each is rated on a 0–5 scale. The *SemanticScore* is the normalized average:

$$SemanticScore = \frac{1}{20} \sum_{i=1}^4 s_i \quad (5)$$

110 We apply blind prompting to prevent model self-bias, and optionally average across multiple LLM
111 reviewers.

112 **H.3 Final metric: Ladder-Score**

113 The final score is a weighted sum of both components:

$$\text{Ladder-Score} = \alpha \cdot \text{TermScore} + \beta \cdot \text{SemanticScore} \quad (6)$$

114 This formulation ensures both domain fidelity and semantic quality are captured in evaluating
115 long-form generative answers under TCM tasks.

116 **H.4 LLM-Based evaluation prompt template**

117 To elicit structured ratings from LLMs, we use the following prompt template for each generated
118 answer:

119 “*You are a TCM expert evaluating a model-generated answer. Please score the answer according to*
120 *the following rubric, from 0 (very poor) to 5 (excellent):*

- 121 1. *Diagnostic Accuracy: Is the TCM syndrome or pathology properly inferred?*
- 122 2. *Treatment Appropriateness: Are the therapeutic suggestions clinically relevant and reasonable?*
- 123 3. *Logical Coherence: Is the explanation consistent and well-structured?*
- 124 4. *Medical Language Quality: Is the response clearly written with accurate terminology?*

125 *Here is the candidate answer:*

126 *[candidate answer here]*

127 *Return your scores in the following format:*

128 *Diagnostic Accuracy: [0–5]*

129 *Treatment Appropriateness: [0–5]*

130 *Logical Coherence: [0–5]*

131 *Medical Language Quality: [0–5]”*

132 **H.5 Terminology extraction strategy**

133 To compute the *TermScore*, we extract domain-specific terms from both reference and candidate
134 answers using the following procedure:

135 **Lexicon construction:**

136 We build a term dictionary combining:

- 137 • The *TCMBank*[31] herb/syndrome/formula corpus
138 • Standardized clinical terminology lists
139 • Manual augmentation with high-frequency terms from training corpora

140 **Tokenization and filtering:**

- 141 • Use *jieba*[33] for Chinese segmentation.
142 • Retain terms satisfying:
143 – Length ≥ 2
144 – Occurrence in the constructed lexicon or inclusion of suffixes.

145 **Semantic similarity (for Term Consistency):**

- 146 • Compute sentence/term embeddings using pretrained Chinese BERT[34].
147 • Measure cosine similarity between aggregated term vectors from references and candidates.
148 • Use thresholded or averaged score as term semantic match score.

149 **H.6 Case study**

150 To assess the effectiveness of the proposed *Ladder-Score*, we compared its performance against
151 commonly used automatic evaluation metrics including BLEU-4[35], ROUGE-L[36], METEOR[37],
152 and BERTScore[38]. While these traditional metrics primarily rely on surface-level lexical overlap or
153 embedding similarity, they fail to capture the nuanced reasoning, terminology accuracy, and domain
154 alignment that are critical in TCM long-form question answering.

155 For instance, for the question (English translation):“A patient presents with a red tongue, yellow
156 coating, bitter taste in the mouth, restlessness, and dark yellow urine. What is your diagnosis and
157 treatment plan according to TCM?”

158 Reference answer (Expert-written) is:“Based on the symptoms of red tongue with yellow coating,
159 bitter taste, and restlessness, the pattern can be identified as Liver-Gallbladder Damp-Heat. The
160 treatment principle is to clear heat and drain dampness while soothing the Liver. The classical
161 formula Long Dan Xie Gan Tang is recommended, with herbs like Gentiana, Scutellaria, and
162 Gardenia. Patients should avoid greasy, spicy food to prevent worsening of symptoms.”

163 A high-quality candidate answer A generated by a large language model is:“The clinical presentation
164 suggests Liver and Gallbladder damp-heat syndrome, often due to emotional stagnation and internal
165 dampness. Treatment should focus on clearing heat, resolving dampness, and soothing the Liver.
166 Long Dan Xie Gan Tang is appropriate, including herbs like Gentiana, Alisma, Plantago Seed, and
167 Scutellaria. In cases with headaches or pronounced bitterness in the mouth, Bupleurum and Mint can
168 be added.”

169 Another candidate answer B, generated by a different model, exhibits high lexical overlap but lacks
170 substantive content:“Red tongue, yellow coating, and bitter taste indicate Liver-Gallbladder damp-
171 heat. You can use Long Dan Xie Gan Tang for treatment. This includes herbs like Gentiana, Gardenia,
172 and Scutellaria. The patient should avoid spicy food.”

173 We computed the evaluation scores of the above responses using different metrics, as shown in Table 6.
174 In this diagnostic task, the two model-generated responses may share a significant amount of surface-
175 level terminology (e.g., “Long Dan Xie Gan Tang,” “bitter taste,” “avoid spicy food”). However,
176 one response merely paraphrases the reference answer, while the other provides a precise pattern
177 differentiation, pharmacological rationale, and therapeutic strategy. Despite the substantial difference

178 in clinical quality, metrics such as BLEU and ROUGE may yield similar scores due to their reliance
 179 on lexical overlap. In contrast, *Ladder-Score* demonstrates a significant advantage. By integrating
 180 terminology matching with a rubric-based semantic evaluation mechanism, it can accurately capture
 181 the following aspects: (1) the accuracy of pattern differentiation; (2) the appropriateness of formula
 182 and herb recommendations; (3) the normative use of TCM-specific terminology; (4) the presence
 183 of templated or generic content; and (5) the ability to robustly assess complex clinical responses in
 184 long-form text.

185 This dual evaluation mechanism enables *Ladder-Score* to more reliably reflect the practical utility
 186 of model outputs in professional clinical settings. It is particularly well suited for specialized tasks
 187 that demand not only medical accuracy but also reasoning capability and contextual alignment,
 188 capabilities that are often lacking in conventional evaluation metrics.

Table 6: Performance of the two candidate responses under different evaluation metrics.

Metric	Candidate A	Candidate B	Difference	Can it distinguish quality?	Remarks
BLEU-4	0.45	0.43	+0.02	No	Rewards word overlap; ignores reasoning depth
ROUGE-L	0.52	0.51	+0.01	No	Captures sequence overlap; surface-level only
METEOR	0.44	0.45	-0.01	No	Slightly penalizes varied but correct expressions
BERTScore	0.87	0.85	+0.02	No	Sensitive to semantics, but lacks clinical awareness
Ladder-Score	0.91	0.75	+0.16	Yes	Correctly rewards richer clinical reasoning and term use

189 I Question design protocol and selection

190 When manually composing questions, TCM practitioners adhered to the criteria outlined in Table 7.
 191 The question stems were constructed based on content from official TCM textbooks, the *National*
 192 *Medical Licensing Examination syllabus*, or relevant online knowledge bases. During the question
 193 development process, care was taken to ensure that the stems were concise, unambiguous, and clearly
 194 stated. The questions were designed to cover multiple cognitive levels, including factual recall,
 195 conceptual understanding, and clinical application, while maintaining a diverse range of difficulty
 196 levels. Upon completion, the questions were further reviewed and filtered according to the screening
 197 criteria specified in Table 8 to eliminate redundancy and ensure quality.

Table 7: Guidelines and standards for constructing questions.

Category	Key points
Source and reference	<ul style="list-style-type: none"> - Base questions on authoritative sources such as official TCM textbooks, national curricula, and licensing examination syllabi. - Ensure coverage of fundamental TCM subjects, including fundamentals, diagnostics, herbal formulas, and pediatrics. - Utilize reputable databases for reference, such as <i>TCM Bank</i>[31], <i>ETCM</i>[39], <i>HERB 2.0</i>[40], etc.
Question stem design	<ul style="list-style-type: none"> - Ensure clarity, precision, and absence of ambiguity in wording. - Focus each question on a single knowledge point, avoiding irrelevant information. - Employ diverse question types to avoid formulaic stems.
Option construction	<ul style="list-style-type: none"> - Ensure only one correct answer for single-choice items; provide clearly defined correct combinations for multiple-choice items. - Include plausible distractors that test the depth of understanding. - Avoid leading cues such as key terms from the stem repeated in the correct answer.
Difficulty level	<ul style="list-style-type: none"> - Include questions across a range of difficulties. - Encourage the use of comprehension and application-based questions.
Common pitfalls to avoid	<ul style="list-style-type: none"> - Avoid content that is beyond the TCM curriculum scope or lacks clear relevance. - Prioritize clinical reasoning and critical analysis. - Ensure grammatical consistency between stem and options. - Avoid ambiguous, logically flawed, or duplicate questions.
Review and validation	<ul style="list-style-type: none"> - All items should undergo peer review by qualified TCM educators or clinicians. - Employ automated tools to detect duplications and evaluate ambiguity.

198 In addition to manual duplication checks, we employed three duplicate detection methods: string edit
 199 distance[41], TF-IDF[42, 43] with cosine similarity, and BERT [44, 45] based semantic encoding, to
 200 achieve multi level duplicate detection from surface to semantic levels. We considered two candidate
 201 questions to be duplicates if their similarity scores exceeded a threshold of 0.9 in any of the three
 202 methods mentioned above.

Table 8: Criteria for the question selection.

Dimension	Detailed explanation
Duplicate question	If two or more questions within the same category assess the same knowledge point with only slight variations in phrasing, retain only the version with clearer and more complete wording.
Ambiguous question	Any question with ambiguous wording, unclear intent, poor grammar, or syntactic confusion should be considered invalid and removed.
Answer uncertainty handling	If a single-choice question has multiple plausible answers or lacks clear differentiation between correct and distractor options, it should be removed. Similarly, multi-choice questions without a clear key or option logic should be discarded.
Option structure issues	Questions should be removed or revised if options are repetitive, nested within each other, or display length bias (e.g., one option is significantly longer than others).
Format and convention check	Questions missing prompts, unnumbered options, or nonstandard language/punctuation should be revised. If unfixable, they should be eliminated.
Outdated question context	Questions that conflict with the latest clinical guidelines or research, such as outdated treatment plans, should be revised or removed.

203 J Limitations and societal impact

204 TCM-Ladder serves as an evaluation benchmark for both general domain LLMs and TCM-specific
 205 models, offering a multimodal dataset to support future training and evaluation efforts in the field
 206 of TCM. However, several limitations remain and warrant further improvement and refinement in
 207 subsequent work.

208 **1. Limitations in multimodal evaluation.** The current models’ multimodal capabilities have been
 209 evaluated only on visual tasks involving herb images and tongue diagnosis images. Other data
 210 modalities that are also crucial in TCM, such as pulse diagnosis signals, audio (e.g., patient voice),
 211 and video, have not yet been incorporated into the evaluation framework in this study.

212 **2. Semantic deviations in cross-lingual translation.** During the construction of the multilingual
 213 version of the dataset, some Chinese questions may have incurred semantic deviations or ambiguities
 214 during the translation into English. Given the context-dependent and culturally embedded nature of
 215 TCM terminology, some expressions do not have direct equivalents in English, which can impair
 216 the model’s comprehension and reasoning capabilities. Such translation errors are particularly
 217 pronounced in complex dialectical reasoning tasks and may reduce the reliability of model evaluation.
 218 Future iterations of the dataset should involve closer collaboration between linguistic experts and
 219 TCM practitioners to ensure higher consistency and fidelity in cross-lingual representation.

220 **3. Imbalanced question categories.** The current version of the TCM-Ladder dataset exhibits
 221 noticeable imbalance across task types and knowledge subdomains. For instance, basic theoretical
 222 questions are overrepresented, while more complex or specialized areas such as clinical case reasoning,
 223 tongue image interpretation, and herbal property analysis are underrepresented. This imbalance
 224 may cause the models to overfit to high-frequency categories, leading to suboptimal generalization
 225 in low-resource tasks and ultimately affecting the fairness and representativeness of the overall
 226 evaluation. Future versions of the dataset should aim for more balanced coverage across task types
 227 and knowledge dimensions to ensure more comprehensive and challenging model assessment.

228 Despite the aforementioned limitations, TCM-Ladder, as the first multimodal dataset in the field of
 229 TCM, provides a foundational resource for the training and evaluation of multimodal TCM models.
 230 Evaluation results indicate that current TCM-specific large language models still exhibit substantial
 231 room for improvement and may lead to the following societal impacts.

232 **1. Facilitating the digitization and global dissemination of TCM.** By systematically organizing
 233 textual and visual modalities of TCM, the dataset contributes to the digitization of TCM knowledge.

234 **2. Enhancing education and clinical decision support.** This dataset provides a foundational
 235 resource for building intelligent TCM educational tools and diagnostic support systems. It may help
 236 address challenges such as the shortage of trained TCM professionals and the heavy reliance on
 237 experiential knowledge, ultimately supporting more standardized and scalable applications in clinical
 238 and educational settings.

239 **3. Privacy and fairness considerations.** Visual data such as tongue and herb images raise privacy
 240 concerns and require informed consent and anonymization. Additionally, the dataset may exhibit
 241 demographic imbalances (e.g., in gender, age, or region), which could introduce bias into trained

242 models. Future versions should incorporate fairness-aware data collection and annotation practices to
243 mitigate these risks and support equitable downstream applications.

244 References

- 245 [1] Lara Orlandic, Tomas Teijeiro, and David Atienza. The coughvid crowdsourcing dataset, a corpus for the
246 study of large-scale cough analysis algorithms. *Scientific Data*, 8(1):156, 2021.
- 247 [2] Debarpan Bhattacharya, Neeraj Kumar Sharma, Debottam Dutta, Srikanth Raj Chetupalli, Pravin Mote,
248 Sriram Ganapathy, C Chandrakiran, Sahiti Nori, KK Suhail, Sadhana Gonuguntla, et al. Coswara: A
249 respiratory sounds and symptoms dataset for remote screening of sars-cov-2 infection. *Scientific data*,
250 10(1):397, 2023.
- 251 [3] Harry Coppock, George Nicholson, Ivan Kiskin, Vasiliki Koutra, Kieran Baker, Jobie Budd, Richard
252 Payne, Emma Karoune, David Hurley, Alexander Titcomb, Sabrina Egglestone, Ana Tendero Cañadas,
253 Lorraine Butler, Radka Jersakova, Jonathon Mellor, Selina Patel, Tracey Thornley, Peter Diggle, Sylvia
254 Richardson, Josef Packham, Björn W. Schuller, Davide Pigoli, Steven Gilmour, Stephen Roberts, and
255 Chris Holmes. Audio-based ai classifiers show no evidence of improved covid-19 screening over simple
256 symptoms checkers. *Nature Machine Intelligence*, 2024.
- 257 [4] BM Rocha, Dimitris Filos, Lea Mendes, Ioannis Vogiatzis, Eleni Perantoni, Evangelos Kaimakamis,
258 P Natsiavas, Ana Oliveira, C Jácome, A Marques, et al. A respiratory sound database for the development
259 of automated classification. In *Precision Medicine Powered by pHealth and Connected Health: ICBHI*
260 2017, Thessaloniki, Greece, 18-21 November 2017, pages 33–37. Springer, 2018.
- 261 [5] Konstantia Zarkogianni, Edmund Dervakos, George Filandrianos, Theofanis Ganitidis, Vasiliki Gkatzou,
262 Aikaterini Sakagianni, Raghu Raghavendra, CL Nikias, Giorgos Stamou, and Konstantina S Nikita. The
263 smarty4covid dataset and knowledge base: a framework enabling interpretable analysis of audio signals.
264 *arXiv preprint arXiv:2307.05096*, 2023.
- 265 [6] Elijah Moothedan, Micah Boyer, Stephanie Watts, Yassmeen Abdel-Aty, Satrajit Ghosh, Anaïs Rameau,
266 Alexandros Sigaras, Olivier Elemento, Bridge2AI-Voice Consortium, and Yael Bensoussan. The bridge2ai-
267 voice application: initial feasibility study of voice data acquisition through mobile health. *Frontiers in*
268 *Digital Health*, 7:1514971, 2025.
- 269 [7] Ugo Cesari, Giuseppe De Pietro, Elio Marciano, Ciro Niri, Giovanna Sannino, and Laura Verde. A new
270 database of healthy and pathological voices. *Computers & Electrical Engineering*, 68:310–321, 2018.
- 271 [8] Patrick R Walden. Perceptual voice qualities database (pvqd): database characteristics. *Journal of Voice*,
272 36(6):875–e15, 2022.
- 273 [9] Jobie Budd, Kieran Baker, Emma Karoune, Harry Coppock, Selina Patel, Ana Tendero Cañadas, Alexander
274 Titcomb, Richard Payne, David Hurley, Sabrina Egglestone, et al. A large-scale and pcr-referenced vocal
275 audio dataset for covid-19. *arXiv preprint arXiv:2212.07738*, 2022.
- 276 [10] Raffaele Dubbioso, Myriam Spisto, Laura Verde, Valentina Virginia Iuzzolino, Gianmaria Senerchia, Elena
277 Salvatore, Giuseppe De Pietro, Ivano De Falco, and Giovanna Sannino. Voice signals database of als
278 patients with different dysarthria severity and healthy controls. *Scientific Data*, 11(1):800, 2024.
- 279 [11] Jobie Budd, Kieran Baker, Emma Karoune, Harry Coppock, Selina Patel, Richard Payne, Ana Ten-
280 dero Cañadas, Alexander Titcomb, David Hurley, Sabrina Egglestone, et al. A large-scale and pcr-
281 referenced vocal audio dataset for covid-19. *Scientific data*, 11(1):700, 2024.
- 282 [12] Weinan Wang, Pedram Mohseni, Kevin L Kilgore, and Laleh Najafizadeh. Pulsedb: A large, cleaned
283 dataset based on mimic-iii and vitaldb for benchmarking cuff-less blood pressure estimation methods.
284 *Frontiers in Digital Health*, 4:1090854, 2023.
- 285 [13] Ivandro Sanches, Victor V Gomes, Carlos Caetano, Lizeth SB Cabrera, Vinicius H Cene, Thomas Beltrame,
286 Wonkyu Lee, Sanghyun Baek, and Otávio AB Penatti. Mimic-bp: A curated dataset for blood pressure
287 estimation. *Scientific Data*, 11(1):1233, 2024.
- 288 [14] Ruoqi Liu, Yuelin Bai, Xiang Yue, and Ping Zhang. Teach multimodal llms to comprehend electrocardio-
289 graphic images. *arXiv preprint arXiv:2410.19008*, 2024.
- 290 [15] Taha Samavati, Mahdi Farvardin, and Aboozar Ghaffari. Efficient deep learning-based estimation of the
291 vital signs on smartphones. *arXiv preprint arXiv:2204.08989*, 2022.

- 292 [16] Ahsan Mehmood, Asma Sarouji, M Mahboob Ur Rahman, and Tareq Y Al-Naffouri. Your smartphone
 293 could act as a pulse-oximeter and as a single-lead ecg. *Scientific Reports*, 13(1):19277, 2023.
- 294 [17] Andrea Nemcova, Enikö Vargova, Radovan Smisek, Lucie Marsanova, Lukas Smital, and Martin Vitek.
 295 Brno university of technology smartphone ppg database (but ppg): Annotated dataset for ppg quality
 296 assessment and heart rate estimation. *BioMed Research International*, 2021(1):3453007, 2021.
- 297 [18] Jianquan Li, Xidong Wang, Xiangbo Wu, Zhiyi Zhang, Xiaolong Xu, Jie Fu, Prayag Tiwari, Xiang Wan,
 298 and Benyou Wang. Huatuo-26m, a large-scale chinese medical qa dataset. *arXiv preprint arXiv:2305.01526*,
 299 2023.
- 300 [19] Ping Yu, Kaitao Song, Fengchen He, Ming Chen, and Jianfeng Lu. Tcmd: A traditional chinese medicine
 301 qa dataset for evaluating large language models. *arXiv preprint arXiv:2406.04941*, 2024.
- 302 [20] Toyhom. Chinese-medical-dialogue-data. <https://github.com/Toyhom/Chinese-medical-dialogue-data>, 2025. Accessed: May 22, 2025.
- 304 [21] Gokul Yenduri, Gautam Srivastava, Praveen Kumar Reddy Maddikunta, Rutvij H Jhaveri, Weizheng
 305 Wang, Athanasios V Vasilakos, Thippa Reddy Gadekallu, et al. Generative pre-trained transformer: A
 306 comprehensive review on enabling technologies, potential applications, emerging challenges, and future
 307 directions. *arXiv preprint arXiv:2305.10435*, 2023.
- 308 [22] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
 309 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
 310 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 311 [23] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan
 312 Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language
 313 models. *arXiv preprint arXiv:2402.03300*, 2024.
- 314 [24] Junying Chen, Xidong Wang, Ke Ji, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang, Dingjie
 315 Song, Wenya Xie, Chuyi Kong, et al. Huatuogpt-ii, one-stage training for medical adaption of llms. *arXiv
 316 preprint arXiv:2311.09774*, 2023.
- 317 [25] Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan.
 318 Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback
 319 and real-world multi-turn dialogue. In *Proceedings of the AAAI conference on artificial intelligence*,
 320 volume 38, pages 19368–19376, 2024.
- 321 [26] Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: Tuning
 322 llama model with chinese medical knowledge. *arXiv preprint arXiv:2304.06975*, 2023.
- 323 [27] IT Research Support Solutions Wiki. Getting started with hpc. <https://docs.itrss.umsystem.edu/pub/hpc/start>, 2025. Accessed: May 22, 2025.
- 325 [28] Shanghai kangqiao pharmaceutical co., ltd. <http://www.shkangqiao.com/>, 2025. Accessed: May 22,
 326 2025.
- 327 [29] Jun Li, Pei Yuan, Xiaojuan Hu, Jingbin Huang, Longtao Cui, Ji Cui, Xuxiang Ma, Tao Jiang, Xinghua Yao,
 328 Jiacai Li, et al. A tongue features fusion approach to predicting prediabetes and diabetes with machine
 329 learning. *Journal of biomedical informatics*, 115:103693, 2021.
- 330 [30] Jiacheng Xie, Congcong Jing, Ziyang Zhang, Jiatuo Xu, Ye Duan, and Dong Xu. Digital tongue image
 331 analyses for health assessment. *Medical Review*, 1(2):172–198, 2021.
- 332 [31] Qiuje Lv, Guanxing Chen, Haohuai He, Ziduo Yang, Lu Zhao, Kang Zhang, and Calvin Yu-Chian Chen.
 333 Tcmbank-the largest tcm database provides deep learning-based chinese-western medicine exclusion
 334 prediction. *Signal Transduction and Targeted Therapy*, 8(1):127, 2023.
- 335 [32] World Health Organization et al. *WHO international standard terminologies on traditional Chinese
 336 medicine*. World Health Organization, 2022.
- 337 [33] Xianwei Zhang, Peng Wu, Jiuming Cai, and Kun Wang. A contrastive study of chinese text segmentation
 338 tools in marketing notification texts. In *Journal of Physics: Conference Series*, volume 1302, page 022010.
 339 IOP Publishing, 2019.
- 340 [34] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. Pre-training with whole word masking
 341 for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3504–3514,
 342 2021.

- 343 [35] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation
344 of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational*
345 *Linguistics*, pages 311–318, 2002.
- 346 [36] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches*
347 *out*, pages 74–81, 2004.
- 348 [37] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved
349 correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation*
350 *measures for machine translation and/or summarization*, pages 65–72, 2005.
- 351 [38] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text
352 generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- 353 [39] Hai-Yu Xu, Yan-Qiong Zhang, Zhen-Ming Liu, Tong Chen, Chuan-Yu Lv, Shi-Huan Tang, Xiao-Bo Zhang,
354 Wei Zhang, Zhi-Yong Li, Rong-Rong Zhou, et al. Etcm: an encyclopaedia of traditional chinese medicine.
355 *Nucleic acids research*, 47(D1):D976–D982, 2019.
- 356 [40] Kai Gao, Liu Liu, Shuangshuang Lei, Zhinong Li, Peipei Huo, Zhihao Wang, Lei Dong, Wenxin Deng,
357 Dechao Bu, Xiaoxi Zeng, et al. Herb 2.0: an updated database integrating clinical and experimental
358 evidence for traditional chinese medicine. *Nucleic Acids Research*, 53(D1):D1404–D1414, 2025.
- 359 [41] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals. In
360 *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union, 1966.
- 361 [42] Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information retrieval*,
362 volume 39. Cambridge University Press Cambridge, 2008.
- 363 [43] Gerard Salton, Anita Wong, and Chung-Shu Yang. A vector space model for automatic indexing. *Communi-*
364 *cations of the ACM*, 18(11):613–620, 1975.
- 365 [44] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidi-
366 rectional transformers for language understanding. In *Proceedings of the 2019 conference of the North*
367 *American chapter of the association for computational linguistics: human language technologies, volume*
368 *1 (long and short papers)*, pages 4171–4186, 2019.
- 369 [45] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks.
370 *arXiv preprint arXiv:1908.10084*, 2019.