Should I Believe in What Medical AI Says? A Chinese Benchmark for Medication Based on Knowledge and Reasoning

Anonymous ACL submission

Abstract

Large language models (LLMs) show potential in healthcare but often generate hallucinations, especially when handling unfamiliar information. In medication, a systematic benchmark to evaluate model capabilities is lacking, which is critical given the high-risk nature of medical information. This paper introduces a Chinese benchmark aimed at assessing models in medication tasks, focusing on knowledge and reasoning across six datasets: indication, dosage and administration, contraindicated population, mechanisms of action, drug recommendation, and drug interaction. We evaluate eight closedsource and five open-source models to identify knowledge boundaries, providing the first systematic analysis of limitations and risks in proprietary medical models.

1 Introduction

017

019

024

027

Large language models (LLMs) have made significant strides in various domains, including medication, where they provide information and recommendations related to medical treatments (Singhal et al., 2022; Nori et al., 2023). However, a significant challenge remains: these models are prone to generating hallucinations and confidently providing incorrect or incomplete information, especially in cases where they lack adequate knowledge (Stefansson and Johansson, 2021; Shukla et al., 2022). In the context of medication and drug usage, such hallucinations can lead to critical errors, particularly in high-risk situations like identifying contraindicated populations or recommending unsafe drug combinations. Despite the progress made in medical AI, there is a notable gap in the development of systematic benchmarks to evaluate the full range of a model's capabilities in medication applications.

> In this paper, we construct a Chinese benchmark called **ChiDrug**, specifically designed to assess LLMs' knowledge and reasoning abilities in the



Figure 1: Our benchmark involves four datasets that directly examine model parametric knowledge and two datasets that examine model reasoning ability.

medication domain. As shown in Figure 1, our benchmark is structured into two key subdimensions: **parametric knowledge** and **reasoning capability**. We construct six diverse datasets that cover crucial aspects of drug information—dosage and administration, indication, contraindicated populations, mechanisms of action, medication recommendations, and drug interactions. 041

042

043

044

045

047

051

053

055

057

060

061

062

063

To evaluate the capabilities of existing models, we apply our benchmark to eight closed-source and five open-source models. Our work also explores various methods for expressing knowledge boundaries, providing insights into the potential risks of overconfident but inaccurate AI-generated responses.

Our contributions include: (1) This benchmark serves as the first systematic tool for analyzing the capabilities of LLMs in the field of medicine across various dimensions. (2) We are pioneers in conducting knowledge boundary analysis on medical models within medicine, providing a comprehensive overview of their performance in real-world medical applications.

Related Work 2

065

069

071

077

081

090

101

102

103

104

106

107

108

110

111

112

Chinese Benchmark in Medication 2.1

Assessing the capabilities of Large Language Models (LLMs) in the medical field requires specialized benchmarks, especially when dealing with Chinese medical texts. Recent efforts have led to the development of several Chinese-specific medical benchmarks, focusing on various domains such as clinical question answering, knowledge recall, and medication recommendations (Singhal et al., 2023; Liu et al., 2024; Wang et al., 2024; Yue et al., 2024).

MedExpQA (Liu et al., 2024) proposes a multilingual benchmark evaluating models on medical question answering tasks, including drug-related and clinical guideline questions. DialMed (He et al.) focuses on dialogue-based medication recommendations, testing models on handling patient symptom queries and drug interactions. However, existing datasets do not have a dedicated benchmark built in the field of medication in Chinese to evaluate the model's ability in this area.

2.2 Abstention in LLMs

The ability of Large Language Models (LLMs) to refrain from providing answers when uncertain-is crucial for enhancing model reliability and safety. Studies have explored various methods to improve this capability (Wen et al., 2024):

Currently, methods to guide models in refusing to answer include: Calibration-Based Methods: After the model provides an answer, continue by asking, "Are you sure about your answer?" to verify its confidence (Tian et al.). Training-Based Methods: Construct a training set containing both questions the model can answer and those it cannot, training the model to refuse to answer questions with unfamiliar knowledge (Slobodkin et al., 2023; Zhang et al., 2023; Stengel-Eskin et al., 2024). Consistency-Based Methods: Perform multiple samplings and calculate the consistency score of the model's responses to assess reliability(Kuhn et al.; Feng et al., 2024). Token Probability Methods: Ensemble the probability of each token generated by the model to determine the uncertainty of the response (Liang et al., 2024; Malinin and Gales, 2021).

Dataset 3

ChiDrug is designed to assess models' parametric knowledge and reasoning ability in handling critical medication-related tasks. Below, we outline the dataset construction process and the verification procedures used to ensure the quality and reliability of the data. The entire benchmark construction process is shown in Figure 2

3.1 Dataset Construction

We began by collecting official drug brochures for existing medications from the internet¹. We organized this information into a table that includes details on 8,000 drugs, encompassing their generic names, ingredients, specifications, indications, dosages, contraindications, drug interactions, adverse reactions, and mechanisms of action. This structured dataset served as the foundation for developing questions that evaluate the model's parametric knowledge in four areas: Indication, Dosage and Administration, Contraindicated Population, and Mechanism of Action. We extracted the relevant sections from each drug brochure and utilized Spark² to generate multiplechoice question stems and answer options. In constructing these questions, we ensured that the incorrect options did not overlap with the correct ones.

The second step involved constructing questions for Medication Recommendation. We collected doctor-patient dialogues from the existing DIALMED dataset (He et al.), where Spark transformed these dialogues into question formats, using the recommended medication by the doctor as the correct option. To generate distractor options that could confuse the model, we first used Spark to extract the patient's symptoms and demographic information, then searched the drug brochures for medications that treat the same symptoms but are not suitable for the patient's demographic group, thereby creating incorrect options (e.g., "symptom in indication and demographic in contraindicated population").

In the third step, we constructed a dataset for Drug Interaction. First, doctors defined three risk levels for drug interaction (high, medium, and low). We then randomly selected a drug from the brochures and identified its combination guidelines. From there, we extracted the ingredients involved in drug interactions and further searched for medications that contained the same ingredients. Finally, we input the two drugs and the interaction documentation into Spark to generate the appropriate risk level as the correct answer.

113

114

117

118 119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

¹https://drugs.dxy.cn

²https://xinghuo.xfyun.cn



Figure 2: Overview of our benchmark construction process

3.2 Verification

162

163

164

165

166

167

168

170

171

172

174

175

176

177

178

179

181

185

189

190

191

Since we automatically generated the questions for the dataset, we used a double-check process to make sure the questions were reasonable. Each question was tested by three large models (GPT-4³, Qwen-max⁴, ERNIE bot⁵). We gave these models the question, options, and document sources and asked them to check the following: (1) If the question makes sense. (2) If the answer is correct. (3) If the answer is unique. A question was considered valid only if all three models agreed it was correct. Additionally, We hire doctors with licensed qualifications to examine all the datasets we construct.

In the end, we created a benchmark dataset with a total of 5,243 samples, covering the following categories: Indication (705), Dosage and Administration (651), Contraindicated Population (659), Mechanism of Action (773), Medication Recommendation (838), and Drug Interaction (1,617).

4 Experiment

In this section, we evaluate the performance of large language models (LLMs) on our benchmark. We assess both closed-source and open-source models, using our benchmark to examine their capabilities in handling medication-related queries and their ability to identify knowledge gaps and overconfidence. Table 1 presents the results for the model ability, while the second table focuses on the methods to express the knowledge boundaries in seven different methods.

⁴https://tongyi.aliyun.com/qianwen

4.1 Model Performance Evaluation

We selected models with strong Chinese language capabilities, including GPT40 (Hurst et al., 2024), Claude3.5-Sonnet⁶, Qwen-max⁷, Doubao⁸, GLM4 (GLM et al., 2024), Baichuan4⁹, XiaoYi¹⁰, and ERNIE Bot¹¹, for evaluation of closed-source models. For open-source models, we chose Bencao (Wang et al., 2023), MedGLM (Haochun Wang, 2023), MedicalGPT (Xu, 2023), ChiMedical (Tian et al., 2024), and HuatuoGPT2 (Chen et al., 2024) for evaluation.

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

The results summarized in Table 1, The closedsource models generally outperformed the opensource models across all dimensions, with XiaoYi leading in overall performance, followed closely by GPT40 and ERNIE Bot. In Open-source models, Bencao and MedicalGPT demonstrated lower performance, particularly in complex tasks like Contraindicated Populations and Drug Interactions, while HuatuoGPT2 generally outperformed other models. We will provide a more detailed analysis of each model in the Appendix A.

4.2 Methods for Knowledge Boundary Detection

In this subsection, we apply seven methods to explore their impact on expressing uncertainty or abstention, using HuatuoGPT2 as the backbone.

³https://chatgpt.com

⁵https://yiyan.baidu.com

⁶https://claude.ai

⁷https://tongyi.aliyun.com/qianwen

⁸https://www.doubao.com/chat

¹⁰https://chatdr.iflyhealth.com

¹¹https://yiyan.baidu.com

Close-source Models										
	Dosage and Administration	Indication	Contraindicated Population	Mechanism of Action	Medication Recommendation	Drug Interaction	Avg.			
XiaoYi	81.1	77.87	66.71	92.85	65.31	63.27	73.52			
GPT40	66.41	73.65	69.35	92.13	59.79	59.93	70.21			
ERNIE	67.64	65.3	57.97	92.76	51.43	38.59	62.28			
Qwen-max	69.02	72.13	68.19	<u>93.28</u>	61.22	54.73	<u>69.76</u>			
Doubao	71.32	71.24	54.17	92.77	63.25	55.35	68.02			
GLM4	71.32	75.71	71.02	94.16	59.79	54.92	71.15			
Claude3.5	54.59	74.53	70.29	89.92	54.06	<u>60.73</u>	67.24			
Baichuan4	62.14	69.97	69.24	90.35	52.98	52.81	66.25			
Open-source Models										
Bencao	28.92	<u>19.88</u>	12.2	40.71	16.23	38.28	26.04			
MedGLM	38.92	13.21	8.75	44.86	20.17	34.59	26.75			
MedicalGPT	<u>33.51</u>	10.14	3.18	<u>49.41</u>	13.84	30.98	23.51			
ChiMedical	<u>33.51</u>	16.04	<u>14.32</u>	38.54	<u>24.71</u>	<u>36.05</u>	27.20			
HuatuoGPT2	55.83	47.03	18.66	77.16	25.18	25.60	41.58			

Table 1: This table presents the performance of 8 closed-source models and 5 open-source models across various medication-related tasks. Bold indicates the best performance, while underlining denotes the second-best.

	Dosage an	d Administration	Indica	tion	Contraindicated Population		Mechanism of Action		Medication Recommendation		Drug Interaction		Avg.	
	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc
Baseline	55.83	55.83	47.03	47.03	18.66	18.66	77.16	77.16	25.18	25.18	25.60	25.60	41.58	41.58
Post-calibration	55.94	57.91	47.01	47.62	18.05	18.17	77.21	77.37	25.23	26.12	24.76	25.67	41.37	42.12
IDK	54.26	54.68	47.18	47.18	17.90	17.30	78.68	80.28	24.47	24.66	24.55	26.64	41.17	41.79
LNS	53.99	50.46	50.78	50.85	18.18	23.03	79.97	71.07	28.22	26.21	24.94	23.58	42.68	40.87
Probing	66.11	43.90	66.11	51.86	22.83	27.43	85.24	80.01	29.82	19.80	31.31	30.54	50.24	42.26
R-tuning	53.90	57.61	53.66	56.00	22.12	31.00	81.35	83.00	19.80	12.45	30.34	31.68	43.53	45.29
Self-Consistency	66.85	46.20	68.03	48.46	15.24	10.67	89.37	78.60	30.15	20.84	20.90	30.00	48.42	39.13
Semantic Entropy	64.79	55.32	70.28	52.71	28.11	26.95	86.24	83.09	29.22	24.95	38.56	31.76	52.87	45.79

Table 2: This table displays the performance of 7 different methods on the models' ability to detect knowledge boundaries and manage uncertainty.

Post-Calibration (Tian et al.): Enhances model confidence by prompting it to verbalize its certainty after providing an answer.

IDK (I Don't Know): Trains models to acknowledge uncertainty by explicitly stating when they lack knowledge, thereby reducing hallucinations.

LNS (Malinin and Gales, 2021): Utilizes probabilistic ensemble-based techniques to assess uncertainty in structured prediction tasks, aiding in more reliable outputs.

Probing (Slobodkin et al., 2023): Analyzes internal model representations to understand how they encode information about answerability, helping detect overconfidence and hallucinations.

R-tuning (Zhang et al., 2023): Instructs models to explicitly state when they lack knowledge, reducing the generation of hallucinated information.

Self-Consistency (Kuhn et al.): Enhances reasoning by generating multiple reasoning paths and selecting the most consistent answer, improving response reliability.

Semantic Entropy (Feng et al., 2024): Estimates uncertainty in natural language generation by considering linguistic invariances, allowing models to better assess the reliability of their outputs.

In this section, two evaluation metrics are used: **Precision**: This metric measures the proportion of correct answers out of the total predictions made, without abstaining. **Abstain Accuracy**: This metric evaluates when models correctly answer or choose not to respond due to uncertainty.

248

249

250

251

252

253

254

256

257

258

259

261

262

263

264

265

266

267

269

270

271

272

273

Results are shown in Table 2. Post-calibration and IDK cannot achieve good results in HuatuoGPT2 with weak instruction capabilities. Self-Consistency improved accuracy in complex tasks like Medication Recommendations. Probing refined uncertainty estimations, with varying effectiveness. R-tuning reduced hallucinations but sometimes sacrificed performance on hard tasks, while LNS showed mixed results, improving Medication Recommendations but hindering performance on Drug Interaction. Overall, Semantic Entropy has achieved good results in both metrics, and we further analyze the effectiveness of this method on multiple models in Appendix B.

5 Conclusion

We present ChiDrug, a benchmark designed to evaluate LLMs (Large Language Models) in medication-related tasks, with an emphasis on their knowledge and reasoning abilities. Both GLM4 and XiaoYi performed exceptionally well; however, even these advanced models exhibited gaps in drug knowledge. This highlights the need for effective methods to align the knowledge boundaries of LLMs, particularly for high-risk tasks. 274

275

2

28

281

- 28
- 28

28

- 290
- 291 292
- 293 294
- 295 296
- 290
- 2
- 300 301
- 302 303
- 304 305
- 306 307

3

3

312

- 313 314
- 315
- 310
- 318 319

320 321 322

- 323 324
- 32

325 326 6 Limitations

This study mainly focuses on Chinese medical texts, which may affect generalizability. The benchmark doesn't fully capture real-world medical decision-making complexities. Additionally, model generalization to new knowledge, handling uncertainty, and reliance on high-quality, up-todate data are ongoing challenges for AI in healthcare.

References

- Junying Chen, Xidong Wang, Ke Ji, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang, Dingjie Song, Wenya Xie, Chuyi Kong, Jianquan Li, Xiang Wan, Haizhou Li, and Benyou Wang. 2024. Huatuogpt-ii, one-stage training for medical adaption of llms. *Preprint*, arXiv:2311.09774.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. 2024. Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration. *arXiv preprint arXiv:2402.00367*.
- Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Jingyu Sun, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. Preprint, arXiv:2406.12793.
 - Sendong Zhao Bing Qin Ting Liu Haochun Wang, Chi Liu. 2023. Chatglm-med: chatglm. https: //github.com/SCIR-HI/Med-ChatGLM.
- Zhenfeng He, Yuqiang Han, Zhenqiu Ouyang, Wei Gao, Hongxu Chen, Guandong Xu, Jian Wu, Haochun Wang, Chi Liu, Nuwa Xi, et al. Dialmed: A dataset for dialogue-based medication recommendation.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In

The Eleventh International Conference on Learning Representations.

327

328

329

331

332

333

334

335

336

339

340

341

342

343

345

346

347

350

351

352

353

354

355

356

358

360

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

- Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Yi Wang, Zhonghao Wang, Feiyu Xiong, et al. 2024. Internal consistency and self-feedback in large language models: A survey. *arXiv preprint arXiv:2407.14507*.
- Mianxin Liu, Jinru Ding, Jie Xu, Weiguo Hu, Xiaoyang Li, Lifeng Zhu, Zhian Bai, Xiaoming Shi, Benyou Wang, Haitao Song, Pengfei Liu, Xiaofan Zhang, Shanshan Wang, Kang Li, Haofen Wang, Tong Ruan, Xuanjing Huang, Xin Sun, and Shaoting Zhang. 2024. Medbench: A comprehensive, standardized, and reliable benchmarking system for evaluating chinese medical large language models. *Preprint*, arXiv:2407.10990.
- Andrey Malinin and Mark Gales. 2021. Uncertainty estimation in autoregressive structured prediction. In International Conference on Learning Representations.
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. *Preprint*, arXiv:2303.13375.
- Abhay Shukla, Paheli Bhattacharya, Soham Poddar, Rajdeep Mukherjee, Kripabandhu Ghosh, Pawan Goyal, and Saptarshi Ghosh. 2022. Legal case document summarization: Extractive and abstractive methods and their evaluation. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1048–1064.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large language models encode clinical knowledge. *Preprint*, arXiv:2212.13138.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. 2023. Towards expert-level medical question answering with large language models. *Preprint*, arXiv:2305.09617.

- 398
- 399 400 401
- 402 403 404 405 406
- 407 408 409 410
- 411
- 412 413

414

- 415 416 417
- 418 419

420 421

422 423

424

425 426

- 427 428
- 429 430
- 431

432 433

434 435 436

- 437
- 438 439

440

- Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. 2023. The curious case of hallucinatory (un) answerability: Finding truths in the hidden states of over-confident large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 3607–3625.
- Elis Stefansson and Karl H. Johansson. 2021. Computing complexity-aware plans using kolmogorov complexity. Preprint, arXiv:2109.10303.
- Elias Stengel-Eskin, Peter Hase, and Mohit Bansal. 2024. Lacie: Listener-aware finetuning for confidence calibration in large language models. arXiv preprint arXiv:2405.21028.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In The 2023 Conference on Empirical Methods in Natural Language Processing.
- Yuanhe Tian, Ruyi Gan, Yan Song, Jiaxing Zhang, and Yongdong Zhang. 2024. Chimed-gpt: A chinese medical large language model with full training regime and better alignment to human preferences. Preprint, arXiv:2311.06025.
- Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. 2023. Huatuo: Tuning llama model with chinese medical knowledge. arXiv preprint arXiv:2304.06975.
- Xidong Wang, Guiming Chen, Song Dingjie, Zhang Zhiyi, Zhihong Chen, Qingying Xiao, Junying Chen, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, and Haizhou Li. 2024. CMB: A comprehensive medical benchmark in Chinese. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 6184-6205, Mexico City, Mexico. Association for Computational Linguistics.
- Bingbing Wen, Jihan Yao, Shangbin Feng, Chenjun Xu, Yulia Tsvetkov, Bill Howe, and Lucy Lu Wang. 2024. Know your limits: A survey of abstention in large language models. arXiv preprint arXiv:2407.18418.
- Ming Xu. 2023. Medicalgpt: Training medical https://github.com/shibing624/ gpt model. MedicalGPT.
- Wenjing Yue, Xiaoling Wang, Wei Zhu, Ming Guan, Huanran Zheng, Pengfei Wang, Changzhi Sun, and Xin Ma. 2024. Tembench: A comprehensive benchmark for evaluating large language models in traditional chinese medicine. Preprint, arXiv:2406.01126.
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi R Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2023. R-tuning: Teaching large

language models to refuse unknown questions. arXiv *preprint arXiv:2311.09677.*

Model Performance Analysis Α

A.1 **Visualization of Model Performance**

In this section, we present radar chart visualiza-445 tions to highlight the performance of both closed-446 source and open-source models across different 447 medication-related tasks. The radar charts provide 448 a clear, comparative view of how various models 449 handle tasks such as Indication, Dosage and Admin-450 istration, Contraindicated Population, and Mecha-451 nisms of Action. Notably, models like GLM4 and 452 XiaoYi stand out for their excellent performance, 453 with XiaoYi leading the closed-source models and 454 GLM4 showing remarkable consistency. On the 455 other hand, HuatuoGPT2 significantly outperforms 456 the other open-source models, as shown in Figure 3 457 and Figure 4. These findings underscore the impor-458 tance of model selection in high-stakes domains 459 like healthcare, where the quality of responses di-460 rectly impacts patient safety.



Figure 3: Radar Chart Representation of Close-Source Models Performance.



Figure 4: Radar Chart Representation of Open-Source Models Performance.

461

441

442

443

462 463 464

465

466

467

468

469

470

471

472

473

481

491

493

494

496

497 498

499

502

A.2 **Knowledge Mastery Assessment of Common Drugs**

To further evaluate model capabilities, we focus on a subset of 282 commonly used drugs. For each drug, we constructed questions pertaining to Indication, Dosage and Administration, Contraindicated Population, and Mechanism of Action, drawing from the benchmark dataset. The knowledge boundary of models was then assessed by visualizing their performance on these tasks, as shown in the radar charts for GLM4, XiaoYi, and GPT40 and the result is shown in Figure 3.

The radar charts depict knowledge boundaries 474 by showing the areas where models could answer 475 correctly once (orange area), 5 times (yellow area), 476 and areas where errors could be corrected after 5 477 times attempts (yellow with red area). These visu-478 alizations emphasize that while some models show 479 robustness in their knowledge, significant gaps re-480 main in certain drug-related tasks. The results indicate that GLM4 and XiaoYi exhibit stronger 482 consistency in answering these questions correctly 483 484 across the four tasks compared to GPT40. However, there were cases where even the most advanced 485 models struggled to demonstrate comprehensive 486 knowledge across all aspects of these drugs. This 487 highlights a key issue-despite their advanced capabilities, large models still fall short in areas of 489 490 medication-related knowledge. This reinforces the importance of exercising caution when deploying such models in high-risk areas like medication us-492 age.

	Dosage and Administration	Indication	Contraindicated Population	Mechanisms of Action	Avg.
XiaoYi	82.27	78.01	63.12	92.20	78.90
GPT40	67.02	74.11	70.2	92.91	76.24
ERNIE	69.15	64.89	58.51	91.13	70.92
Qwen-max	67.73	74.11	71.63	92.55	76.51
Doubao	75.89	74.47	64.54	92.55	76.86
GLM4	74.11	75.89	71.99	93.97	78.99
Claude3.5	67.09	67.73	54.26	93.46	70.64
Baichuan4	59.93	70.21	51.06	91.49	68.17

Table 3: Performance of various model on Common Drugs

B Semantic Entropy (SE) Method for **Knowledge Boundary Expression**

In this section, we explore the Semantic Entropy (SE) method used to detect knowledge boundaries, as introduced in Section 4.2. The SE method is particularly noteworthy for its effectiveness in expressing model uncertainty and improving response reliability, as demonstrated in our experiments. We applied this method to HuatuoGPT2 and XiaoYi, observing that it significantly enhanced the models' performance on challenging tasks, such as Medication Recommendations and Drug Interactions.

503

504

505

506

507

508

509

510

511

512

513

514

515

516

As shown in Table 4, the SE method proved to be robust and consistent across different model architectures and sizes. It improved both Precision and Abstain Accuracy, regardless of the model's scale. This reinforces the notion that SE is an effective tool for managing uncertainty, making it an essential method for enhancing the reliability of models in real-world medical applications.

С **Case Study**

In Figure 6, we present a case to illustrate the practical section of the ChiDrug.



Figure 5: Knowledge boundary chart for GLM4, XiaoYi, and GPT40 across 282 common drugs. The orange area indicates that the model answered correctly once, while the yellow area indicates 5 times opportunities to answer correctly. The red areas in the yellow-covered region represent cases where a model made an error in a single attempt but was able to recover after multiple tries.

Model	Method	Dosage and Administration		Indication		Contraindicated Population		Mechanisms of Action		Medication Recommendation		Drug Interaction	
		Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc	Precision	A-Acc
HuatuoGPT2	w/o SE	55.83	55.83	47.03	47.03	18.66	18.66	77.16	77.16	25.18	25.18	25.60	25.60
	SE	64.79	55.32	70.28	52.71	28.11	26.95	86.24	83.09	29.22	24.95	38.56	31.72
XiaoYi	w/o SE	81.1	81.1	77.87	77.87	66.71	66.71	92.85	92.85	59.31	59.31	63.27	3.27
	SE	83.42	84.46	94.01	94.15	80.14	83.28	91.71	92.1	67.71	75.71	63.85	69.02

Table 4: Application of the SE method on HuatuoGPT2 and XiaoYi models, showcasing the performance improvements achieved through the SE method. This method enhances precision and uncertainty handling, effectively reducing hallucinations.

Tasks	Question	Answer
Dosage and Administration	复方锌布颗粒剂的推荐服用方式为多少包,多少次一天? (A)儿童5岁以下一次2包,一天3次 (B)成人一次2包,一天3次 (C)成人一次1包,一天3次 (D)6~14岁儿童一次1包,一天2次 (E)6~14岁儿童一次1包,一天3次	B 用法用量: 口服。3~5岁儿童,一次半包:6~14岁儿童,一 次1包,成人,一次2包,一日3次。
Indication	复方锌布颗粒剂主要用于缓解以下哪些症状? (A)普通感冒引起的发热 (B)急性肠胃炎 (C)普通感冒引起的四肢酸痛 (D)普通感冒引起的打喷嚏	ACD 适应症: 用于缓解普通感冒或流行性感冒引起的发热、头痛、 四肢酸痛、鼻塞、流涕、打喷嚏等症状。
Contraindicated Population	复方锌布颗粒剂不适用于以下哪些人群? (A)心脏病患者 (B)哺乳期妇女 (C)对阿司匹林过敏的哮喘患者 (D)高血压患者	BC 茶忌: 1.对其他非甾体抗炎药过敏者禁用。2.孕妇及哺乳 期妇女禁用。3.对阿司匹林过敏的哮喘患者禁用。
Mechanism of Action	关于复方锌布颗粒剂各组分的主要药理作用是: (A)布洛芬具有抗炎作用,葡萄糖酸锌促进蛋白质合成,马来酸氯苯那敏 为解热镇痛药 (B)布洛芬具有解热镇痛作用,葡萄糖酸锌能增强免疫功能,马来酸氯苯 那敏为抗组胺药 (C)布洛芬为抗组胺药,葡萄糖酸锌具有解热功能,马来酸氯苯那敏具 有镇痛作用 (D)布洛芳为解热镇痛药,葡萄糖酸锌参与多种酶的合成与激活,马来 酸氯苯那敏为抗组胺药	D 作用机制: 布洛芬能抑制前列腺素合成,具有解热镇痛作用: 葡萄糖酸锌中锌离子能参与多种酶的合成与激活, 有增强吞噬细胞的吞噬能力的作用;马来酸氯苯那 敏为抗组胶药,能减轻由感冒或流感引起的鼻塞、 流涕、打喷嚏等症状。
Drug Recommendatio	回答以下不定项选择题(可能包含1个或多个正确选项): 一位孕晚期患者因为感冒出现咳嗽、喉咙里有异物感以及扁桃体发炎, 可以考虑推荐的药物是: (A)双黄连口服液 (B)热毒宁注射液 (C)贝美前列素滴眼液 (D)银苓胶囊	A 双黄连口服液的适应症: 疏风解表,清热解毒。用于外感风热所致的感冒, 症见发热、咳嗽、咽痛。
Drug Interaction	注射用降纤酶与抗纤溶药联用的风险等级是?	高风险 注射用降纤酶 使用本品应避免与水杨酸类药物(如:阿司匹林) 合用。抗凝血药可加强本品作用,引起意外出血; 抗纤溶药可抵消本品作用,禁止联用。

Figure 6: Partial cases of ChiDrug on 6 sub datasets.