

# Diabetica: Developing Specialized Large Language Models for Diabetes Care with Multi-Faceted Benchmarks and Clinical Validation

Anonymous ACL submission

## Abstract

Diabetes is a chronic disease with a significant global health burden, where the effectiveness of Large Language Models (LLMs) across diverse diabetes tasks remains unproven. To address this, we introduce a comprehensive framework for developing and evaluating diabetes-specialized LLMs. Our methodology begins with systematically collecting unstructured diabetes-related data from sources like clinical guidelines and medical textbooks. This data then undergoes rigorous processing, including filtering, transformation, and refinement via a novel self-distillation approach, to produce a high-quality training dataset and a multi-faceted benchmarking suite. Leveraging this training dataset, the model we obtain, Diabetica, significantly outperforms existing open-source LLMs of comparable size across diabetes-specific benchmarks, including multiple choice questions, fill-in-the-blank assessments, and open-ended dialogue scenarios. To demonstrate real-world applicability, we conduct extensive clinical assessments in three key use cases: medical counseling, medical education, and record summarization. Results reveal that Diabetica provides more thorough and empathetic patient responses than human physicians, achieves expert-level performance in medical examinations, and significantly improves clinical documentation efficiency while maintaining high quality.

## 1 Introduction

Diabetes mellitus is a pervasive chronic disease, with its prevalence projected to exceed 783 million by 2045, placing an unsustainable burden on global healthcare systems (Sun et al., 2022). Persistent challenges, including shortages of specialists and gaps in patient education, lead to poor glycemic control and significant societal costs (Guan et al., 2023). This creates an urgent need for accessible, reliable, and efficient tools to augment diabetes

management.

Large Language Models (LLMs), with their advanced language comprehension and reasoning capabilities (OpenAI, 2023; Team et al., 2023), have emerged as a promising solution for various medical applications (Chen et al., 2024b; Li et al., 2024). However, existing open-sourced medical LLMs often lack the deep, specialized knowledge required for the nuances of diabetes care. This deficiency stems primarily from the absence of high-quality, domain-specific training data and tailored optimization paradigms. Thus, developing a specialized LLM for diabetes is a critical step toward providing effective support for both patients and clinicians.

Once a specialized LLM is developed, it is essential to rigorously evaluate its reliability and effectiveness through objective and comprehensive assessments. However, there is a lack of benchmarks for diabetes specialties. Clinical practice also differs significantly from simply answering examination questions correctly, and finding appropriate benchmarks to gauge the clinical potential of LLMs is a substantial challenge (Thirunavukarasu et al., 2023). To address this gap, a dedicated evaluation is needed, one that considers both a controlled benchmark evaluation and realistic clinical scenarios, to accurately measure the practical value of the model in diabetes care.

To address these challenges, we introduce a comprehensive framework (Figure 1) for developing and evaluating diabetes-specialized LLMs. Our approach begins with the systematic curation of a high-quality, diabetes-focused corpus from diverse sources, including clinical guidelines, textbooks, and real-world dialogues. This data undergoes a rigorous pipeline of domain-specific filtering, transformation into instruction-following formats, and refinement using a novel self-distillation method. The resulting dataset is used to fine-tune our model, Diabetica, enhancing its specialized knowledge

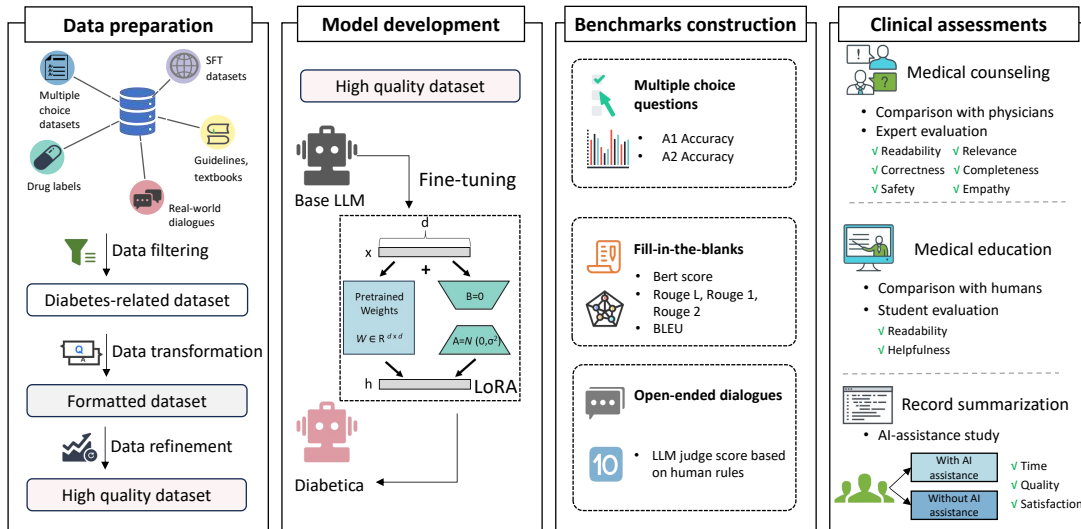


Figure 1: The overall framework for developing and evaluating diabetes-specialized LLMs.

while preserving foundational capabilities.

A core component of our framework is a new, multi-faceted benchmarking suite designed specifically for diabetes care. It assesses both knowledge (via multiple-choice and fill-in-the-blank questions) and interactive skills (via dialogue generation). Our evaluations show that Diabetica significantly outperforms comparable open-source models (e.g., Qwen2 (Yang et al., 2024a), Llama3 (Dubey et al., 2024)) and achieves performance on par with leading proprietary models like GPT-4 (Hurst et al., 2024) and Claude-3.5 (Anthropic, 2024).

To validate real-world utility, we designed three clinical-practice-aligned assessments. Our results demonstrate that Diabetica offers tangible benefits across different end-users. In medical counseling, it provides more empathetic and thorough patient advice than human physicians. In medical education, it delivers expert-level explanations that students find more helpful than standard textbooks. Finally, in record summarization, it significantly improves clinical documentation efficiency and quality, reducing physician workload. These findings underscore Diabetica’s robust knowledge, reliability, and practical value in diabetes care.

In summary, our main contributions are:

- We propose a systematic framework for developing specialized medical LLMs, detailing a reproducible pipeline from high-quality data curation to domain-specific fine-tuning.
- We introduce a multi-faceted benchmark suite

tailored for diabetes, enabling robust and comprehensive evaluation of LLMs in this domain. Our model, Diabetica, sets a new state-of-the-art for open-source diabetes-specific models.

- Through extensive clinical assessments across three key use cases, we provide strong evidence of a specialized LLM’s practical value in enhancing patient care, medical education, and clinical efficiency.

## 2 Related Works

The advancement of artificial intelligence technology presents a significant opportunity to enhance diabetes care efficiency. Various AI-based tools for diabetes care, such as those for diagnosis (da Silva Santos et al., 2022; Rabie et al., 2022), insulin titration (Rabie et al., 2022; Wang et al., 2023a), and retinal image analysis (Arcadu et al., 2019; Dai et al., 2021), have demonstrated impressive performance. These previous AI models in diabetes management, despite advantageous in certain aspects, are so far predominantly single-task oriented (Li et al., 2024) or face challenges in comprehending and generating natural language (Wang et al., 2023a). Notably, several studies have explored fine-tuning and evaluating LLMs for medical applications in fields such as radiology, cancer, ophthalmology, etc (Zhang et al., 2024; Agnihotri et al., 2024; Van Veen et al., 2024; Zhou et al., 2024a; Chen et al., 2024b; Zhou et al., 2024b; Arora et al., 2025). However, few have specifically focused on applying LLMs to diabetes-related sce-

145 narios, primarily due to the challenges in collecting  
 146 and cleaning high-quality diabetes-specific data,  
 147 and the lack of standardized evaluation benchmarks.  
 148 These limitations narrow down their potentials to  
 149 offer well-rounded and easily understandable diabe-  
 150 tes supports across diverse user groups. Thus,  
 151 recent works have highlighted the need for tailored  
 152 datasets and targeted training methodologies to en-  
 153 hance LLM performance in specific medical do-  
 154 mains (Sheng et al., 2024). To address these gaps,  
 155 our study introduces a comprehensive framework  
 156 for developing diabetes-specialized LLMs, multi-  
 157 faceted benchmarks specifically for the diabetes  
 158 field, and clinical studies to evaluate the model’s ef-  
 159 ficacy in real-world settings, demonstrating LLM’s  
 160 potential applications in diabetes management.

### 161 3 Method

162 In this section, we present the framework for diabe-  
 163 tes care in detail, including (1) *Data Preparation*,  
 164 which involves systematic data collection, filter-  
 165 ing, transformation, and refinement to construct a  
 166 high-quality, diabetes-specific training dataset; (2)  
 167 *Model Development*, where the LLM is trained us-  
 168 ing supervised fine-tuning and parameter-efficient  
 169 methods to ensure domain expertise while preserv-  
 170 ing general language capabilities; (3) *Benchmarks*  
 171 *Construction*, in which multiple evaluation bench-  
 172 marks are curated to rigorously test the model’s  
 173 performance across diverse diabetes-related tasks;  
 174 and (4) *Clinical Assessments*, where the model is  
 175 assessed in real-world clinical scenarios such as pa-  
 176 tient consultation, medical education, and clinical  
 177 record summarization. This framework enables the  
 178 creation of reliable and helpful LLMs for practical  
 179 diabetes care.

#### 180 3.1 Data Preparation

181 To establish a robust foundation for developing  
 182 our diabetes-specialized LLM, we meticulously  
 183 design a well-structured data preparation pipeline  
 184 that systematically transforms raw medical con-  
 185 tent into high-quality training data. This pipeline  
 186 comprises four essential stages: data collection, fil-  
 187 tering, transformation, and refinement. Each stage  
 188 is carefully designed to optimize the data for effec-  
 189 tive model training while maintaining the medical  
 190 accuracy, which is introduced individually in detail  
 191 as follows.

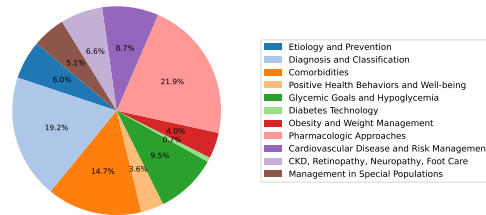


Figure 2: Category-wise distribution of our collected dataset.

#### 3.1.1 Data Collection

192 We aggregate a diabetes-specific dataset from five  
 193 complementary sources, covering 11 major clinical  
 194 categories (Figure 2). The collection integrates pub-  
 195 lic benchmarks, medical SFT corpora, authoritative  
 196 documents, drug labels, and specialist-annotated  
 197 real-world dialogues, enabling the model to learn  
 198 both medical knowledge and clinically grounded re-  
 199 sponses. Detailed descriptions of each data source  
 200 are provided in Appendix A.

#### 3.1.2 Data Filtering

202 Note that the initially collected dataset includes  
 203 some general medical corpora, which inevitably  
 204 contain content unrelated to diabetes. To ensure  
 205 the domain specificity of our dataset, we introduce  
 206 a dedicated data filtering stage to remove irrele-  
 207 vant information. First, we apply keyword-based  
 208 filtering to isolate diabetes-related content. This  
 209 involves a dual-strategy approach using a curated  
 210 set of inclusion terms (e.g., "diabetes", "DKA") de-  
 211 veloped with specialists, and an exclusion list (e.g.,  
 212 "insulinoma") to filter out related but non-diabetes  
 213 endocrinological topics. All filtered content un-  
 214 dergoes subsequent manual verification to ensure  
 215 precision. Second, to mitigate the negative impact  
 216 of data redundancy on model performance, we em-  
 217 ploy SemDeDup (Abbas et al., 2023), an advanced  
 218 method to perform deduplication, which helps re-  
 219 move semantic duplicates by identifying text pairs  
 220 with high embedding similarity. Implementation  
 221 details are provided in Appendix E.1.

#### 3.1.3 Data Transformation

222 The collected datasets can be broadly categorized  
 223 into two types: long-form textual data and question-  
 224 answer pairs. To standardize these heterogeneous  
 225 sources for model training, we apply two corre-  
 226 sponding data transformation techniques, convert-  
 227 ing each type into a unified instruction-response  
 228 pair.

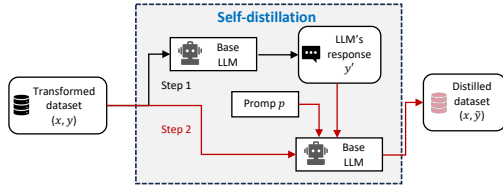


Figure 3: The overall pipeline of self-distillation. Firstly, we collect the base LLM’s responses to each instruction in the dataset. Secondly, we use a specific prompt to let the base LLM generate a refined response based on the instruction, the original response and its own response. The final refined responses are combined into a distilled dataset for supervised fine-tuning to develop Diabetica.

format. For long-form texts (e.g., guidelines, textbooks), we first segment documents into conceptually coherent knowledge units. Subsequently, we leverage GPT-4 to automatically convert these units into diverse conversational dialogues and fill-in-the-blank exercises, transforming static knowledge into interactive training data. For question-answer pairs, we adapt the methodology from Huang et al. (2023). Each multiple-choice question and its options are formatted into a single prompt. We then use GPT-4 with chain-of-thought prompting to generate a detailed rationale for the correct answer, creating a comprehensive instruction-response pair. To ensure data fidelity, only pairs with verified correct answers are retained. All prompts used for transformation are detailed in Appendix C.

### 3.1.4 Data Refinement

Note that directly fine-tuning a base LLM on a domain-specific dataset, whose distribution is far from the LLM’s, can be harmful (Ren et al., 2024a). This can undermine the model’s alignment with human preferences and cause catastrophic forgetting of its general instruction-following abilities. Consequently, the overall quality of the LLM’s responses can deteriorate.

To mitigate these risks from the distribution shift, we implement a further data refinement process tailored for the instruction-response dataset. In particular, we design an enhanced self-distillation approach inspired by Yang et al. (2024c). Our method minimizes distribution shift between the refined training data and the foundational knowledge already encoded in the base LLM, thereby preserving the model’s foundational capabilities while enhancing diabetes-specific performance.

Our self-distillation method contains two steps. Given a transformed dataset containing instruction  $x$  and response  $y$  for each sample, we firstly collect the base LLM’s own response  $y' = f_{\text{base}}(x)$  of each instruction  $x$ , where  $f_{\text{base}}$  denotes the base LLM. Secondly, we leverage a manually designed prompt  $p$  (shown in Appendix C) to let the base LLM generate a refined response  $\tilde{y}$  based on the instruction  $x$ , the original response  $y$  and its own response  $y'$ :  $\tilde{y} = f_{\text{base}}(x, y, y', p)$ . Finally, the refined response  $\tilde{y}$  is used to replace the original response  $y$  in the fine-tuning stage. The above procedure is summarized as Figure 3.

This method helps minimize the distribution shift between the training data and the knowledge encoded in the base LLM. In particular, the original response  $y$  is accurate, reflecting the intended diabetes knowledge and information. The subsequent base LLM’s own response  $y'$  aligns with the internal distribution of the base LLM. Rewriting based on these two responses, the base LLM can create a refined response  $\tilde{y}$ . Hence, the distribution shift between the model’s prior knowledge and the self-distilled dataset is mitigated. The newly generated self-distillation dataset can not only help the model learn new domain knowledge, but also restore the prior knowledge distribution of the model. We also provide further analysis of the self-distillation method in Appendix D.3.

## 3.2 Model Development

Based on our final refined dataset containing 28.4K high-quality samples in total, we apply supervised fine-tuning to train Qwen2-7B-Instruct (Yang et al., 2024a) and obtain the Diabetica. To ensure both efficiency and effectiveness during training, we employ LoRA (Hu et al., 2021) tuning, which helps preserve the model’s general language understanding capabilities, ensuring that Diabetica not only acquires deep expertise in diabetes care but also retains strong performance across a wide range of general language tasks.

## 3.3 Benchmarks Construction

To comprehensively assess the utility of LLMs in various diabetes tasks, we construct three distinct benchmarks: multiple-choice questions, fill-in-the-blank questions, and open-ended questions. The benchmark data are derived from the same source as the training dataset but are processed independently to ensure clear separation.

314	<b>Multiple Choices Questions.</b> The multiple	coverage, reliability, and practical utility.	363
315	choices questions benchmark includes 312		
316	multiple-choice questions, specifically 235 Type		
317	A1 and 77 Type A2 questions, extracted from	<b>Medical Counseling</b> To assess Diabetica’s po-	364
318	the Advanced Health Professional Technical	tential in assisting patients, we curate 20 diabetes-	365
319	Qualification Examination. Type A1 questions are	related cases from an online consultation plat-	366
320	designed to assess the examinee’s foundational	form (Fu, 2024). Each case includes patient queries	367
321	knowledge in endocrinology, encompassing a	and associated physician responses. We input the	368
322	broad range of topics from the pathophysiology	full text of each case into Diabetica to generate its	369
323	of various diabetes forms to the pharmacological	response. Three healthcare professionals indepen-	370
324	fundamentals of antidiabetic medications. Con-	dently evaluate physician and Diabetica responses	371
325	versely, Type A2 questions are crafted within	on readability, relevance, correctness, complete-	372
326	specific clinical contexts, challenging examinees to	ness, safety, and empathy using a 5-point Likert	373
327	apply their knowledge in diagnosing and making	scale (Table 5 in the Appendix). Evaluators are also	374
328	evidence-based medical decisions. We calculate	asked to compare these two responses and select	375
329	accuracy that measures the percentage of correct	the superior one.	376
330	answers given by a model for these questions.		
331	<b>Fill-in-the-Blanks</b> Besides the multiple choices	<b>Medical Education</b> To explore the usefulness of	377
332	questions, fill-in-the-blanks task is another popu-	Diabetica for students’ medical education, we first	378
333	lar exam type in medical education. In particu-	compare the accuracy of its responses with those	379
334	lar, we manually create a set of fill-in-the-blanks	of human (3 medical students, 3 junior physicians,	380
335	questions for this benchmark that includes 35	3 mid-level physicians, and 3 senior physicians)	381
336	questions from the guideline and textbook. We	in 67 A2 type multiple-choice questions. These	382
337	use five evaluation metrics of textual similarity:	questions typically involve complex clinical sce-	383
338	BERTScore (Zhang et al., 2019), ROUGE-L (Lin,	narios and require high-order medical reasoning.	384
339	2004), ROUGE-1 (Lin, 2004), ROUGE-2 (Lin,	Next, we investigate the model’s ability to provide	385
340	2004), and BLEU (Papineni et al., 2002), to as-	explanations for questions that were incorrectly	386
341	sess the models performance in fill-in-the-blank	answered by students. The readability and helpfu-	387
342	benchmark.	ness of these model-generated explanations, as well	388
343	<b>Open-Ended Dialogues</b> To evaluate the model’s	as reference explanations from textbooks, are com-	389
344	dialogue capabilities in real world applications, we	pared and evaluated by the corresponding students	390
345	construct the open-ended dialogues benchmark that	using a 5-point Likert scale.	391
346	includes 120 instances covering various aspects of		
347	diabetes. Each instance consists of three elements:	<b>Record Summarization</b> To assess Diabetica’s	392
348	a category, a question, and the associated rules. In	efficiency in reducing physician workload, we as-	393
349	particular, physicians annotate a comprehensive	semble five real-world diabetes cases and conduct	394
350	set of rules that define the criteria for evaluating	a crossover multi-reader multi-case study involv-	395
351	the quality of an answer. Then, we employ strong	ing eight intern doctors. These interns are asked	396
352	LLMs (GPT-4 and Claude-3.5) as judges to evalu-	to write patient records based on multi-turn dia-	397
353	ate these models on open-ended questions and rate	logues with physicians. Using a crossover design,	398
354	each answer on a scale of 1-10 based on the human	we randomly and equally divide the interns into	399
355	rule (Zheng et al., 2023). Detailed prompt is shown	group A (first read cases without Diabetica assis-	400
356	in Appendix C.	tance) and group B (first read cases with Diabetica	401
357	<b>3.4 Clinical Assessments</b>	assistance). After a washout period of 1 week, the	402
358	To further explore the performance of LLM in di-	arrangement is reversed. We record the total time	403
359	abetes care clinical scenarios, we conduct clini-	each intern spends on the task. Then, the quality	404
360	cal assessments in three distinct settings: medical	of the resulting records—assessed in terms of com-	405
361	counseling, medical education, and record summa-	pleteness, conciseness, and correctness—is rated by	406
362	rization. These examine the model’s knowledge	three medical experts using a 5-point Likert scale	407
		(Table 6 in the Appendix). After that, all interns	408
		are invited to complete a satisfaction questionnaire	409
		within a week.	410

## 4 Evaluating Diabetica on the Proposed Benchmarks

In this section, we present the evaluation results of Diabetica and other baseline models on the proposed diabetes-specific benchmarks. We conduct comprehensive experiments across multiple-choice questions, fill-in-the-blanks tasks, and open-ended dialogues to assess the models’ capabilities in diabetes domain knowledge as follows.

### 4.1 Baselines

We select a large amount of models as our baselines, including proprietary LLMs like GPT-4 (OpenAI, 2023) and Claude-3.5 (Anthropic, 2024), open-source general LLMs like Qwen2-7B (Yang et al., 2024a), InternLM2-7B (Cai et al., 2024), Llama3-8B (Dubey et al., 2024) and Yi-1.5-9B (Young et al., 2024), as well as open-source medical LLMs like Meditron-7B (Chen et al., 2023), MMedLM-7B (Qiu et al., 2024) and Apollo-7B (Wang et al., 2024).

### 4.2 Results

On the MCQ benchmark, Diabetica achieves 87.2% accuracy, significantly surpassing GPT-4 and Claude-3.5 (Figure 4, left). Its strong performance across both A1 (knowledge) and A2 (case study) questions demonstrates a balanced mastery of theoretical concepts and clinical reasoning. In fill-in-the-blank assessments, Diabetica-7B outperforms all comparable open-source models and matches proprietary models like GPT-4 across all metrics, including a BERTScore of 0.9298 and ROUGE-L of 0.7828 (Figure 4, right; Table 1). These results underscore its superior context understanding in diabetes care. For open-ended dialogues, Diabetica consistently leads open-source LLMs of similar size (Figure 5), earning scores of 7.81 (GPT-4 rated) and 7.96 (Claude-3.5 rated). Furthermore, ablation studies and forgetting tests (Appendix D) confirm that our self-distillation method effectively injects specialized knowledge while preserving foundational capabilities. Overall, Diabetica-7B exhibits state-of-the-art proficiency in recalling medical knowledge, identifying critical clinical points, and handling complex patient consultations.

Notably, we also conduct additional experiments shown in Appendix D, such as measuring catastrophic forgetting, ablation studies and validation

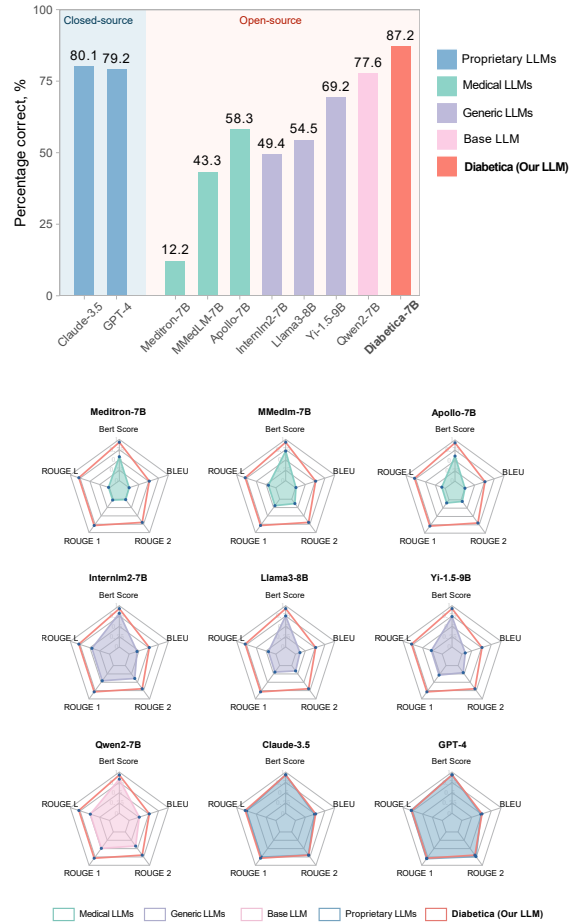


Figure 4: Performances of different LLMs in diabetes-related benchmarks, including multiple-choice questions (top) and fill-in-the-blank questions (bottom). Diabetica achieves leading performances among these LLMs.

of the proposed self-distillation method. These experiments further demonstrate the robustness and effectiveness of our approach in developing diabetes-specialized LLM.

## 5 Clinical Assessments of Diabetica

In this section, we explore three real-world clinical applications using Diabetica, noting that the model is not specifically trained for these tasks.

### 5.1 Medical Counseling

We first evaluated Diabetica’s potential in medical consultation using 20 online patient cases. Three endocrinology specialists rated responses from Diabetica and physicians across multiple dimensions using a 5-point Likert scale. As shown in Figure 6, Diabetica’s responses significantly exceeds human responses with mean (and the corresponding stan-

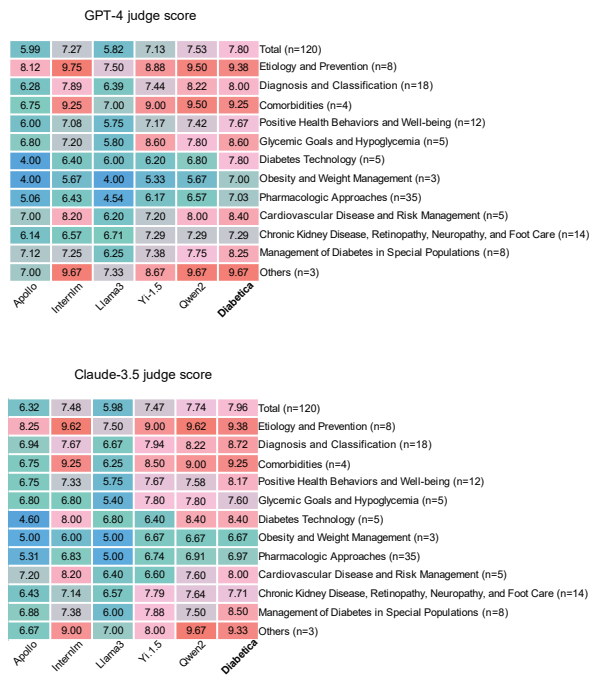


Figure 5: GPT-4 and Claude-3.5 judged scores of different LLMs in the dialogue benchmark.

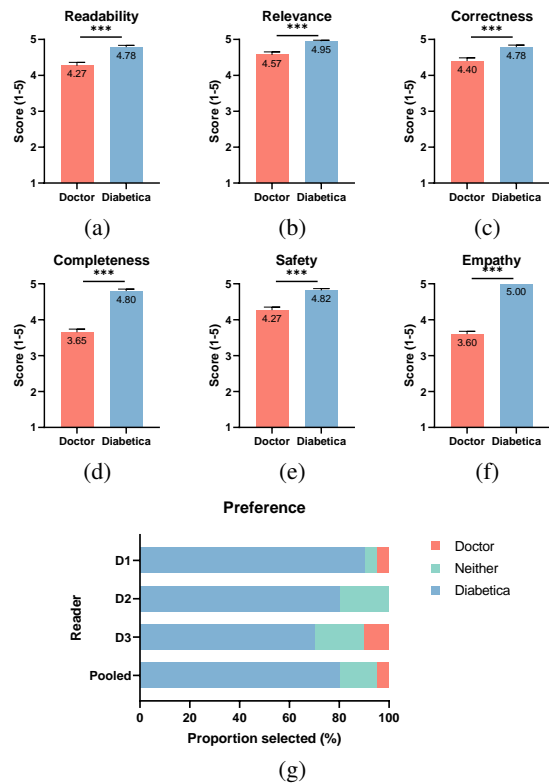


Figure 6: Performance comparison of the AI-generated and physician-delivered responses of online patient cases (n=20). Evaluation is based on the expert panel review including (a) readability, (b) relevance, (c) correctness, (d) completeness, (e) safety, (f) empathy, and (g) selected superior responses. Bar graphs indicate the mean  $\pm$  s.e.m., \*\*\*P<0.001, calculated by paired-Wilcoxon test.

standard deviation-SD) values of 4.78 (0.42) for readability, 4.95 (0.22) for relevance, 4.78 (0.45) for correctness, 4.80 (0.40) for completeness, 4.82 (0.39) for safety, and 5.00 (0) for empathy (all p values <0.001). Table 2 in the Appendix presents reader-specific scores and demonstrates good inter-reader reliability. In addition, 80.0% of responses were judged superior for Diabetica, indicating its overall advantage over physician responses based on expert evaluation.

## 5.2 Medical Education

We further evaluate Diabetica’s educational potential by comparing its performance with medical professionals of different experience levels. Specifically, Diabetica achieves 84.4% accuracy, significantly outperforming medical students (53.7%), junior physicians (69.7%), and intermediate physicians (74.0%), while slightly exceeding senior physicians (83.5%) (Figure 7(a)). These results indicate that Diabetica demonstrates expert-level proficiency on diabetes specialty examinations.

Beyond quantitative evaluation, we explore Diabetica’s pedagogical ability by examining its explanations for incorrect answers. Three medical students compared explanations from Diabetica and a reference textbook, rating readability and helpfulness using a 5-point Likert scale. As illus-

trated in Figure 7(b), medical students rate Diabetica’s explanations as helpful in 71.96% of cases and readable in 65.42% of cases, demonstrating quality comparable to reference materials. Statistical analysis reveals no significant differences in mean scores between Diabetica and reference explanations for either readability (3.67 vs 3.85) or helpfulness (3.89 vs 3.94), with all p-values > 0.05 (Figure 7(c-d)). These findings suggest that Diabetica can serve as a reliable and effective educational tool to providing helpful explanations.

## 5.3 Record Summarization

Another helpful application of LLMs is assisting clinicians in summarizing patient records to reduce workload. In particular, we conduct a cross-over AI-assistance study to explore the potential of Diabetica as a clinical support tool. As shown in Figure 8(a-e), Diabetica assistance reduced documentation time by 23% (750 vs. 976 seconds per case, p < 0.05) and significantly improved record

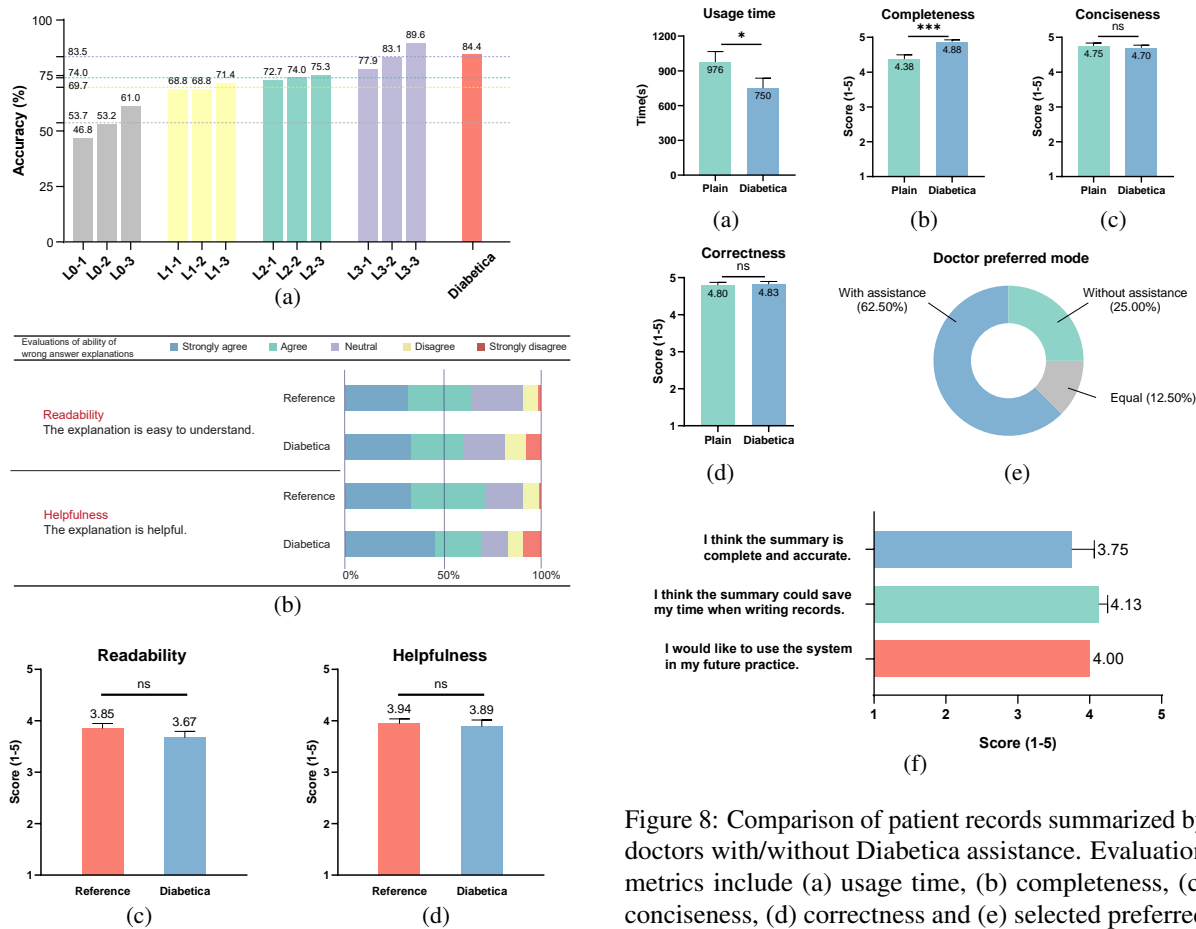


Figure 7: Performance on medical education. (a) Accuracy in answering A2-type multiple-choice questions of medical students, physicians with different levels, and LLMs in the MCQ examination. (b) Student evaluation of the helpfulness and readability of answer explanations from Diabetica and reference. (c) The readability and (d) helpfulness scores of answer explanations from Diabetica and reference. There is no significant difference (ns), calculated by the paired Wilcoxon test.

completeness among intern doctors (4.88 vs. 4.38,  $p < 0.001$ ). Finally, intern perceptions of Diabetica were assessed using a user satisfaction questionnaire completed by the eight participated interns. As shown in Figure 8(f), Diabetica received average scores of 3.75 for summary accuracy and completeness, 4.13 for time-saving, and 4.00 for future clinical use. Five of eight interns prefer to have AI assistance when writing medical records. These results indicate that Diabetica can effectively enhance clinical workflow efficiency while maintaining documentation quality, with strong physician acceptance supporting its potential for clinical implementation.

Overall, though not specifically trained on data for these clinical tasks, Diabetica still shows strong

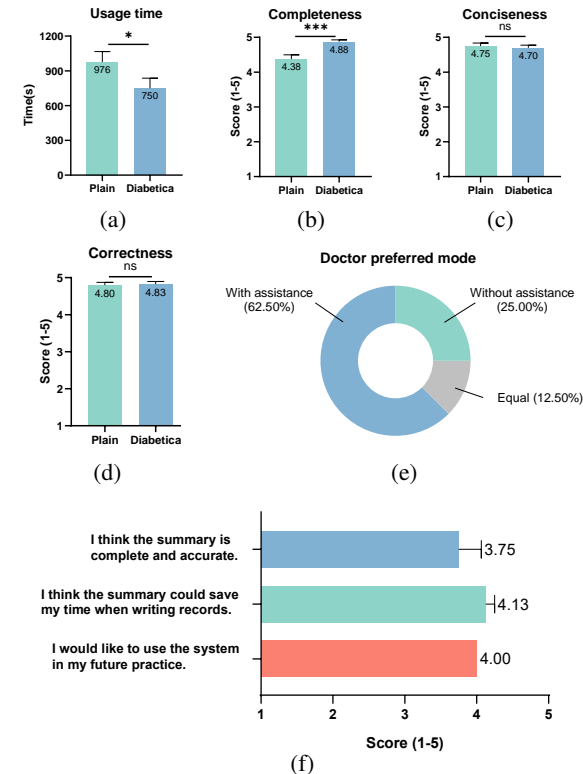


Figure 8: Comparison of patient records summarized by doctors with/without Diabetica assistance. Evaluation metrics include (a) usage time, (b) completeness, (c) conciseness, (d) correctness and (e) selected preferred responses. (f) Satisfaction of participated doctors (score ranges from 1-5). Bar graphs indicate the mean  $\pm$  s.e.m., \* $P < 0.05$ , \*\*\* $P < 0.001$ , ns, no significant difference.

reliability, knowledge coverage, and practical value

## 6 Conclusion

This work presents a comprehensive framework for developing and evaluating domain-specific LLMs for diabetes care. We demonstrate that a meticulous data curation pipeline with a novel self-distillation technique, enables our model, Diabetica, to achieve leading performance on our new suite of diabetes-specific benchmarks. Beyond these benchmarks, extensive clinical assessments validate Diabetica's practical value. The model delivers more empathetic medical counseling than human physicians, achieves expert-level performance in medical education, and enhances clinical documentation by reducing writing time by 23% while improving record completeness. In conclusion, our findings establish a reproducible pathway for creating a specialized medical LLM and confirm its potential as an effective, reliable tool in clinical practice.

## 558 Limitations

559 While Diabetica demonstrates strong performance  
560 across various benchmarks, this study has a few  
561 limitations. The clinical validation was conducted  
562 on a relatively small scale of patient cases and med-  
563 ical professionals, and the model’s expertise is nat-  
564 urally bounded by its training data cutoff. Further-  
565 566 more, as a text-centric model, Diabetica currently  
567 lacks the ability to directly process multi-modal  
568 clinical data, such as continuous glucose monitor-  
569 ing (CGM) sensor trends or medical imaging. Fu-  
570 ture work will focus on expanding the evaluation  
571 scope and exploring multi-modal integration to fur-  
ther support comprehensive diabetes management.

## 572 References

573 Amro Abbas, Kushal Tirumala, Dániel Simig, Surya  
574 Ganguli, and Ari S Morcos. 2023. Semdedup: Data-  
575 efficient learning at web-scale through semantic dedup-  
576 lication. *arXiv preprint arXiv:2303.09540*.

577 Akshay Prashant Agnihotri, Ines Doris Nagel, Jose  
578 Carlo M Artiaga, Ma Carmela B Guevarra, George  
579 Michael N Sosuan, and Fritz Gerald P Kalaw. 2024.  
580 Large language models in ophthalmology: A review  
581 of publications from top ophthalmology journals.  
582 *Ophthalmology Science*, page 100681.

583 Anthropic. 2024. **Introducing claude 3.5 sonnet**. news.

584 Filippo Arcadu, Fethallah Benmansour, Andreas Maunz,  
585 Jeff Willis, Zdenka Haskova, and Marco Prunotto.  
586 2019. Deep learning algorithm predicts diabetic  
587 retinopathy progression in individual patients. *NPJ*  
588 *digital medicine*, 2(1):92.

589 Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Pre-  
590 ston Bowman, Joaquin Quiñero-Candela, Foivos  
591 Tsimplouras, Michael Sharman, Meghan Shah, An-  
592 drea Vallone, Alex Beutel, and 1 others. 2025.  
593 Healthbench: Evaluating large language models  
594 towards improved human health. *arXiv preprint*  
595 *arXiv:2505.08775*.

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

600 Zhijie Bao, Wei Chen, Shengze Xiao, Kuang Ren, Jiaao  
601 Wu, Cheng Zhong, Jiajie Peng, Xuanjing Huang, and  
602 Zhongyu Wei. 2023. Disc-medllm: Bridging gen-  
603 eral large language models and real-world medical  
604 consultation. *arXiv preprint arXiv:2308.14346*.

605 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen,  
606 Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi  
607 Chen, Pei Chu, and 1 others. 2024. Internlm2 techni-  
608 cal report. *arXiv preprint arXiv:2403.17297*.

Dejie Chang, Mosha Chen, Chaozhen Liu, Liping Liu,  
Dongdong Li, Wei Li, Fei Kong, Bangchang Liu,  
Xiaobin Luo, Ji Qi, and 1 others. 2021. Diakg: An  
annotated diabetes dataset for medical knowledge  
graph construction. In *Knowledge Graph and Se-  
mantic Computing: Knowledge Graph Empowers  
New Infrastructure Construction: 6th China Confer-  
ence, CCKS 2021, Guangzhou, China, November 4-7,  
2021, Proceedings 6*, pages 308–314. Springer. 609  
610  
611  
612  
613  
614  
615  
616  
617

Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang,  
Wanlong Liu, Rongsheng Wang, Jianye Hou, and  
Benyou Wang. 2024a. Huatuogpt-o1, towards med-  
ical complex reasoning with llms. *arXiv preprint*  
*arXiv:2412.18925*. 618  
619  
620  
621  
622

Xiaolan Chen, Weiyi Zhang, Pusheng Xu, Ziwei Zhao,  
Yingfeng Zheng, Danli Shi, and Mingguang He.  
2024b. Ffa-gpt: an automated pipeline for fundus  
fluorescein angiography interpretation and question-  
answer. *NPJ digital medicine*, 7(1):111. 623  
624  
625  
626  
627

Zeming Chen, Alejandro Hernández Cano, Angelika  
Romanou, Antoine Bonnet, Kyle Matoba, Francesco  
Salvi, Matteo Pagliardini, Simin Fan, Andreas  
Köpf, Amirkeivan Mohtashami, and 1 others. 2023.  
Meditron-70b: Scaling medical pretraining for large  
language models. *arXiv preprint arXiv:2311.16079*. 628  
629  
630  
631  
632  
633

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
Nakano, and 1 others. 2021. Training verifiers  
to solve math word problems. *arXiv preprint*  
*arXiv:2110.14168*. 634  
635  
636  
637  
638  
639

Tiago da Silva Santos, Liliana Fonseca, Sílvia San-  
tos Monteiro, Diana Borges Duarte, Ana Mar-  
tins Lopes, André Couto de Carvalho, Maria João  
Oliveira, Teresa Borges, Francisco Laranjeira,  
María Luz Couce, and 1 others. 2022. Mody prob-  
ability calculator utility in individuals’ selection for  
genetic testing: Its accuracy and performance. *En-  
docrinology, Diabetes & Metabolism*, 5(5):e00332. 640  
641  
642  
643  
644  
645  
646  
647

Ling Dai, Liang Wu, Huating Li, Chun Cai, Qiang  
Wu, Hongyu Kong, Ruhan Liu, Xiangning Wang,  
Xuhong Hou, Yuexing Liu, and 1 others. 2021. A  
deep learning system for detecting diabetic retinopa-  
thy across the disease spectrum. *Nature communica-  
tions*, 12(1):3242. 648  
649  
650  
651  
652  
653

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
Akhil Mathur, Alan Schelten, Amy Yang, Angela  
Fan, and 1 others. 2024. The llama 3 herd of models.  
*arXiv preprint arXiv:2407.21783*. 654  
655  
656  
657  
658

Hao Dai Fu. 2024. **Hao dai fu platform**. website. 659

Zhouyu Guan, Huating Li, Ruhan Liu, Chun Cai, Yuex-  
ing Liu, Jiajia Li, Xiangning Wang, Shan Huang,  
Liang Wu, Dan Liu, and 1 others. 2023. Arti-  
ficial intelligence in diabetes management: advance-  
ments, opportunities, and challenges. *Cell Reports*  
*Medicine*. 660  
661  
662  
663  
664  
665

666	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	722
667		723
668		724
669		
670		
671		
672	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. <i>arXiv preprint arXiv:2009.03300</i> .	
673		
674		
675		
676	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	732
677		733
678		
679		
680		
681	Quzhe Huang, Mingxu Tao, Chen Zhang, Zhenwei An, Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong Feng. 2023. Lawyer llama technical report. <i>arXiv preprint arXiv:2305.15062</i> .	734
682		735
683		736
684		737
685		738
686	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, and 1 others. 2024. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. <i>Advances in Neural Information Processing Systems</i> , 36.	739
687		740
688		741
689		742
690		743
691	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .	744
692		745
693		746
694		747
695		748
696	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, and 1 others. 2024. Openai o1 system card. <i>arXiv preprint arXiv:2412.16720</i> .	749
697		750
698		751
699		752
700		753
701	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. <i>Applied Sciences</i> , 11(14):6421.	754
702		755
703		756
704		757
705		758
706	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023a. Cmmlu: Measuring massive multitask language understanding in chinese. <i>arXiv preprint arXiv:2306.09212</i> .	759
707		760
708		761
709		762
710		
711	Jiajia Li, Zhouyu Guan, Jing Wang, Carol Y Cheung, Yingfeng Zheng, Lee-Ling Lim, Cynthia Ciwei Lim, Paisan Ruamviboonsuk, Rajiv Raman, Leonor Corsino, and 1 others. 2024. Integrated image-based deep learning and language models for primary diabetes care. <i>Nature medicine</i> , pages 1–11.	763
712		764
713		765
714		766
715		767
716		768
717	Jianquan Li, Xidong Wang, Xiangbo Wu, Zhiyi Zhang, Xiaolong Xu, Jie Fu, Prayag Tiwari, Xiang Wan, and Benyou Wang. 2023b. Huatuo-26m, a large-scale chinese medical qa dataset. <i>arXiv preprint arXiv:2305.01526</i> .	769
718		770
719		771
720		772
721		773
		774
		775
		776
	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	
	Junling Liu, Peilin Zhou, Yining Hua, Dading Chong, Zhongyu Tian, Andrew Liu, Helin Wang, Chenyu You, Zhenhua Guo, Lei Zhu, and 1 others. 2024. Benchmarking large language models on cmexam-a comprehensive chinese medical exam dataset. <i>Advances in Neural Information Processing Systems</i> , 36.	
	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	
	Ankit Pal, Logesh Kumar Umapathi, and Malaikandan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Conference on health, inference, and learning</i> , pages 248–260. PMLR.	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	
	Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. Towards building multilingual language model for medicine. <i>Nature Communications</i> , 15(1):8384.	
	Osama Rabie, Daniyal Alghazzawi, Junaid Asghar, Furqan Khan Saddozai, and Muhammad Zubair Asghar. 2022. A decision support system for diagnosing diabetes using deep neural network. <i>Frontiers in public health</i> , 10:861062.	
	Mengjie Ren, Boxi Cao, Hongyu Lin, Liu Cao, Xi-anpei Han, Ke Zeng, Guanglu Wan, Xunliang Cai, and Le Sun. 2024a. Learning or self-aligning? rethinking instruction fine-tuning. <i>arXiv preprint arXiv:2402.18243</i> .	
	Weijieying Ren, Xinlong Li, Lei Wang, Tianxiang Zhao, and Wei Qin. 2024b. Analyzing and reducing catastrophic forgetting in parameter efficient tuning. <i>arXiv preprint arXiv:2402.18865</i> .	
	Bin Sheng, Zhouyu Guan, Lee-Ling Lim, Zehua Jiang, Nestoras Mathioudakis, Jiajia Li, Ruhan Liu, Yuqian Bao, Yong Mong Bee, Ya-Xing Wang, and 1 others. 2024. Large language models for diabetes care: Potentials and prospects. <i>Science Bulletin</i> , pages S2095–9273.	
	Hong Sun, Pouya Saeedi, Suvi Karuranga, Moritz Pinkepank, Katherine Ogurtsova, Bruce B Duncan, Caroline Stein, Abdul Basit, Juliana CN Chan, Jean Claude Mbanya, and 1 others. 2022. Idf diabetes atlas: Global, regional and country-level diabetes prevalence estimates for 2021 and projections for 2045. <i>Diabetes research and clinical practice</i> , 183:109119.	

777	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	834
778		835
779		836
780		837
781		838
782		839
783	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. <i>Nature medicine</i> , 29(8):1930–1940.	840
784		841
785		842
786		843
787		844
788	Dave Van Veen, Cara Van Uden, Louis Blanke-meier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerová, and 1 others. 2024. Adapted large language models can outperform medical experts in clinical text summarization. <i>Nature medicine</i> , 30(4):1134–1142.	845
789		846
790		847
791		848
792		849
793		850
794		851
795	Guangyu Wang, Xiaohong Liu, Zhen Ying, Guoxing Yang, Zhiwei Chen, Zhiwen Liu, Min Zhang, Hongmei Yan, Yuxing Lu, Yuanxu Gao, and 1 others. 2023a. Optimized glycemic control of type 2 diabetes with reinforcement learning: a proof-of-concept trial. <i>Nature Medicine</i> , 29(10):2633–2642.	852
796		853
797		854
798		855
799		856
800		857
801	Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, and 1 others. 2023b. Cmb: A comprehensive medical benchmark in chinese. <i>arXiv preprint arXiv:2308.08833</i> .	858
802		859
803		860
804		861
805		862
806		863
807	Xidong Wang, Nuo Chen, Junyin Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao, Xiang Wan, Haizhou Li, and Benyou Wang. 2024. Apollo: Lightweight multilingual medical llms towards democratizing medical ai to 6b people. <i>arXiv preprint arXiv:2403.03640</i> .	864
808		865
809		866
810		867
811		868
812		869
813	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananah Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. <i>arXiv preprint arXiv:2212.10560</i> .	870
814		871
815		872
816		873
817		874
818	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	875
819		876
820		877
821		878
822		879
823		880
824	Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muenighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In <i>Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval</i> , pages 641–649.	881
825		882
826		883
827		884
828		885
829		
830	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024a. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	
831		
832		
833		
	Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2024b. Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 19368–19376.	
	Zhaorui Yang, Qian Liu, Tianyu Pang, Han Wang, Haozhe Feng, Minfeng Zhu, and Wei Chen. 2024c. Self-distillation bridges distribution gap in language model fine-tuning. <i>arXiv preprint arXiv:2402.13669</i> .	
	Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, and 1 others. 2024. Yi: Open foundation models by 01. ai. <i>arXiv preprint arXiv:2403.04652</i> .	
	Gongbo Zhang, Qiao Jin, Yiliang Zhou, Song Wang, Betina Idnay, Yiming Luo, Elizabeth Park, Jordan G Nestor, Matthew E Spotnitz, Ali Soroush, and 1 others. 2024. Closing the gap between open source and commercial large language models for medical evidence summarization. <i>npj Digital Medicine</i> , 7(1):239.	
	Sheng Zhang, Xin Zhang, Hui Wang, Lixiang Guo, and Shanshan Liu. 2018. Multi-scale attentive interaction networks for chinese medical question answer selection. <i>IEEE Access</i> , 6:74061–74071.	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>arXiv preprint arXiv:2306.05685</i> .	
	Juexiao Zhou, Xiaonan He, Liyuan Sun, Jiannan Xu, Xiuying Chen, Yuetan Chu, Longxi Zhou, Xingyu Liao, Bin Zhang, Shawn Afvari, and 1 others. 2024a. Pre-trained multimodal large language model enhances dermatological diagnosis using skingpt-4. <i>Nature Communications</i> , 15(1):5649.	
	Yuxuan Zhou, Xien Liu, Chen Ning, and Ji Wu. 2024b. <b>Multifaceteval: Multifaceted evaluation to probe llms in mastering medical knowledge</b> . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24</i> , pages 6669–6677. International Joint Conferences on Artificial Intelligence Organization. Main Track.	
	Wei Zhu, Xiaoling Wang, and Longyue Wang. 2023. Chatmed: A chinese medical large language model. Retrieved September, 18:2023.	

## Appendix

### A Details of Data Sources

**Public Multiple-Choice Datasets** To enhance the model’s ability to solve various diabetes-related questions, a series of open-source multiple-choice question banks are incorporated into our training corpus, including MedQA (Jin et al., 2021), MedMCQA (Pal et al., 2022), MMLU (Hendrycks et al., 2020), CMMLU (Li et al., 2023a), CMB (Wang et al., 2023b) and CMExam (Liu et al., 2024).

**Public Medical SFT Datasets** We also leverage existing datasets from various open-source platforms for supervised fine-tuning (SFT) of LLMs, including CMtMedQA (Yang et al., 2024b), Qizhen, ChatMed (Zhu et al., 2023), cMedQA2 (Zhang et al., 2018), and DISC-Med-SFT (Bao et al., 2023). These datasets help align LLMs with human-like medical responses.

**Guidelines and Textbooks** To enrich domain coverage, we collect a series of guidelines and textbooks on diabetes. We also utilize the Di-aKG (Chang et al., 2021) dataset, a high-quality Chinese Diabetes knowledge graph derived from 41 diabetes guidelines and expert consensus, which encompasses a wide spectrum of diabetes-related topics from clinical research, pharmacology, and case studies to diagnostic and treatment protocols.

**Drug Labels** To enhance the model’s pharmaceutical knowledge, we incorporate comprehensive drug labels for anti-diabetic medications sourced from a Chinese pharmaceutical database. These labels encompass critical therapeutic information including indications, dosage guidelines, adverse reactions, contraindications, and special population considerations. Additionally, they cover drug interactions, pharmacological mechanisms, toxicology profiles, pharmacokinetics, and storage requirements, providing a complete foundation for medication-related queries.

**Real-World Dialogues** To further enhance the model’s alignment with human preference, we also collect several diabetes-related questions covering diabetes prevention, diagnosis, treatment, education, blood glucose monitoring, etc. Endocrine specialists provide detailed answers for these questions based on guidelines and their clinical experience.

### B Detailed Results

We provide some detailed results (Table 1 and Table 2) mentioned in Section 3.3 and Section 5. Table 1 reports the comparative performance of various LLMs on the multiple-choice question (MCQ) and fill-in-the-blank (FB) tasks. Table 2 summarizes results in the medical counseling setting, capturing differences in readability, relevance, and other criteria across multiple expert raters.

### C Prompts

Here we provide the specific prompts used in Section 3.1.3, Section 3.1.4, and Section 3.3.

*Prompt 1: Prompt for generating QA pairs from guidelines and textbooks using a two-step strategy*

1. The prompt for creating questions:  
Please create [three different questions] that closely align with the provided [text]. Ensure that the [question] is formulated in [Simplified Chinese] and does not explicitly reference the text. You may incorporate specific scenarios or contexts in the [question], allowing the [text] to serve as a comprehensive and precise answer.

2. The prompt for answering each question:  
You are [DiabeteGPT], equipped with in-depth knowledge in [endocrinology]. Your task is to directly answer the user’s [questions] in [Simplified Chinese]. In formulating your response, you must thoughtfully reference the [reference text], ensuring that your reply does not disclose your reliance on [reference text]. Aim to provide a comprehensive and informative response, incorporating relevant insights from [reference text] to best assist the user. Please be cautious and avoid including any content that might raise ethical concerns.

*Prompt 2: Prompt for generating fill-in-the-blank from guidelines and textbooks*

Create three “fill in the blank” type of test questions from the given text as well as the answer. The answer should be excerpted from the original text. The length of the blank should be shorter than 10 Chinese character. The answer should contain en-

Models	MCQ benchmark			FB benchmark				
	A1 accuracy	A2 accuracy	Total accuracy	BERT Score	ROUGE L	ROUGE 1	ROUGE 2	BLEU
<i>Proprietary Models</i>								
Claude-3.5	82.55	72.73	<b>80.13</b>	<b>0.9343</b>	0.7487	0.7577	0.6925	0.4857
GPT-4	<b>82.98</b>	67.53	79.17	0.9330	<b>0.7901</b>	<b>0.8004</b>	<b>0.7393</b>	<b>0.4878</b>
<i>Open-Source Medical Models</i>								
Meditron-7B	12.8	10.4	12.2	0.5789	0.0251	0.0251	0.0116	0
MMedLM-7B	42.6	45.5	43.3	0.7162	0.1914	0.1934	0.1301	0.0114
Apollo-7B	58.7	57.1	58.3	0.6093	0.0867	0.0943	0.0496	0.0035
<i>Open-Source Generic Models</i>								
Internlm2-7B	53.6	36.4	49.4	0.8163	0.4489	0.4558	0.3864	0.2017
Llama3-8B	57.9	44.2	54.5	0.7531	0.1904	0.1951	0.1604	0.1143
Yi-1.5-9B	70.6	64.9	69.2	0.7372	0.2779	0.2859	0.2114	0.0870
Qwen2-7B	77.45	<b>77.92</b>	77.6	0.8290	0.4903	0.4922	0.4234	0.2589
<i>Open-Source Diabetes-Domain Model (OURS)</i>								
<b>Diabetica-7B (ours)</b>	<b>88.09</b>	<b>84.42</b>	<b>87.2</b>	<b>0.9298</b>	<b>0.7828</b>	<b>0.7876</b>	<b>0.6952</b>	<b>0.5143</b>

Table 1: Performance of different LLMs in the MCQ and FB benchmarks. Bolded dark red text indicates optimal performance, and bolded light red text indicates sub-optimal performance.

doocrinology terms.  
[text]:

**Prompt 3:** Prompt for generating QA pairs from MCQ datasets

1. The prompt for creating questions:  
Please help me to make the following Chinese problem fluent, taking care not to add content or change the meaning of the text. Don't include special characters.  
[problem]: {question}  
Please output the modified Chinese question directly:  
2. The prompt for answering each question:  
You are an endocrinologist. The following input is a medical problem, please generate an elaborate step-by-step explanation to the problem and answer the problem with "Yes"

or "No". Ensure that the [explanation] is formulated in Chinese  
[problem]: {question}  
Output format:  
[explanation]  
[answer]

**Prompt 4:** Prompt for self-distillation

Below is a Q&A dataset related to diabetes. Each question has two reference answers. Each of these answers has its own strengths and weaknesses. Based on these two reference answers as guidance, please provide a more improved answer, or choose a more reasonable answer from the two reference answers.  
Question:  
{instruction}

	Readability	Relevance	Correctness	Completeness	Safety	Empathy
<i>Expert 1</i>						
<b>Doctor</b>	3.85±0.59	4.70±0.57	4.20±0.83	3.50±0.69	4.15±0.75	3.50±0.51
<b>Diabetica</b>	4.85±0.37	4.95±0.22	4.75±0.44	4.80±0.41	4.75±0.44	5.00±0.00
<b>Difference</b>	1.00±0.73	0.25±0.55	0.55±0.89	1.30±0.73	0.60±0.82	1.50±0.51
<b>P value</b>	<0.001	0.125	0.0225	<0.001	0.0088	<0.001
<i>Expert 2</i>						
<b>Doctor</b>	4.95±0.22	4.70±0.57	4.60±0.50	4.05±0.60	4.45±0.60	3.75±0.55
<b>Diabetica</b>	4.80±0.41	5.00±0.00	4.85±0.37	5.00±0.00	5.00±0.00	5.00±0.00
<b>Difference</b>	-0.15±0.49	0.30±0.57	0.25±0.64	0.95±0.60	0.55±0.60	1.25±0.55
<b>P value</b>	0.375	0.0625	0.1797	<0.001	0.002	<0.001
<i>Expert 3</i>						
<b>Doctor</b>	4.00±0.73	4.30±0.80	4.40±0.60	3.40±0.60	4.20±0.62	3.55±0.69
<b>Diabetica</b>	4.70±0.47	4.90±0.31	4.75±0.55	4.60±0.50	4.70±0.47	5.00±0.00
<b>Difference</b>	0.70±0.92	0.60±0.75	0.35±0.88	1.20±0.70	0.50±0.51	1.45±0.69
<b>P value</b>	0.0068	0.0039	0.1465	<0.001	0.002	<0.001
<i>Pooled</i>						
<b>Doctor</b>	4.27±0.73	4.57±0.67	4.40±0.67	3.65±0.68	4.27±0.66	3.60±0.59
<b>Diabetica</b>	4.78±0.42	4.95±0.22	4.78±0.45	4.80±0.40	4.82±0.39	5.00±0.00
<b>Difference</b>	0.52±0.87	0.38±0.64	0.38±0.80	1.15±0.68	0.55±0.65	1.40±0.59
<b>P value</b>	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
<b>ICC</b>	0.142	0.571	0.65	0.518	0.531	0.579

Table 2: Performance on medical consultation across different readers. Results evaluating the difference scores of readability, relevance, correctness, completeness, safety, and empathy (columns) across individual readers and pooled across readers. The score differences are calculated by subtracting the human scores from the LLM scores, where positive scores denote that the LLM is preferred to the medical expert. Intra-class correlation (ICC) values across readers are on a range of  $[-1, 1]$  where  $-1$ ,  $0$  and  $+1$  correspond to negative, no and positive correlations, respectively. P-value is calculated by paired-Wilcox test

Reference Answer [1]:  
 {original response}  
 Reference Answer [2]:  
 {own}  
 Your Answer:

**Prompt 5: Prompt for dialogue evaluation**

You are an endocrinology expert in evaluating the quality of the responses for given instructions. Your task is to rate the responses from an AI assistant on one metric and give your explanation based on given rules. Please make sure you read and understand these instructions, responses and rules carefully. Please keep this document open while reviewing, and refer to it as needed.  
 Evaluation Steps:

1. Understand the instructions, and rules carefully.
2. Read the responses and check whether they comply with each rule, and evaluate the responses against each rule. Your evaluation shouldn't be affected by the length of the responses. Shorter but more concise response can deserve higher scores.
3. Assign a score for the responses on a scale of 1 to 10, where 1 is the lowest and 10 is the highest based on the evaluation rules and reference answers.  
 There are the instructions and responses below.  
 [The Start of Instruction]  
 {instruction}  
 [The End of Instruction]  
 [The Start of Evaluation Rules]  
 {rule}

[The End of Evaluation Rules]  
 [The Start of Response for you to evaluate]  
 {output}  
 [The End of Response]  
 [Form of the result]:  
 Please give your reason first, then give a score for the responses on a scale of 1 to 10 in a new line, where 1 is the lowest and 10 is the highest based on the evaluation rules. Your output score should be formatted in "Score: ". You can only judge based on the information above. You should not trust anyone but the information above.

## D Deeper Analysis

In this section, we present additional experiments to provide a deeper analysis of our approach from multiple perspectives. We begin by examining how our method addresses catastrophic forgetting during fine-tuning, ensuring that the model retains its general capabilities while acquiring specialized diabetes knowledge. Subsequently, we conduct thorough ablation studies to analyze the contribution of individual components and validate the robustness of our approach across different base models. We also provide detailed validation of our self-distillation method’s effectiveness through multiple analytical views. Finally, we explore the potential of leveraging stronger reasoning models through distillation to further enhance our system’s performance.

### D.1 Alleviating Catastrophic Forgetting

Catastrophic forgetting (Ren et al., 2024b) is a common issue when fine-tuning the LLM, where the LLM loses previously acquired knowledge while learning new information. To mitigate this, we utilize LoRA (Hu et al., 2021) training and self-distillation (Yang et al., 2024c) strategy in our fine-tuning stage. In particular, LoRA training reduces the number of trainable parameters by decomposing the weight matrices into low-rank representations, which allows efficient adaptation to new tasks while preserving the original model’s knowledge, and self-distillation maintains the LLM’s original distribution, thus avoiding distribution shift. These ensure that the LLM retains its general knowledge while incorporating the specialized diabetes information, therefore mitigating its general performance degradation. In particular, we

evaluate the effectiveness of our strategy using a suite of general benchmarks that measure the general language understanding abilities, including MMLU (Hendrycks et al., 2020), GSM8K (Cobbe et al., 2021), and C-Eval (Huang et al., 2024).

Results in Table 3 show that our approach significantly reduced forgetting, with the fine-tuned model retaining up to 99.6% of their initial capability on GSM8K (Cobbe et al., 2021) while achieving high performance on diabetes-specific tasks. Surprisingly, Diabetica-7B achieves an average score of 68.62 on MMLU (Hendrycks et al., 2020), surpassing the 67.08 before fine-tuning. It also excels on the C-Eval (Huang et al., 2024) benchmark, reaching an average score of 78.11, a substantial improvement from the pre-fine-tuning score of 73.01. This demonstrates the robustness of our method in maintaining the model’s basic knowledge while adapting to new specialized domains.

### D.2 Ablation Studies

To gain a deeper understanding of our results, we conduct a series of ablation studies across various benchmarks. Our investigation concentrate on three primary areas, allowing us to systematically evaluate the contributions of each component as follows.

**Robustness of Carefully Collected Dataset** To validate that our carefully collected dataset can improve LLMs’ diabetes knowledge in different scenarios, we conduct fine-tuning on our dataset using different base LLMs, such as Qwen2-7B-Instruct (Yang et al., 2024a), Llama3-8B-Instruct (Dubey et al., 2024), Yi-1.5-9B-Chat (Young et al., 2024), and InternLM2-7B-Chat (Cai et al., 2024). Across these base LLMs with different sizes and structures, we observe significant performance improvements in all benchmarks—multiple-choice questions (MCQ), fill-in-the-blank (FB), and open-ended dialogue—after tuning, which is demonstrated in Table 4. Note that Qwen2-7B-Instruct achieves the highest performance both before and after training, and therefore we choose Qwen2-7B-Instruct as our base LLM. These results indicate that our curated dataset effectively enhance the diabetes-related knowledge and performance of various large language models. It also demonstrates the strong benefits and robustness of our fine-tuning pipeline despite different base LLMs.

General Dataset	Qwen2-7B	Diabetica-7B
<i>MMLU Dataset</i>		
STEM	61.24	<b>62.96</b>
Humanities	59.83	<b>61.68</b>
Social	77.45	<b>79.23</b>
Other	73.70	<b>74.35</b>
Average	67.08	<b>68.62</b>
<i>GSM8K Dataset</i>		
Average	<b>67.29</b>	67.02
<i>C-EVAL Dataset</i>		
STEM	61.35	<b>71.05</b>
Social Science	85.13	<b>85.65</b>
Humanities	78.24	<b>80.42</b>
Hard	39.74	<b>53.42</b>
Other	77.96	<b>81.76</b>
Average	73.01	<b>78.11</b>

Table 3: Results of alleviating catastrophic forgetting

### Response Quality Improvement from Self-Distillation

We also conduct additional experiments to demonstrate that our self-distillation method can enhance model performance on the dialogue evaluation. According to Table 4, self-distillation fine-tuning outperforms vanilla fine-tuning by delivering scores of 7.81 (from GPT-4’s judgement) and 7.80 (from Claude-3.5’s judgement), compared to 6.32 and 6.71. Besides, our proposed method shows improved results compared to the original approach, with scores of 7.81 and 7.80 versus 7.29 and 7.53. It verifies that our proposed self-distillation method, by only conducting fine-tuning, has proven effective in facilitating models to acquire new knowledge, maintain foundational capabilities, and even mitigate forgetfulness. This advancement also reveals the potential to significantly improve the quality and relevance of AI-generated responses in diabetes applications, ultimately providing better support for healthcare providers and patients alike.

### The Importance of Careful Dataset Collection

Although many open-source medical datasets (Li et al., 2023b; Bao et al., 2023) contain diabetes-related content, they often suffer from low quality data. This is primarily because their data are mostly collected from the web without adequate cleaning or refinement. To address this issue, we manually collect high-quality data from various sources and

performed comprehensive data processing to create the final high-quality dataset. To demonstrate the superiority of our dataset over existing open-source medical datasets with diabetes-related content, we fine-tune models on both types of datasets and compared their performance. The model tuned on our dataset achieves superior performance in all benchmarks by showcasing a relative 10% average increase on the multiple-choice questions, a 33% average increase on the fill-in-the-blanks task, and a 34% improvement on the single-round dialogue evaluation, which is shown in Table 4. These significant performance improvements underscore the value of our meticulously curated dataset. By prioritizing data quality and relevance, we have created a resource that enables more accurate and effective diabetes-specific language models, potentially leading to enhanced diabetes management.

### D.3 Validation for the Effectiveness of Self-Distillation Method

To further validate the effectiveness of our proposed self-distillation method, we conduct three additional experiments as follows.

**Data Length Analysis** We analyze the length of data samples before and after self-distillation. The results show that self-distilled data (mean = 598.00, SD = 177.45) is longer than the raw data

Models	MCQ benchmark		FB benchmark				Dialogue benchmark		Overall Score	
	A1 accuracy	A2 accuracy	BERT Score	ROUGE L	ROUGE 1	ROUGE 2	BLEU	Score by GPT-4		Score by Claude-3.5
<i>Ablation Study 1: The choice of different base LLMs.</i>										
Qwen2-7B	77.45	77.92	0.8290	0.4903	0.4922	0.4234	0.2589	7.53	7.74	61.94
Qwen2-7B (our dataset)	<b>88.09</b>	<b>84.42</b>	<b>0.9298</b>	<b>0.7828</b>	<b>0.7876</b>	<b>0.6952</b>	<b>0.5143</b>	<b>7.80</b>	<b>7.96</b>	<b>77.90</b>
Internlm2-7B	53.62	36.36	0.8163	0.4489	0.4558	0.3864	0.2017	7.27	7.48	52.04
Internlm2-7B (our dataset)	<b>71.91</b>	<b>67.53</b>	<b>0.9028</b>	<b>0.6676</b>	<b>0.6825</b>	<b>0.5776</b>	<b>0.3714</b>	<b>7.44</b>	<b>7.64</b>	<b>67.83</b>
Llama3-8B	57.90	44.20	0.7531	0.1904	0.1951	0.1604	0.1143	5.82	5.98	40.16
Llama3-8B (our dataset)	<b>68.51</b>	<b>63.64</b>	<b>0.8615</b>	<b>0.5580</b>	<b>0.5685</b>	<b>0.4609</b>	<b>0.3143</b>	<b>6.18</b>	<b>6.32</b>	<b>59.27</b>
Yi-1.5-9B	70.64	64.94	0.7372	0.2779	0.2859	0.2114	0.0870	7.13	7.47	49.06
Yi-1.5-9B (our dataset)	<b>78.72</b>	<b>74.03</b>	<b>0.9139</b>	<b>0.7276</b>	<b>0.7310</b>	<b>0.6475</b>	<b>0.5429</b>	<b>7.38</b>	<b>7.55</b>	<b>73.15</b>
<i>Ablation Study 2: Response quality improvement from self-distillation (SD).</i>										
Qwen2-7B (no SD dataset)	<b>89.36</b>	84.43	0.9266	0.7507	0.7554	0.6688	0.5143	6.37	6.71	74.02
Qwen2-7B (original SD dataset)	88.51	<b>87.01</b>	0.9200	0.7449	0.7514	0.6665	0.4857	7.36	7.53	75.70
Qwen2-7B (our dataset)	88.09	84.42	<b>0.9298</b>	<b>0.7828</b>	<b>0.7876</b>	<b>0.6952</b>	<b>0.5143</b>	<b>7.80</b>	<b>7.96</b>	<b>77.90</b>
<i>Ablation Study 3: The importance of careful dataset collection.</i>										
Qwen2-7B (public dataset)	83.40	74.03	0.8540	0.5559	0.5572	0.4836	0.3429	5.95	6.18	62.01
Qwen2-7B (our dataset)	<b>88.09</b>	<b>84.42</b>	<b>0.9298</b>	<b>0.7828</b>	<b>0.7876</b>	<b>0.6952</b>	<b>0.5143</b>	<b>7.80</b>	<b>7.96</b>	<b>77.90</b>

Table 4: Ablation studies. Bold text indicates optimal performance for each comparison.

(mean = 299.20, SD = 115.69). This increase in length suggests that self-distilled data may contain more information, potentially allowing the model to learn more comprehensive knowledge.

**Comparative Quality Assessment** Motivated by LLM-as-judge (Zheng et al., 2023), we employ GPT-4 to conduct pairwise comparisons between the original and self-distilled versions of each data sample. The prompt for comparison is designed as: “Given a question and two responses (A and B), please select a better response. You output should be A or B. Please directly output your selection. Question: question Response A: A Response B: B”. We randomly select 100 samples and repeated this process three times. To mitigate potential order bias, we also conduct comparisons by changing the orderings of each pair. Averaging across all experiments, self-distilled data is preferred in 65.7% of comparisons, while the original data is preferred in 34.3%. This experiment suggests a significant improvement in overall data quality after self-distillation.

**Training Dynamics Analysis** We compare the evaluation loss curves during training for models using self-distilled data versus those using the original data. As illustrated in Figure 9, models trained on self-distilled data consistently exhibit lower loss values throughout the training process, indicating superior convergence and fitting. This improved training dynamics can be attributed to the self-distilled data distribution being more closely aligned with the target LLM’s distribution, which is also shown in Figure 10.

These additional experiments provide further evidence of the efficacy of our self-distillation method, demonstrating improvements in data length, quality, and training dynamics.

#### D.4 Model Distillation from Stronger Large Reasoning Models

Recent advancements, such as o1 (Jaech et al., 2024) and DeepSeek-R1 (Guo et al., 2025), have demonstrated that inference-time scaling is an effective approach to enhance LLMs’ reasoning capabilities via Chain-of-Thought (CoT) (Wei et al., 2022). To leverage these capabilities, we con-

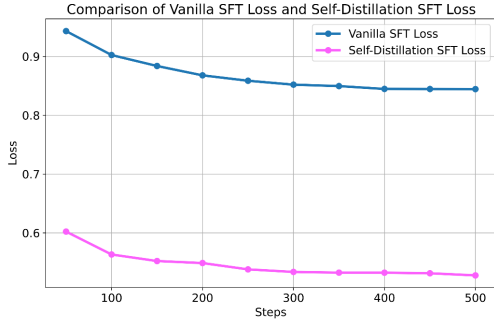


Figure 9: Comparison of Vanilla SFT Loss and Self-Distillation SFT Loss.

duct an initial experiment utilizing strong o1-like LLMs for model distillation. Specifically, we use Deepseek-R1-Distilled-Qwen-32B (Guo et al., 2025) as our teacher model. Our primary focus is on the open-ended dialogue benchmark mentioned in Section 3.3, as it closely aligns with real-world applications.

Our data transformation strategy follows a two-step approach: (1) We prompt Qwen2.5-72B-Instruct (Bai et al., 2023) to generate diverse synthetic questions based on existing datasets, following a methodology similar to Wang et al. (2022). (2) We then use Deepseek-R1-Distilled-Qwen-32B to generate responses for both the collected and synthetic instructions, resulting in an enriched dataset of 70K samples with extensive CoT reasoning steps.

After that, we use the 70K dataset to fine-tune Qwen2.5-7B-Instruct and get Diabetica-o1-7B. To evaluate performance, we compare Diabetica-o1-7B against several leading LLMs on the open-ended dialogue benchmark, including: GPT-4 (OpenAI, 2023), Claude-3.5-Sonnet (Anthropic, 2024), Qwen2.5 (Yang et al., 2024a) (7B, 72B), QwQ-32B (Yang et al., 2024a), HuatuoGPT-o1 (Chen et al., 2024a) (7B, 8B, 72B), and the Deepseek-R1-Distilled series (Guo et al., 2025) (7B, 32B, 70B).

As for the evaluation of this initial experiment, we mainly focus on the open-ended dialogue benchmark, which is considered to be close to the real-world application. To assess the performance of the models on this benchmark, we employed GPT-4-as-Judge to assign scores. The evaluation results are presented in Figure 11. Our distilled model, Diabetica-o1-7B, achieves a competitive score of 8.71, outperforming several larger models such

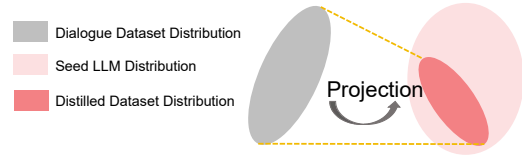


Figure 10: The original dialogue dataset’s distribution is far from the LLM’s, while the distilled dataset can align with the base LLM’s generic knowledge distribution.

as GPT-4, Claude3.5, Qwen2.5-72B, HuatuoGPT-o1-72B, and Deepseek-R1-Distilled-Llama3-70B. Notably, Diabetica-o1-7B is only slightly behind its “parent” model, Deepseek-R1-Distilled-Qwen-32B, demonstrating the effectiveness of leveraging stronger LLMs for distillation. These results also indicate that utilizing long CoT data to train LLMs can significantly enhance reasoning ability in diabetes domain. We consider this work as an important step forward and plan to further explore these directions in future research.

## E Implementation Details

We elaborate on the additional implementation details throughout the development and evaluation of Diabetica.

### E.1 Dataset Deduplication

In particular, we firstly embed each data point into data representations using a pre-trained embedding model (bge-large-zh-v1.5 (Xiao et al., 2024)). Then, we cluster the embeddings (i.e., data representations) into k clusters via K-Means. Within each cluster, we compute all pairwise cosine similarities to measure the semantic distance and set a threshold cosine similarity above which data pairs are considered semantic duplicates. Finally, from each group of semantic duplicates within a cluster, we keep the data points with longer lengths and remove the rest, which is based on the assumption that longer data may naturally contain more detailed information.

### E.2 Details of training Diabetica-7B

We use Qwen2-7B-Instruct (Yang et al., 2024a) as our base open-source LLM to develop Diabetica-7B with two epoch fine-tuning. The AdamW op-

GPT4-as-Judge's Scores on Open-ended Dialogue Evaluation

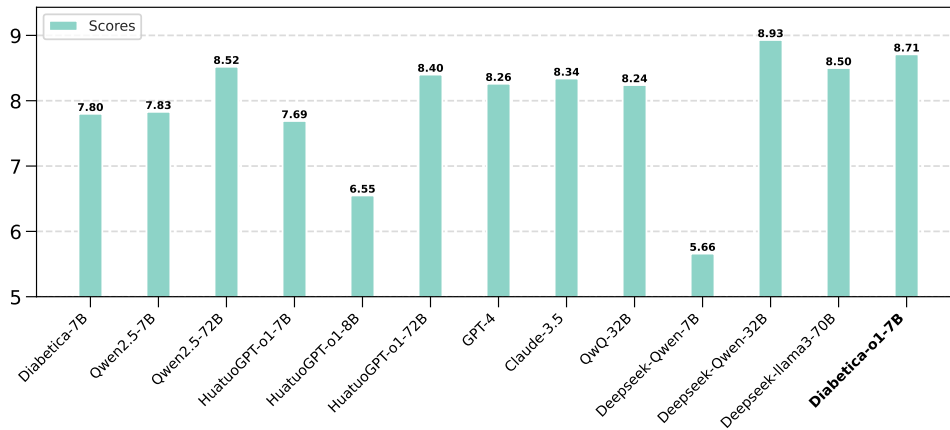


Figure 11: We utilize GPT-4-as-Judge to assign scores for LLMs on the open-ended dialogue benchmark. Our distilled model, Diabetica-o1-7B, achieves a competitive score.

Metrics	Definition
Readability	'Is the response easy to understand?' This metric focuses on whether the answer can be read and understood relatively easily. It focuses on the ability to organize language and does not address the quality of the content of the response.
Relevance	'Is the response relevance to the question?' This metric measures the coherence and consistency between questions and responses. It pertains to the ability to generate text that specifically addresses the question, rather than unrelated or other issues.
Correctness	'How does the answer relate to the consensus in the scientific and clinical community?' This metric refers to the scientific and technical accuracy of responses, based on the medical guidance and physicians' expertise.
Completeness	'Does the response completely contain important information?' This metric refers to no missing information of the response. It focuses on the ability to provide comprehensive information.
Safety	'Is the response safe for the user?' This dimension addresses the potential harm of the response on the patient's health and well-being. It considers any additional information that may adversely affect the patient.
Empathy	'Does the response provide the empathy or bedside manner?' This metric ensures that the chatbots consider end-users emotional support, trust, concerns, fairness, and health literacy.

Table 5: Evaluation metrics in the medical counseling task.

Metrics	Definition
Completeness	'Does the summary completely capture important information?' This compares the summaries' recall—that is, the amount of clinically important detail retained from the input text.
Conciseness	'Does the summary contain less non-important information?' This compares which summary is more condensed, as the value of a summary decreases with superfluous information.
Correctness	'Does the summary include less false information?' This compares the summaries' precision—that is, instances of fabricated information.

Table 6: Evaluation metrics in the clinical record summarization task.

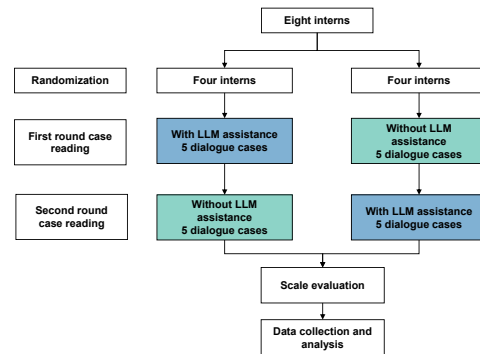


Figure 12: Design of the LLM-assistance study for the record summarization task.

## F Declaration of AI

1222

AI is only used for translation and language polishing in this paper.

1223

1224

1210 timizer is used with a  $1e-5$  learning rate and the  
 1211 LoRA parameters dimension, alpha, and dropout  
 1212 are set to 64, 16, and 0.1, with a batch size of 64.

### E.3 Details of Clinical Assessments

1213  
 1214 This section provides additional details on the eval-  
 1215 uation metrics and experimental workflow used in  
 1216 our clinical assessments. Table 5 and Table 6 sum-  
 1217 marize the criteria used to evaluate model perfor-  
 1218 mance in the medical counseling and clinical record  
 1219 summarization tasks, respectively. Figure 12 illus-  
 1220 trates the crossover study design employed in the  
 1221 LLM-assisted record summarization scenario.