

# ClinDoc Copilot: Enhancing Efficiency and Normalization in Chinese Outpatient Settings

Anonymous ACL submission

## Abstract

Clinical documentation in Chinese outpatient settings is a legally regulated and time-intensive task, requiring physicians to produce records that are both efficient and auditable. Existing medical large language models demonstrate strong performance on medical benchmarks. However, their generated content cannot be directly adopted in clinical documentation due to legal risk, as responsibility for medical records cannot be delegated to or assumed by autonomous models. We present ClinDoc Copilot, a physician-centered assistant designed to support outpatient documentation. We evaluate ClinDoc Copilot in realistic simulated outpatient scenarios with licensed physicians and clinical trainees. Results show that the system improves documentation efficiency, terminology standardization, and evidence-grounded traceability. Our work highlights the value of interactive, physician-led AI assistance for clinical documentation in regulated healthcare settings. All code is available at an anonymous repository: <https://anonymous.4open.science/r/ClinDoc-Copilot/README.md>.

## 1 Introduction

In Chinese outpatient settings, clinical documentation is a strictly regulated professional activity. Outpatient medical records are **legally binding documents** subject to routine audits by the national health insurance reimbursement system. Physicians bear full legal responsibility for the accuracy, completeness, and clinical validity of all documented content. As a result, every diagnosis and treatment order must be explicitly supported by documented evidence and remain traceable throughout the record.

At the same time, outpatient documentation is completed under **severe time constraints**. Physicians are expected to conduct patient interviews, perform clinical reasoning, and finalize medical records within a limited consultation window

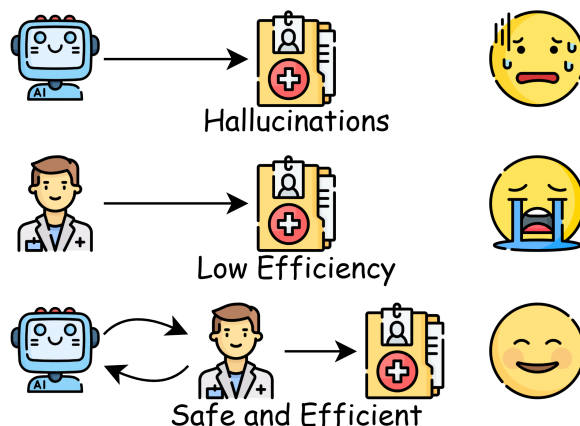


Figure 1: Comparison of Different HCI Modes in Chinese Clinical Documentation.

(Shanafelt et al., 2016). The coexistence of documentation rigor and time pressure places substantial cognitive burden on physicians, particularly during real-time writing.

Motivated by current situations, we introduce **ClinDoc Copilot**, a physician-centered interactive assistant for outpatient clinical documentation. ClinDoc Copilot operates in a human-in-the-loop manner, providing real-time, non-binding assistance during the writing process while ensuring that the physician remain the sole legal owner of medical records. By prioritizing physician control and accountability, the system aims to improve documentation efficiency and writing quality without compromising regulatory requirements.

We evaluate ClinDoc Copilot in realistic simulated outpatient scenarios involving both licensed physicians and clinical trainees. Experimental results show that the system substantially reduces documentation effort while improving terminology standardization and evidence-grounded traceability of diagnoses and treatment orders. A user study further indicates that clinicians perceive ClinDoc Copilot as practical, trustworthy, and well aligned with real outpatient documentation workflows.

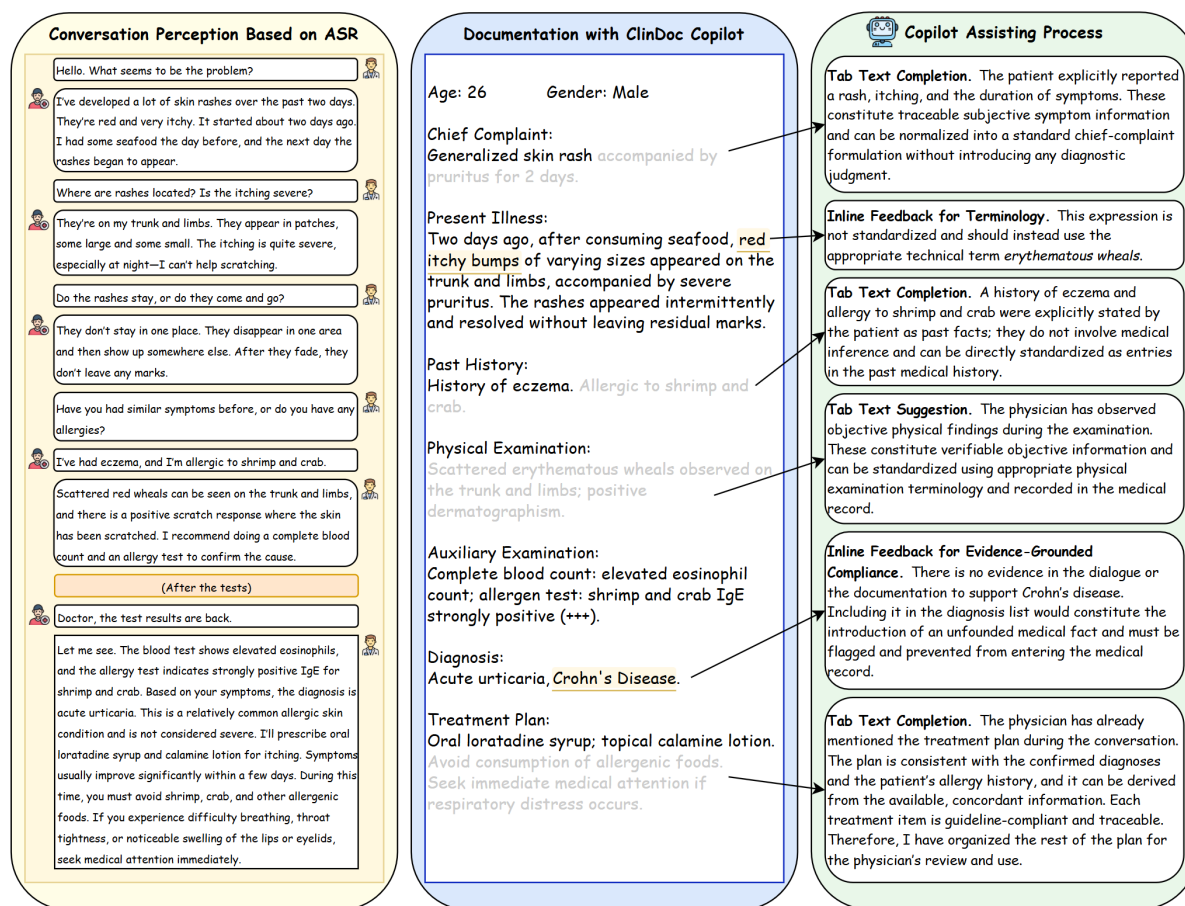


Figure 2: Overview of the ClinDoc Copilot Workflow. Black text represents content entered by the physician, while grey text indicates suggested completions. For readers' convenience, the illustration is in English.

## 2 Background

This section provides background for the design of ClinDoc Copilot. Section 2.1 discusses why existing medical large language models are unsuitable. Section 2.2 summarizes key user needs in Chinese outpatient documentation. Next, section 2.3 reviews existing documentation assistance systems and their limitations. Finally, section 2.4 outlines principles from responsible and human-centered AI that inform the system design.

### 2.1 Limitations of Medical LLMs for Clinical Documentation

Recent advances in medical large language models demonstrate near-human or even super-human performance on exam-style benchmarks (Singhal et al., 2025; Kung et al., 2023; Pal et al., 2022; Xie et al., 2025). Such capabilities make LLMs effective for medical education (Gilson et al., 2023), knowledge dissemination (Thirunavukarasu et al., 2023), and general-purpose decision support (Oniani et al., 2024). However, clinical documenta-

tion poses fundamentally different requirements. Beyond medical correctness, documentation demands efficiency under time pressure, normalized medical expression, and clear attribution of clinical responsibility. Recent studies highlight that benchmark-driven LLMs remain insufficient for legally accountable documentation, due to risks such as hallucination, reasoning errors, and misalignment with real clinical workflows (Nori et al., 2023; Asgari et al., 2025; Vladika et al., 2025).

### 2.2 User Needs in Chinese Outpatient Documentation

**Interviews.** To better understand real-world outpatient documentation challenges, we conducted semi-structured interviews with licensed physicians across both modern Western medicine and traditional Chinese medicine (TCM). Despite differences in clinical experience and medical systems, participants reported a highly consistent set of documentation pain points, with these issues being particularly pronounced among trainees.

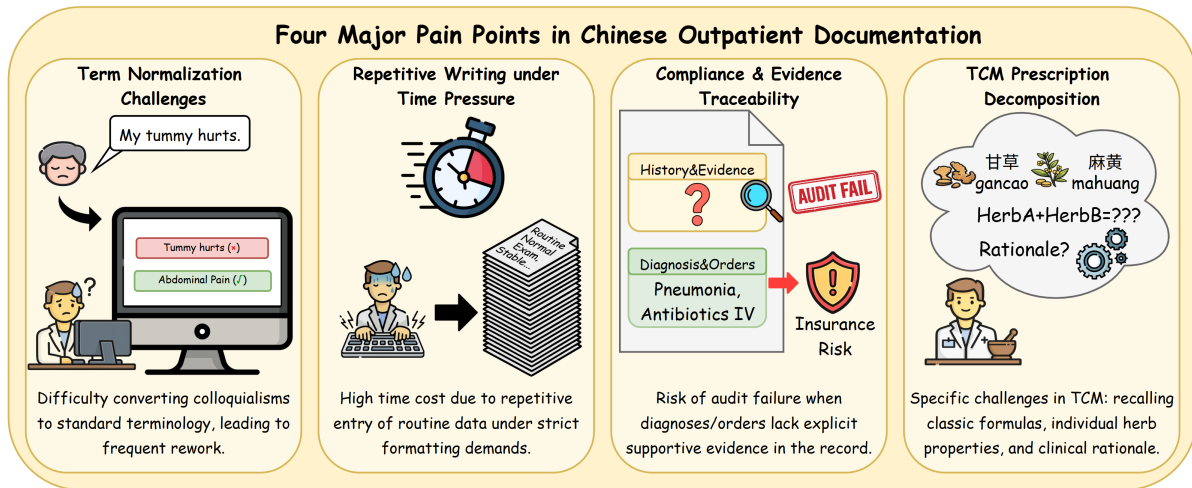


Figure 3: Four major pain points in Chinese Outpatient Documentation

**Pain points.** Four recurring pain points in Chinese outpatient clinical documentation that hinder efficiency and compliance are identified:

- **Terminology normalization.** Participants frequently reported difficulties in translating patients' colloquial descriptions into standardized medical terminology, resulting in records that require substantial post hoc revision.
- **Repetitive writing under time pressure.** Physicians emphasized the significant time cost associated with repeatedly drafting routine documentation sections while adhering to strict formatting requirements.
- **Compliance and evidence-grounded traceability.** Many trainees struggled to ensure documentation compliance, particularly the traceability requirement of diagnoses and treatment orders. This requirement is closely tied to national health insurance auditing practices.
- **TCM prescription decomposition.** Physicians practicing modern traditional Chinese medicine highlighted domain-specific challenges in prescription writing, including recalling classical formulas, individual herbal functions, and the clinical rationale underlying herb combinations.

Taken together, these pain points reveal two fundamental needs in Chinese outpatient documentation: improving writing efficiency by reducing repetitive manual effort, and enforcing writing norms that ensure terminology accuracy and evidence-grounded traceability.

## 2.3 Existing Systems

### 2.3.1 Automated Generation

A substantial body of prior work has framed clinical documentation as an automatic text generation or summarization task. Representative approaches include dialogue-to-SOAP note generation pipelines (Krishna et al., 2021; Ben Abacha et al., 2023). Subsequent studies have explored modular summarization, retrieval augmentation, and large language models for generating summaries and discharge notes from clinical dialogues or transcripts (Giorgi et al., 2023; Xie et al., 2023).

While these systems demonstrate promising performance on benchmarks, they typically emphasize autonomous note generation. Such designs implicitly treat documentation as an output artifact rather than a physician-controlled legal record, making them less suitable for outpatient settings where clinicians remain fully accountable for every documented statement. In regulated environments, silently generated content may introduce unacceptable risks when diagnoses or treatment decisions are insufficiently grounded in explicitly documented evidence.

### 2.3.2 Commercial Systems

In parallel with academic research, commercial systems based on ambient clinical intelligence have been deployed to support clinical documentation (Healthcare, 2025). Despite demonstrated efficiency benefits, these commercial systems are proprietary, English-centric, and optimized for healthcare systems with documentation norms that differ from Chinese outpatient practice.

175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
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  
219  
220  
221  
222  
223  
224  
225

## 2.4 Responsible AI in Healthcare

A growing body of research has emphasized that deploying AI systems in healthcare requires not only high predictive or generative performance, but also careful consideration of responsibility, accountability, and human oversight. Unlike many consumer-facing AI applications, clinical systems operate in safety-critical environments where errors may carry legal, ethical, and clinical consequences, and where responsibility for decisions cannot be delegated to opaque automated systems (Wiens et al., 2019; Sendak et al., 2020).

Prior work in responsible and human-centered medical AI consistently highlights several recurring principles. First, clinical AI systems should preserve clinician agency and decision authority, ensuring that human professionals remain accountable for all medically relevant actions and documentation (Amershi et al., 2019). Second, automation should be carefully scoped to avoid overreliance, silent failure modes, or the unintentional introduction of unsupported or fabricated information into clinical workflows (Bansal et al., 2021; Nori et al., 2023). Third, AI assistance must be aligned with existing clinical workflows, regulatory constraints, and auditing practices, rather than optimizing solely for technical performance metrics (Sculley et al., 2015).

These considerations are particularly salient for clinical documentation, which serves not only as a communication artifact but also as a legally binding record subject to post hoc review and insurance auditing. In such contexts, even fluent and factually plausible AI-generated text may be problematic if it cannot be clearly traced to documented clinical evidence or explicit physician intent.

Taken together, this line of work motivates a design stance in which AI systems support, rather than replace, clinician-led documentation. For ClinDoc Copilot, these insights translate into concrete system requirements: AI assistance must be non-binding, explicitly invoked, and grounded in the encounter context; any content incorporated into the medical record must result from deliberate physician action; and the system must avoid introducing new clinical facts or decisions beyond those already established by the clinician. By operationalizing these responsible AI principles at the interaction and workflow level, ClinDoc Copilot aims to reconcile efficiency gains with the accountability demands of regulated outpatient clinical practice.

## 3 System Overview

ClinDoc Copilot is a physician-centered, human-in-the-loop assistant designed to support outpatient clinical documentation under strict regulatory constraints in Chinese clinical settings. Rather than autonomously generating medical records, the system assists physicians during documentation writing by providing optional, real-time support grounded in the evolving conversation and draft.

### 3.1 Design Goals

Based on interviews with licensed physicians and clinical trainees, we identify two primary design goals for ClinDoc Copilot.

First, the system should **improve documentation efficiency** by reducing repetitive and time-consuming manual writing. In particular, routine documentation sections that follow stable structural patterns should be completed with minimal friction, allowing physicians to focus on clinically meaningful decisions including medical orders.

Second, the system should **support documentation norms required** in regulated clinical practice. This includes helping ensure that medical terminology is standardized and that diagnoses and treatment orders are supported by, and derivable from, information documented earlier in the record. Such evidence-grounded traceability is especially critical in the context of national health insurance auditing. Importantly, both goals must be achieved without weakening physician authority or accountability. ClinDoc Copilot is designed as a supportive tool rather than an autonomous author, with all assistance remaining optional and under explicit physician control.

### 3.2 Core Components

ClinDoc Copilot consists of three integrated components. Together, they address both efficiency and writing norm requirements while preserving physician authorship and responsibility.

**Conversation-Grounded Ghost Text and Tab Completion.** ClinDoc Copilot supports fast drafting through ghost text preview and explicit tab-based completion. During outpatient encounters, the system relies on audio recording and transcription to construct a conversational context. This context, together with the surrounding draft text, is used to generate ghost text suggestions that appear as low-contrast inline previews.

226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273

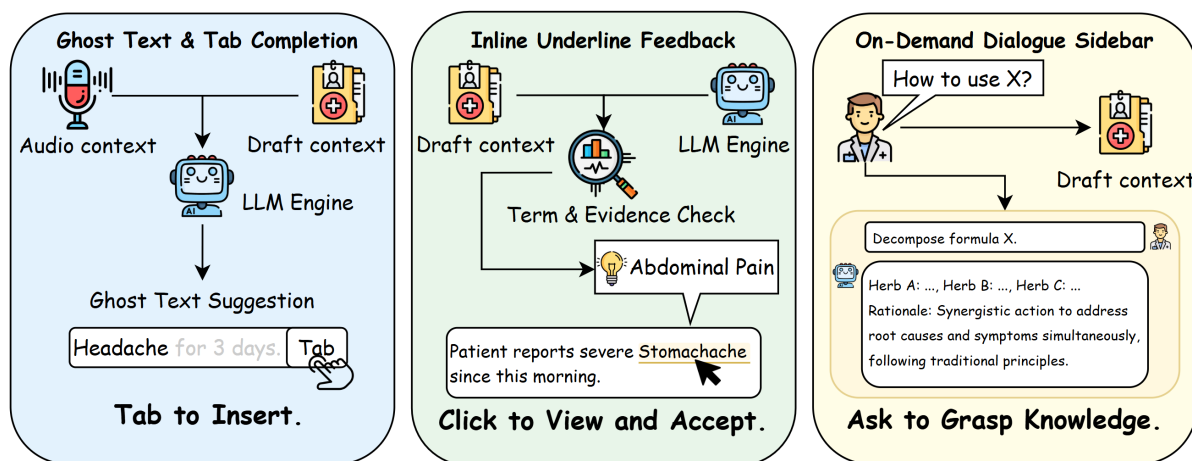


Figure 4: Overview of the ClinDoc Copilot Components.

274 The ghost text **can adapt dynamically to user**  
 275 **input**. And physicians may explicitly accept a  
 276 suggestion by pressing the tab key, at which point  
 277 the preview text is inserted and becomes fully edi-  
 278 table. If a suggestion is ignored, it disappears with-  
 279 out affecting the draft. In this way, conversational  
 280 grounding enables relevant suggestions, while the  
 281 ghost-text-and-tab interaction ensures that insertion  
 282 remains intentional and physician-controlled.

283 **Inline Underline Feedback for Terminology and**  
 284 **Evidence-Grounded Compliance.** To support  
 285 writing norms, ClinDoc Copilot provides inline  
 286 feedback using visual underlines. When colloquial  
 287 expressions in the physician’s draft should be nor-  
 288 malized into standard medical terminology, the sys-  
 289 tem highlights the relevant text spans with a yel-  
 290 low underline. Upon clicking for details, ClinDoc  
 291 Copilot presents localized revision suggestions that  
 292 illustrate how the highlighted expression can be  
 293 rewritten using appropriate clinical terminology.  
 294 These suggestions are scoped to the underlined  
 295 fragment and do not trigger automatic rewriting,  
 296 allowing physicians to decide whether and how to  
 297 apply the proposed changes.

298 In addition, when diagnoses or treatment orders  
 299 appear insufficiently supported by information docu-  
 300 mented earlier in the record, the system similarly  
 301 applies underline cues to prompt physician review.  
 302 These inline underline feedbacks are triggered not  
 303 only by non-standard terminology, but also by in-  
 304 sufficient logical support within the documentation.  
 305 In all cases, the underline serves as a reminder  
 306 rather than an enforcement mechanism, preserving  
 307 physician-led correction and decision-making.

308 **On-Demand Right-Side AI Dialogue Sidebar.**  
 309 ClinDoc Copilot includes an AI dialogue sidebar  
 310 that appears only when the physician explicitly  
 311 invokes it. This design reflects the clinical require-  
 312 ment that assistance must be physician-initiated  
 313 and non-intrusive. When invoked, the sidebar sup-  
 314 ports clarification and just-in-time assistance dur-  
 315 ing documentation. Typical use cases include: ask-  
 316 ing how to properly phrase a section in compliant  
 317 clinical language, requesting terminology clarifi-  
 318 cation, or, in modern traditional Chinese medicine  
 319 practice, querying prescription decomposition such  
 320 as rationale for herb combinations. Outputs in the  
 321 sidebar are advisory. They do not automatically  
 322 insert content into the medical record and must be  
 323 manually incorporated by the physician if desired.

### 3.3 Implementation Details 324

325 ClinDoc Copilot adopts a modular system archi-  
 326 tecture that cleanly separates speech understand-  
 327 ing, language modeling, and interaction logic.  
 328 For speech recognition, the system employs SenseVoice (An et al., 2024) to transcribe outpatient  
 329 doctor-patient conversations in real time. For  
 330 language understanding and generation, ClinDoc  
 331 Copilot uses GPT-4o (OpenAI, 2024) as the back-  
 332 end large language model engine. 333

### 3.4 Accountability and Control 334

335 Preserving clear responsibility boundaries is a cen-  
 336 tral design principle of ClinDoc Copilot. The sys-  
 337 tem does not introduce new clinical facts as com-  
 338 mitted record content, infer diagnoses, or modify  
 339 treatment decisions. By relying on explicit user ac-  
 340 tions for any insertion and by keeping the dialogue  
 341 interface physician-invoked, ClinDoc Copilot en-

342 sures that physicians remain the sole authors and  
 343 legal owners of medical records. This design re-  
 344 duces the risk of uncontrollable hallucination being  
 345 silently committed into documentation and sup-  
 346 ports accountable, auditable practices aligned with  
 347 outpatient regulatory requirements.

## 348 4 Evaluation

349 We evaluate ClinDoc Copilot in realistic simu-  
 350 lated outpatient documentation scenarios. The  
 351 evaluation addresses three research questions: (1)  
 352 whether ClinDoc Copilot improves documentation  
 353 efficiency, (2) whether it improves documentation  
 354 quality and evidence-grounded traceability without  
 355 introducing hallucination risks, and (3) whether  
 356 clinicians perceive the system as usable and suit-  
 357 able for real clinical practice.

### 358 4.1 Experimental Setup

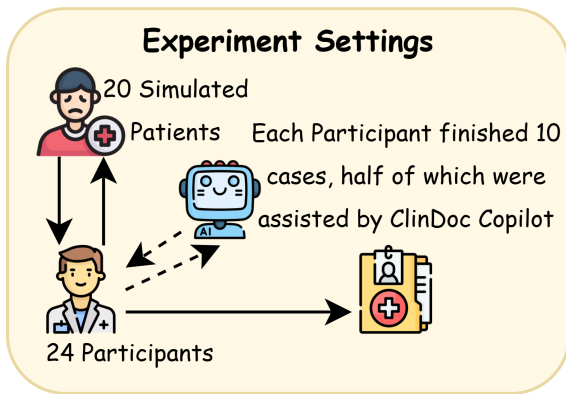


Figure 5: Overview of Evaluation Settings.

359 **Participants.** We recruited a total of 24 partici-  
 360 pants, including 12 licensed outpatient physicians  
 361 and 12 clinical trainees. All physicians had compa-  
 362 rable outpatient experience within their respective  
 363 clinical systems, and all trainees were at a similar  
 364 stage of clinical training. Participants covered two  
 365 independent clinical systems: Western medicine  
 366 (WM) and traditional Chinese medicine (TCM).

367 **Tasks and Scenarios.** The evaluation was con-  
 368 ducted using realistic simulated outpatient encoun-  
 369 ters derived from twenty de-identified clinical cases  
 370 (zhangmingme, 2025; SylvanL, 2025). All use data  
 371 are in Chinese. Ten cases were from WM system  
 372 and another ten were from TCM system. For each  
 373 case, a trained actor played the role of the patient  
 374 following a standardized script, ensuring every par-  
 375 ticipant would receive identical information. Each

376 participant completed documentation for ten outpa-  
 377 tient encounters within their corresponding clinical  
 378 system. Because simulated patients cannot fully re-  
 379 produce all clinical details, a controlled mechanism  
 380 was adopted for auxiliary examinations. When a  
 381 physician requested laboratory tests or imaging  
 382 studies during the encounter, the simulated patient  
 383 provided the corresponding results if such examina-  
 384 tions were present in the original medical record. If  
 385 the requested examination did not exist in the real  
 386 case, the physician was explicitly informed that the  
 387 examination could not be performed.

388 **Conditions.** Participants were randomly as-  
 389 signed five cases under **Copilot** condition, in which  
 390 ClinDoc Copilot was enabled, and five cases un-  
 391 der **Baseline** condition, in which participants used  
 392 the same documentation interface without Copilot  
 393 assistance. This fully randomized assignment miti-  
 394 gates potential confounding effects arising from  
 395 fixed case-condition pairing, ensuring that ob-  
 396 served differences are not attributable to particu-  
 397 lar case characteristics. For each participant, the  
 398 ten documentation tasks were arranged in an inter-  
 399 leaved manner, alternating between the Copilot and  
 400 Baseline conditions rather than being completed  
 401 in separate blocks. This cross-condition ordering  
 402 helps reduce fatigue- and learning-related biases  
 403 that could otherwise skew efficiency and quality  
 404 measurements. This setting resulted in a balanced  
 405  $2 \times 2 \times 2$  design. Table 1 summarizes the participant  
 406 groups and task distribution.

Group	Cond.	#Part.	Cases
WM / P.	Cop.	6	30
WM / T.	Cop.	6	30
TCM / P.	Cop.	6	30
TCM / T.	Cop.	6	30
WM / P.	Base.	6	30
WM / T.	Base.	6	30
TCM / P.	Base.	6	30
TCM / T.	Base.	6	30

Table 1: Participant groups and documentation tasks. P. denotes Physician; T. denotes Trainee; Cop. denotes Copilot condition; Base. denotes Baseline condition; Cond. denotes Condition; #Parts. denotes number of participants.

407 **Ethical Considerations.** All participants pro-  
 408 vided informed consent prior to the study. Audio  
 409 recordings were collected solely for documenta-  
 410 tion assistance in simulated scenarios and did not  
 411 involve real patient care. No identifiable patient  
 412 information was used or stored, and all case materi-

Group	Cond.	Time↓	Keys↓	Edit↓	TN↑	TR↑	SH↑
WM / P.	Cop.	485.50 (68.5)	45.50 (12.5)	0.22 (0.05)	4.75 (0.42)	4.68 (0.45)	4.92 (0.28)
WM / T.	Cop.	530.00 (75.2)	68.00 (15.8)	0.18 (0.06)	4.62 (0.48)	4.55 (0.50)	4.83 (0.38)
TCM / P.	Cop.	550.50 (72.0)	52.50 (14.0)	0.25 (0.07)	4.70 (0.45)	4.65 (0.48)	4.96 (0.20)
TCM / T.	Cop.	610.00 (85.5)	85.00 (18.5)	0.15 (0.05)	4.58 (0.52)	4.42 (0.55)	4.79 (0.41)
WM / P.	Base.	618.75 (58.4)	88.25 (12.5)	–	4.38 (0.49)	4.46 (0.51)	–
WM / T.	Base.	793.33 (75.6)	123.67 (18.2)	–	3.67 (0.61)	3.54 (0.66)	–
TCM / P.	Base.	715.00 (68.2)	110.00 (15.3)	–	4.13 (0.50)	4.21 (0.51)	–
TCM / T.	Base.	920.00 (85.4)	145.50 (20.1)	–	3.21 (0.59)	3.13 (0.61)	–
Overall Effect Size	d	1.94	2.23	–	1.35	1.12	–
Significance	p	<0.001	<0.001	–	<0.001	<0.001	–

Table 2: Efficiency metrics averaged per medical record and rubric-based documentation quality and safety scores (1–5). Clinical cases were randomly assigned to experimental conditions using a counterbalanced design.

als were derived from de-identified clinical records. All collected data were used exclusively for research purposes and handled in accordance with institutional data protection practices.

## 4.2 Evaluation Metrics

**Efficiency Metrics.** Documentation efficiency was evaluated using the following metrics, all reported as averages per medical record:

- **Time (↓):** average seconds needed to complete one outpatient record.
- **Keys (↓):** average number of keyboard inputs per record.
- **Edit (↓):** average ghost text edit rate, defined as the proportion of tab-accepted ghost text suggestions that were subsequently modified.

**Quality and Safety Metrics.** Documentation quality was assessed using three rubric-based metrics:

- **TN (↑):** average Terminology Normalization.
- **TR (↑):** average Evidence-Grounded Traceability of diagnoses and treatment orders.
- **SH (↑):** average Safety and Hallucination level, with higher scores indicating fewer hallucination issues.

All quality metrics were scored on a five-point scale by two senior physicians in a double-blind manner. Detailed rubric definitions are provided in Appendix B.

## 4.3 Statistical Analysis and Overall Effects

To assess the overall impact of ClinDoc Copilot on documentation efficiency and record quality,

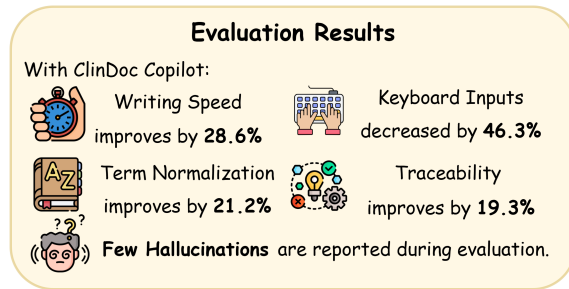


Figure 6: Overview of Evaluation Results.

we conducted a paired-samples t-test on the aggregated per-record metrics across experimental conditions. This analysis compares Copilot-assisted documentation against the baseline condition while controlling for case content via a counterbalanced assignment.

As summarized in the bottom rows of Table 2, ClinDoc Copilot yields statistically significant improvements across all comparative metrics. Beyond statistical significance, we further report Cohen’s *d* to quantify the magnitude of these effects.

Efficiency gains are particularly pronounced. Documentation time exhibits a very large effect size, and keystroke count shows an even stronger effect. These results indicate not merely incremental speedups, but a fundamental shift in documentation workflow, where repetitive manual writing is substantially reduced under Copilot assistance.

In addition to efficiency, documentation quality also improves consistently. Terminology normalization demonstrates a large effect, suggesting that Copilot effectively supports standardized clinical language use. Reasoning traceability similarly shows a substantial effect, indicating improved alignment between documented diagnoses or orders and the supporting clinical evidence recorded earlier in the note.

Meanwhile, the ghost text edit rate in the Copi-

lot condition remained low across groups, indicating that most tab-accepted completions aligned well with physician intent and required minimal post-editing. Importantly, hallucination-related issues are rare under the Copilot condition. This is because all system assistance required explicit physician action and was grounded in the recorded encounter context.

Taken together, these results form a complete evidentiary chain combining descriptive statistics, effect sizes, and significance testing. They demonstrate that ClinDoc Copilot delivers robust efficiency gains while simultaneously enhancing documentation quality.

#### 4.4 User Study

After completing the documentation tasks, all participants who used ClinDoc Copilot completed a post-study questionnaire and a brief semi-structured interview. The questionnaire assessed participants’ subjective perceptions along four dimensions: (1) *efficiency improvement* achieved through text auto-completion and AI dialogue (EFF), (2) *documentation quality and standardization*, particularly terminology normalization (QUAL), (3) *trustworthiness of system assistance* grounded in the encounter context (TRUST), and (4) *willingness to adopt the system* in real outpatient practice (ADOPT). All items were rated on a five-point Likert scale with 1 = Strongly Disagree, 5 = Strongly Agree. The full questionnaire is provided in Appendix C.

Table 3 summarizes the questionnaire results, reported as mean Likert scores across participants. Overall, participants gave consistently positive ratings for efficiency improvement as well as documentation quality and standardization, suggesting that ClinDoc Copilot is perceived as practically useful in outpatient documentation. Trust-related ratings were comparatively lower but remained in the positive range, indicating a generally cautious yet favorable attitude toward system assistance during clinical use.

Item	Score (1–5 ↑)
EFF	4.50 (0.50)
QUAL	4.58 (0.51)
TRUST	4.00 (0.59)
ADOPT	4.13 (0.54)

Table 3: User study questionnaire results. Scores denote mean Likert ratings across participants (1 = Strongly Disagree, 5 = Strongly Agree).

## 5 Discussion

Our results suggest that ClinDoc Copilot provides practical value by aligning AI assistance with the legal and workflow constraints of outpatient clinical documentation. Rather than maximizing automation, the system focuses on supporting physicians while preserving clear responsibility boundaries.

ClinDoc Copilot improves efficiency primarily by reducing repetitive manual writing, rather than by replacing clinical reasoning. At the same time, inline cues help surface terminology and traceability issues that are often overlooked under time pressure. Together, these mechanisms enable faster documentation without sacrificing audit readiness.

More broadly, our findings have implications for the deployment of AI-assisted documentation tools beyond the specific system studied here. As healthcare institutions increasingly explore AI integration, there is a risk that efficiency metrics alone may drive adoption decisions. Our results caution against equating automation with usefulness, particularly in environments where legal accountability and audit readiness are non-negotiable. Interactive, physician-led systems like ClinDoc Copilot may offer a more sustainable path for deployment, especially in early stages of adoption. Over time, this alignment may prove more important than raw generative capability in determining long-term impact.

## 6 Conclusion

We presented ClinDoc Copilot, a physician-centered, human-in-the-loop assistant for outpatient clinical documentation in Chinese clinical settings. By combining conversation-grounded ghost text completion, inline compliance cues, and an on-demand AI dialogue sidebar, ClinDoc Copilot improves documentation efficiency while supporting terminology standardization and evidence-grounded traceability.

Through a controlled evaluation with licensed physicians and clinical trainees, we showed that ClinDoc Copilot reduces documentation effort, improves documentation quality, and is perceived by clinicians as practical and trustworthy. Our results highlight the importance of interactive, physician-led AI assistance for clinical documentation, particularly in settings where legal accountability and audit requirements are central. We believe ClinDoc Copilot represents a promising direction for deploying AI systems that meaningfully support, rather than replace, clinicians in real-world practice.

## 564 **Limitations**

565 This study has several limitations. First, our evaluation  
566 was conducted in simulated outpatient scenarios using trained actors rather than in live clinical  
567 environments. While this design allowed for controlled comparison across conditions, real-world  
568 deployment may introduce additional variability, such as institution-specific workflows.

569 Second, the participant pool, while balanced across roles and clinical systems, was limited in  
570 size. Larger-scale studies across more specialties and hospitals are needed to assess generalizability.  
571

572 Third, ghost text relies on audio recording and transcription of outpatient encounters. Although  
573 recordings were collected with informed consent in this study, privacy considerations and institutional  
574 policies may affect adoption in some settings.  
575

## 581 **Ethical Considerations**

582 **Human Subjects and Informed Consent.** All  
583 evaluations in this study were conducted in simulated outpatient scenarios using trained actors  
584 rather than real patients. Licensed physicians and clinical trainees participated voluntarily and provided  
585 informed consent prior to the study. Audio recordings were collected solely for the purpose  
586 of supporting documentation assistance within the simulated tasks and did not involve real patient care  
587 or identifiable patient data.  
588

592 **Patient Safety and Clinical Responsibility.** ClinDoc Copilot is not designed to diagnose, recommend  
593 treatments, or make clinical decisions. All generated suggestions are non-binding and require  
594 explicit physician action to be incorporated into the medical record. Physicians remain the sole  
595 authors and legal owners of all documentation. This design minimizes the risk of unverified or hallucinated  
596 content being silently introduced into clinical records and preserves clear responsibility boundaries  
597 consistent with regulated outpatient practice.  
598

603 **Data Privacy and Confidentiality.** The study used de-identified clinical case scripts for simulation.  
604 No real patient records were accessed or stored. Audio recordings and written records were  
605 handled in accordance with institutional data protection practices and were used exclusively for  
606 research evaluation. The system design assumes deployment within secure hospital IT environments,  
607 where access control, logging, and data governance policies are enforced by the institution. In this  
608  
609  
610  
611  
612

study, all data, code, and models we use and we 613  
release are under the CC-BY 4.0 license or Apache 614  
license 2.0. We have verified that their usage complies 615  
with the original license agreements and access 616  
conditions. Furthermore, participants involved 617  
in the human evaluation phase were explicitly informed 618  
that their evaluation results would be used solely 619  
for academic purposes. A total of 46 annotators 620  
from one country participated in this study, and 621  
each annotator was equipped with at least basic 622  
medical knowledge. Our study complies with 623  
ethical standards and regulations. 624

**Risk of Misuse and Deployment Considerations.** 625  
While AI-assisted documentation tools may improve 626  
efficiency, improper deployment could risk over-reliance 627  
or inappropriate automation. ClinDoc Copilot 628  
mitigates these risks by requiring explicit user 629  
invocation for all assistance and by avoiding end-to-end 630  
autonomous note generation. We emphasize that the 631  
system is intended as a supportive writing aid rather 632  
than a substitute for clinical judgment, and should be 633  
deployed only in settings that maintain physician 634  
oversight and regulatory compliance. 635  
636

## 637 **References**

- 638 Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam 639  
Fourney, Besmira Nushi, Penny Collisson, Jina Suh, 640  
Shamsi Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, 641  
Ruth Kikin-Gil, and Eric Horvitz. 2019. [Guidelines for human-ai interaction](#). In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, page 1–13, New York, NY, USA. Association for Computing Machinery. 642  
643  
644  
645  
646
- 647 Keyu An, Qian Chen, Chong Deng, Zhihao Du, Changfeng 648  
Gao, Zhifu Gao, Yue Gu, Ting He, Hangrui Hu, Kai Hu, 649  
Shengpeng Ji, Yabin Li, Zerui Li, Heng Lu, Haoneng Luo, 650  
Xiang Lv, Bin Ma, Ziyang Ma, Chongjia Ni, and 14 others. 2024. [Funaudiollm: Voice understanding and generation foundation models for natural interaction between humans and llms](#). *Preprint*, arXiv:2407.04051. 651  
652  
653  
654
- 655 Elham Asgari, Nina Montaña Brown, Magda Dubois, Saleh 656  
Khalil, Jasmine Balloch, Joshua Au Yeung, and Dominic Pimenta. 2025. [A framework to assess clinical safety and hallucination rates of llms for medical text summarisation](#). *npj Digital Medicine*, 8(1):274. 657  
658  
659
- 660 Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, 661  
Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel 662  
Weld. 2021. [Does the whole exceed its parts? the effect of ai explanations on complementary team performance](#). In *Proceedings* 663  
664

665	<i>of the 2021 CHI Conference on Human Factors in Computing Systems, CHI '21, New York, NY, USA.</i>	OpenAI. 2024. Gpt-4o system card. <a href="https://openai.com/index/gpt-4o-system-card/">https://openai.com/index/gpt-4o-system-card/</a> . Accessed: 2024.	721
666	Association for Computing Machinery.		722
667			723
668	Asma Ben Abacha, Wen-wai Yim, Griffin Adams, Neal Snider, and Meliha Yetisgen. 2023. <a href="#">Overview of the MEDIQA-chat 2023 shared tasks on the summarization &amp; generation of doctor-patient conversations</a> . In <i>Proceedings of the 5th Clinical Natural Language Processing Workshop</i> , pages 503–513, Toronto, Canada. Association for Computational Linguistics.	Ankit Pal, Logesh Kumar Umaphathi, and Malaikannan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Proceedings of the Conference on Health, Inference, and Learning</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260. PMLR.	724
669			725
670			726
671			727
672			728
673			729
674			730
675			
676	Aidan Gilson, Conrad W Safranek, Thomas Huang, Vimig Socrates, Ling Chi, Richard Andrew Taylor, David Chartash, and 1 others. 2023. How does chatgpt perform on the united states medical licensing examination (usmle)? the implications of large language models for medical education and knowledge assessment. <i>JMIR medical education</i> , 9(1):e45312.	D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, and Dan Dennison. 2015. <a href="#">Hidden technical debt in machine learning systems</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 28. Curran Associates, Inc.	731
677			732
678			733
679			734
680			735
681			736
682			737
683	John Giorgi, Augustin Toma, Ronald Xie, Sondra Chen, Kevin An, Grace Zheng, and Bo Wang. 2023. <a href="#">WangLab at MEDIQA-chat 2023: Clinical note generation from doctor-patient conversations using large language models</a> . In <i>Proceedings of the 5th Clinical Natural Language Processing Workshop</i> , pages 323–334, Toronto, Canada. Association for Computational Linguistics.	Mark Sendak, Madeleine Clare Elish, Michael Gao, Joseph Futoma, William Ratliff, Marshall Nichols, Armando Bedoya, Suresh Balu, and Cara O'Brien. 2020. "the human body is a black box": supporting clinical decision-making with deep learning. In <i>Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* '20</i> , page 99–109, New York, NY, USA. Association for Computing Machinery.	738
684			739
685			740
686			741
687			742
688			743
689			744
690			745
691	Microsoft Healthcare. 2025. <a href="#">Microsoft dragon copilot</a> . AI clinical assistant to streamline documentation and workflows.	Tait D. Shanafelt, Lotte N. Dyrbye, Christine Sinsky, Omar Hasan, Daniel Satele, Jeff Sloan, and Colin P. West. 2016. <a href="#">Relationship between clerical burden and characteristics of the electronic environment with physician burnout and professional satisfaction</a> . <i>Mayo Clinic Proceedings</i> , 91(7):836–848.	746
692			747
693			748
694	Kundan Krishna, Sopan Khosla, Jeffrey Bigham, and Zachary C. Lipton. 2021. <a href="#">Generating SOAP notes from doctor-patient conversations using modular summarization techniques</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4958–4972, Online. Association for Computational Linguistics.	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R. Pfohl, Heather Cole-Lewis, Darlene Neal, Qazi Mamunur Rashid, Mike Schaeckermann, Amy Wang, Dev Dash, Jonathan H. Chen, Nigam H. Shah, Sami Lachgar, Philip Andrew Mansfield, and 16 others. 2025. <a href="#">Toward expert-level medical question answering with large language models</a> . <i>Nature Medicine</i> , 31(3):943–950.	749
695			750
696			751
697			752
698			753
699			754
700			755
701			756
702			757
703	Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, and 1 others. 2023. Performance of chatgpt on usmle: potential for ai-assisted medical education using large language models. <i>PLoS digital health</i> , 2(2):e0000198.	SylvanL. 2025. <a href="#">Traditional chinese medicine dataset pretrain</a> . Dataset for text pretraining in Traditional Chinese Medicine.	758
704			759
705			760
706			761
707			
708			
709			
710	Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. <a href="#">Capabilities of gpt-4 on medical challenge problems</a> . <i>Preprint</i> , arXiv:2303.13375.	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.	762
711			763
712			764
713			
714	David Oniani, Xizhi Wu, Shyam Visweswaran, Sumit Kapoor, Shravan Kooragayalu, Katelyn Polanska, and Yanshan Wang. 2024. Enhancing large language models for clinical decision support by incorporating clinical practice guidelines. In <i>2024 IEEE 12th International Conference on Healthcare Informatics (ICHI)</i> , pages 694–702. IEEE.	Juraj Vladika, Annika Domres, Mai Nguyen, Rebecca Moser, Jana Nano, Felix Busch, Lisa C. Adams, Keno K. Bressemer, Denise Bernhardt, Stephanie E. Combs, Kai J. Borm, Florian Matthes, and Jan C. Peeken. 2025. <a href="#">Improving reliability and explainability of medical question answering through atomic fact checking in retrieval-augmented llms</a> . <i>Preprint</i> , arXiv:2505.24830.	765
715			766
716			767
717			768
718			769
719			
720			

778	Jenna Wiens, Suchi Saria, Mark Sendak, Marzyeh Ghassemi, Vincent X. Liu, Finale Doshi-Velez, Kenneth Jung, Katherine Heller, David Kale, Mohammed Saeed, Pilar N. Ossorio, Sonoo Thadaney-Israni, and Anna Goldenberg. 2019. <a href="#">Do no harm: a roadmap for responsible machine learning for health care</a> . <i>Nature Medicine</i> , 25(9):1337–1340.		
779			
780			
781			
782			
783			
784			
785			
786	Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng, Jimin Huang, Jeffrey Zhang, Vipina Keloth, Xinyu Zhou, Lingfei Qian, Huan He, Dennis Shung, Lucila Ohno-Machado, Yonghui Wu, Hua Xu, and Jiang Bian. 2025. <a href="#">Medical foundation large language models for comprehensive text analysis and beyond</a> . <i>npj Digital Medicine</i> , 8(1):141.		
787			
788			
789			
790			
791			
792			
793	Qianqian Xie, Zheheng Luo, Benyou Wang, and Sophia Ananiadou. 2023. <a href="#">A survey for biomedical text summarization: From pre-trained to large language models</a> . <i>Preprint</i> , arXiv:2304.08763.		
794			
795			
796			
797	zhangmingme. 2025. <a href="#">Medical record diagnosis chinese</a> . Apache License 2.0.		
798			
799			
	<b>A Human Subjects Information</b>		
	<b>A.1 User Interface for Human Subjects</b>		
800			
801	As show in Fig. 7, the interface supports structured documentation during the initial patient encounter and consists of two components. The medical record editing interface provides free-text fields for patient demographics, chief complaint, history of present illness, past medical history, physical examination, laboratory and imaging studies, diagnosis, and physician orders. The assistant interface includes a speech transcription assistant and a conversational question-answering assistant. With authorization, the speech assistant transcribes physician–patient conversations to support medical record generation. The conversational assistant enables physicians to interactively query and clarify medical record–related annotations during documentation.		
802			
803			
804			
805			
806			
807			
808			
809			
810			
811			
812			
813			
814			
815			
816			
817			
	<b>A.2 Instructions Provided to Participants</b>		
818	Our study involved two categories of participants. To approximate real-world clinical consultations and medical documentation workflows, medical students were instructed to act as simulated patients and interact with licensed physicians during system evaluation. Accordingly, we describe the instructions provided to participants separately by participant category.		
819			
820			
821			
822			
823			
824			
825			
		<b>A.2.1 Instructions for Medical Student Participants Acting as Simulated Patients</b>	
		<b>Role and Study Purpose.</b> You are invited to participate in a research study evaluating <i>ClinDoc Copilot</i> , an AI-assisted clinical documentation system. In this study, you will act as a <b>simulated patient</b> and interact with a licensed physician during a simulated outpatient consultation. The purpose of your participation is to approximate realistic doctor–patient interactions so that we can study clinical documentation workflows. You are not being evaluated on medical knowledge, acting performance, or communication skills.	
		<b>Case Materials and Interaction Guidelines.</b> You will be provided with a predefined clinical case script prior to each session. These case materials are derived from de-identified clinical scenarios and do not contain any real patient information. During the consultation, you should respond strictly according to the provided script and should not introduce additional symptoms, medical history, or test results beyond what is specified. If the physician requests laboratory tests, imaging studies, or examinations, you should provide results only if they are included in the case materials; otherwise, state that the requested information is unavailable. You should communicate naturally but avoid improvisation that deviates from the script.	
		<b>Audio Recording and Data Collection.</b> The consultation will be audio recorded to support transcription and documentation assistance. Audio recordings are collected solely for research purposes, including evaluation of documentation workflows and expert assessment of documentation quality. No video will be recorded. You will not be asked to provide any personal identifying information, including your real name, student ID, contact details, or personal medical history.	
		<b>Data Use, Privacy, and Consent.</b> All collected data will be de-identified and used exclusively for academic research and system evaluation. De-identified data may be released publicly for research purposes under a CC BY 4.0 license. Participation is entirely voluntary, and you may withdraw from the study at any time without penalty. By signing the consent form, you confirm that you understand the study purpose, consent to audio recording under the described conditions, and agree to the use of de-identified data for research and publication.	

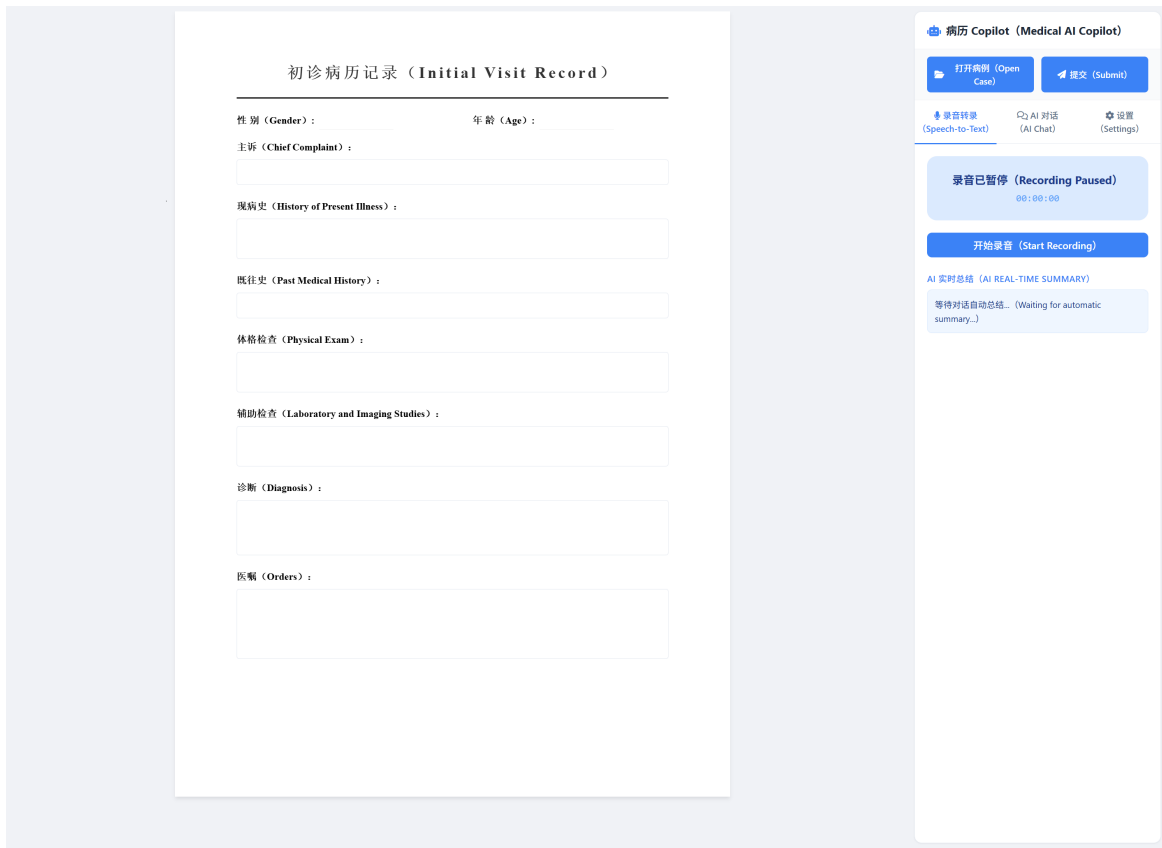


Figure 7: UI of ClinDoc Copilot

**Important Notes.** This study does not involve real patient care. No clinical decisions are made based on your responses, and the AI system does not diagnose or treat patients.

### A.2.2 Instructions for Licensed Physicians and Clinical Trainees

**Study Purpose and Your Role.** You are invited to participate in a research study evaluating *ClinDoc Copilot*, a physician-centered AI assistant designed to support outpatient clinical documentation. You will conduct simulated outpatient consultations with medical students acting as simulated patients and complete the corresponding medical records using the provided documentation interface. Some cases will be completed with AI assistance enabled, while others will be completed without assistance. The study focuses on documentation behavior rather than diagnostic accuracy or clinical performance.

**Use of ClinDoc Copilot.** When enabled, ClinDoc Copilot may provide optional ghost text suggestions, inline cues highlighting terminology or traceability issues, and an on-demand conversational assistant. The system does not diagnose,

recommend treatments, or make clinical decisions. All AI-generated suggestions are non-binding and require explicit physician action to be incorporated into the medical record. You remain the sole author and legal owner of all documentation and may ignore or modify AI assistance at any time.

**Documentation and Interaction Guidelines.** You should document each encounter as you would in routine outpatient practice, based solely on information stated during the consultation and results explicitly provided in the case materials. The AI system should not be relied upon to introduce new clinical facts or unsupported information.

**Audio Recording and Data Handling.** Doctor-patient conversations will be audio recorded to support transcription and research evaluation. No real patient data is involved, and no personal identifying information beyond role and training level will be collected. All data are de-identified and handled in accordance with institutional data protection practices.

**Voluntary Participation and Informed Consent.** Participation is voluntary, and you may withdraw from the study at any time without consequence.

924	By signing the consent form, you acknowledge that	<b>A.4 Data Use Authorization and Informed</b>	972
925	you understand the study procedures, consent to	<b>Consent Protection</b>	973
926	audio recording under the stated conditions, and	During public recruitment, it was clearly commu-	974
927	agree to the use of de-identified data for research	nicated that code collected in this study would be	975
928	and academic publication.	released under the CC BY 4.0 license. Upon en-	976
929	<b>Clarifications.</b> ClinDoc Copilot is a documenta-	rollment, all participants signed hard-copy consent	977
930	tion support tool rather than a clinical decision	forms. Through this process, participants gave in-	978
931	system. The study does not evaluate your medi-	formed consent to the collection and public release	979
932	cal competence, and the simulated setting is de-	of their data for research and academic purposes.	980
933	signed to approximate real workflows while ensur-	<b>B Rubric Definitions for Documentation</b>	981
934	ing safety and regulatory compliance.	<b>Quality Evaluation</b>	982
935	<b>A.3 Recruitment and Compensation</b>	This appendix details the rubric-based evaluation	983
936	<b>A.3.1 Compensation</b>	criteria used to assess documentation quality in	984
937	For professional doctor participants, we based	our study. All rubric scores were assigned by se-	985
938	our compensation on the average annual salary	nior physicians based on the final written medical	986
939	of physicians in the region where the study was	records, with reference to the recorded outpatient	987
940	conducted, which is approximately 490,000 RMB.	conversations. The evaluation focuses on documen-	988
941	Assuming a standard workload of 2,000 working	tation quality rather than medical decision correct-	989
942	hours per year, this corresponds to an average	ness. Two senior medical experts from medical	990
943	hourly wage of about 245 RMB. Given that partici-	institute are hired to judge all the documentation	991
944	ipation in our study required sustained concentra-	with payment of 1000 RMB per hour.	992
945	tion, we provided compensation at more than twice	<b>B.1 Term Normalization</b>	993
946	the average hourly rate, amounting to 800 RMB	<b>Core consideration.</b> This dimension evaluates	994
947	per hour.	the system's ability to assist physicians in convert-	995
948	For student participants, compensation was de-	ing colloquial, patient-facing language into stan-	996
949	termined in accordance with common stipend prac-	dardized clinical terminology. The assessment re-	997
950	tices for graduate students in the region where the	fects how effectively informal expressions in the	998
951	study was conducted. Medical student participants	encounter are transformed into appropriate writ-	999
952	received a stipend of 3,000 RMB, disbursed on a	ten medical language using the system's internal	1000
953	monthly basis. Although the experimental period	language knowledge.	1001
954	did not exceed one month, the full monthly stipend	<b>Scoring guidelines (1–5).</b>	1002
955	was provided in consideration of the workload in-	• <b>5 (Excellent):</b> All colloquial expressions are	1003
956	involved.	consistently normalized into standard medical	1004
957	Overall, participant compensation in this study	terminology. The record is written entirely in	1005
958	exceeded typical income benchmarks in the local	professional clinical language without unnec-	1006
959	context and was designed to be reasonable and com-	essary patient phrasing.	1007
960	mensurate with the effort, attention, and expertise	• <b>4 (Good):</b> Most colloquial expressions are	1008
961	required, without constituting undue inducement.	properly normalized, with only minor resid-	1009
962	<b>A.3.2 Recruitment</b>	ual informal phrasing that does not affect read-	1010
963	Recruitment was conducted primarily through open	ability or professionalism.	1011
964	calls issued by physician participants involved in	• <b>3 (Acceptable):</b> Key medical terms are cor-	1012
965	this study within their respective institutions. In	rect, but multiple patient-style expressions re-	1013
966	addition, several physician participants were affili-	main in the text and would require manual	1014
967	ated with teaching hospitals and medical schools;	revision for formal documentation.	1015
968	these individuals publicly recruited student partici-	• <b>2 (Poor):</b> Colloquial language is frequently	1016
969	pants through their associated medical schools. All	retained, or medical terminology is used in-	1017
970	recruitment processes were conducted in an open	consistently or imprecisely.	1018
971	and transparent manner.		



1111 underline correction hints, and AI explanations  
1112 generated by ClinDoc Copilot are trust-  
1113 worthy and demonstrate clinical logical ratio-  
1114 nality.”

1115 **4. Willingness for Real-World Adoption.**

1116 “Given ClinDoc Copilot’s advantages in im-  
1117 proving documentation efficiency and ensur-  
1118 ing standardization, I would be willing to use  
1119 it as a regular auxiliary tool in my daily out-  
1120 patient practice if deployed in the hospital.”

1121 **C.2 Response Scale**

1122 All questions were answered using the following  
1123 five-point Likert scale:

- 1124 • Strongly Agree
- 1125 • Agree
- 1126 • Neutral
- 1127 • Disagree
- 1128 • Strongly Disagree

1129 The questionnaire focuses on subjective user  
1130 perception rather than objective performance and  
1131 complements the quantitative efficiency and expert-  
1132 based quality evaluations reported in the main pa-  
1133 per.

1134 **D Use of AI Assistants**

1135 AI assistants were employed in a strictly supportive  
1136 capacity during the manuscript preparation process.  
1137 Their use was limited to assisting with English lan-  
1138 guage editing, including improvements to grammar,  
1139 wording, and overall clarity in selected sections of  
1140 the text.

1141 All substantive aspects of the work—including  
1142 the formulation of research questions, methodologi-  
1143 cal development, experimental design, implementa-  
1144 tion, evaluation, and interpretation of results—were  
1145 entirely carried out by the authors. AI assistants  
1146 did not contribute to analysis, result interpretation,  
1147 or any form of scientific decision-making.