

Med-CoT: Towards Proactive Clinical Agentic LLM

Gunnar Von Bergen*, Jingyi Huang*, Maria Kunigk-Bakalar, Liran Ma, Yang Zhang
Department of Computer Science and Software Engineering, Miami University
{vonbergb, huangj84, kunigkmc, mal18, zhang981}@miamioh.edu

Abstract—Current AI-based clinical support agents typically operate in a passive question-answer mode, which limits their effectiveness in real-world clinical settings. To bridge this gap, we introduce Med-CoT, a Medical Chain-of-Thought framework designed to transform LLMs into proactive clinical assistants. By integrating sentiment assessment, strategic retrieval, and multi-step reasoning, Med-CoT enables agents to dynamically clarify ambiguities and intervene in safety-critical scenarios. Experimental results on real-world medical dialogues demonstrate that Med-CoT significantly outperforms direct prompting baselines, achieving a more diverse action distribution, higher sensitivity to safety risks, and reduced repetitive behaviors, thereby aligning more closely with expert clinical decision-making.

I. INTRODUCTION

AI is rapidly reshaping the medical landscape, driving innovation in diagnosis, treatment, and patient management while revealing the limitations of current clinical systems [1]. From decision support tools to predictive analytics, AI has become an indispensable component of modern healthcare, assisting providers in interpreting medical images, optimizing hospital workflows, and personalizing care plans [2]. However, much of the current use of AI in the medical field remains focused on retrospective data analysis and reactive intervention, addressing problems only after they arise [3]. In contrast, proactive clinical support marks a significant shift toward prevention and foresight [4]. This is significant in the medical field because it allows healthcare providers to implement preventive interventions, thereby improving patient outcomes, reducing the frequency of costly emergency care, and overall enhancing the efficiency of healthcare delivery [5]. Beyond predictive modeling, proactive systems can enable augmented data preparation by harvesting information that is otherwise difficult to extract from patient-authored texts, such as symptom logs and digital health communications [6]. Through natural language processing (NLP), common descriptions of symptoms, medication responses, and lifestyle changes can be identified and standardized, creating an enriched data foundation for longitudinal analyses and contextually aware clinical decision-making [7].

Prior work in AI-driven clinical support has established a strong groundwork, ranging from machine-learning-based diagnostic aids to large-scale health monitoring systems. One of the successful proactive healthcare deployments is Greece’s national EVA system, which was used to dynamically assess health risks of patients and guide targeted testing strategies [8]. EVA demonstrated how predictive analytics could be applied

proactively to public health management rather than reacting to them. Another proactive AI system in clinical practice is for Targeted Real-Time Early Warning System (TREWS) for sepsis, in which, on prospective multicenter studies in five different hospital setups, physician utilization of TREWS was correlated with decreased in-hospital deaths and earlier antimicrobials, signifying preventive efforts via real-time AI [9]. Nevertheless, while examples like EVA and TREWS illustrate what is possible at the proaction level, many clinical AI tools remain passive in their interaction with patients and providers. Current AI conversational chatbots typically follow a question-answer pattern by passively answering a patient’s questions, rather than actively leveraging sentiment analysis, retrieval capabilities, and memory of agentic Large Language Models (LLMs) to retrieve medical records, and proactively clarify any detected confusion or address a safety risk [10].

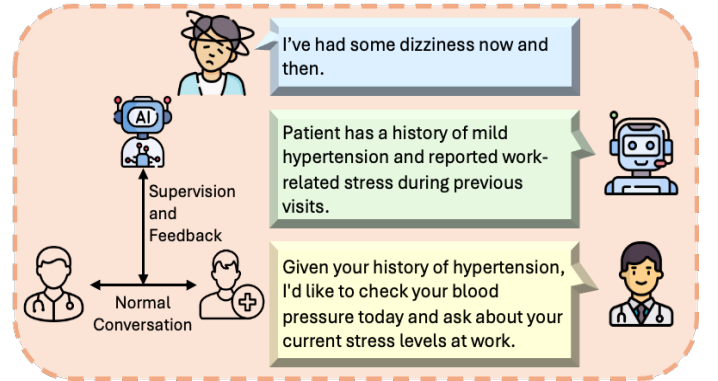


Fig. 1. Med-CoT Overview: Proactive Clinical Dialogue Assistant

To address this limitation, we introduce a proactive clinical support agent powered by a novel Medical Chain-of-Thought (Med-CoT) framework (see Figure 1). Specifically, Med-CoT constructs dynamic reasoning trajectories that integrate restricted dialogue windows with sentiment-aware assessments. This approach enables the agent to not only respond to queries but to actively anticipate patient needs, whether by retrieving missing medical context, clarifying ambiguous information, or redirecting safety-critical cases. We evaluate our proposed model on a real-world medical dialogue dataset. Experiments indicate that Med-CoT significantly shifts the interaction paradigm from passive repetition to active engagement, and successfully achieves a decision-making profile that closely mirrors the balanced and safe judgment of domain experts.

*The first two authors contributed equally to this work.

II. APPROACH

To address the challenges of safety, consistency, and clinical coherence in Large Language Model (LLM) driven doctor-patient dialogues, we introduce Med-CoT (Medical Chain-of-Thought). As depicted in Figure 2, this framework follows strict reasoning process before taking any action or generating any response. Unlike standard end-to-end conversational models that map input utterances directly to responses, Med-CoT decouples clinical assessment from dialogue generation. This ensures that every generated response is grounded in a verified clinical context and the safety assessment.

A. Assessment and SOAP Note Generation

During the inference, the model receives only a restricted dialogue window $D_i = \{u_i\}$ with the most recent exchange $u_i = (d_i, p_i)$ between the doctor d_i and the patient p_i . This one-turn dialogue window enforces the model to focus on the current conversation, which prevents it from being distracted by irrelevant history from previous turns.

At each conversational turn i , Med-CoT performs a rigorous assessment stage ($\langle assess \rangle$). This stage is designed to fully leverage the sentiment analysis ability of the LLM and mimic the cognitive process of a clinician paying attention to the patient’s mood while formulating a plan. This assessment consists of two processes, a binary sentiment state classification (M_i) and the SOAP note generation ($SOAP_i$):

$$\langle M_i, SOAP_i \rangle = assess(D_i). \quad (1)$$

The model first scans the current dialogue window to detect four specific and clinically relevant signals. These are binary flags (True or False) that determine the immediate necessity of intervention:

- **Confusion:** Indicates whether the patient expresses explicit uncertainty, such as asking “why” or “what” questions, or presenting conflicting information. This flag is to differentiate between simple questions and actual comprehension gaps that require extra support.
- **Stress/Anxiety:** Identifies emotional expressions of anxiety, fear, frustration, or urgency. This is vital for assessing the tone of the response, which ensures that clinical advice is appropriately delivered.
- **Safety Risk:** A critical safety layer that flags self-harm ideation, potential medical emergencies, such as chest pain, or dangerous medication use, such as overdose. The activation of this flag overrides most other logic to prioritize immediate redirection.
- **Missing Information:** Set to “True” only if neither the doctor nor the agent can safely or meaningfully respond without specific clinical details, such as dosage, duration, or severity. This prevents the model from hallucination based on incomplete information.

After the binary assessment, the model generates a short note to explain why the specific flags are set. Next, the model generates a structured SOAP (Subjective, Objective,

Assessment, Plan) note to briefly summarize the context within the dialogue window [11]:

- **Subjective (S):** Integrates all information provided from patient-reported symptoms, concerns and history presented in the dialogues.
- **Objective (O):** Extracts verifiable facts, especially vital signs or information retrieved from records.
- **Assessment (A):** Formulates a working hypothesis or summary of the current clinical status from S and O.
- **Plan (P):** Specifies the intended clinical direction (e.g., “Ask about fever” or “Reassure about side effects”, etc.) before the final response is generated.

This intermediate diagnostic SOAP layer forces the model to ground its understanding in objective evidence. By explicitly making the assessment and plan before the final response, the model reduces the inconsistency and ensures that the final output aligns with standard medical reasoning.

B. Proactive Action Design

To effectively implement proactivity in clinical dialogue, we design a structured series of actions to guide the behavior of the LLM. This design serves two essential purposes: i) It motivates the agent to take initiative in patient care by clarifying ambiguities or correcting misconceptions, while simultaneously constraining the output to prevent unstandardized, random actions. ii) It avoids blending conflicting objectives, such as attempting to ask additional questions or deliver complex supportive messages while retrieving missing clinical data in one single turn. At each turn i , the model is required to first perform a required evaluation step and then select exactly one intervention action from the strictly defined action set $a_i \in \mathcal{A}$ based on the given sentiment classification and SOAP note, as specified in Table I. Formally, the decision-making process is modeled as:

$$\langle a_i, R_i \rangle = \arg \max P(a, R \mid M_i, SOAP_i), \quad a_i \in \mathcal{A}. \quad (2)$$

M_i represents the binary assessment flags decided in the first stage, and $SOAP_i$ denotes the generated clinical note. The model maximizes the probability of the action-response pair $\langle a_i, R_i \rangle$ conditioned on these structured insights. The final output R_i corresponds to the natural language utterance delivered to the patient.

C. Cognitive-aware Decision Policy

The transition from assessment to action decision is controlled by a deterministic and hierarchical decision policy rather than only depending on probabilistic sampling. This policy prioritizes patient safety above the conversational flow, which ensures the agent’s behavior is predictable and clinically defensible. It encourages the model to take the action `<retrieve>` at the beginning of each dialogue so that the final responses of the entire session are grounded in the patient’s historical record. Similarly, an `<summarize>` action is triggered by the last turn of the dialogue if without safety risk flagged to consolidate the key findings. During the prior assessment stage, the model evaluates immediate

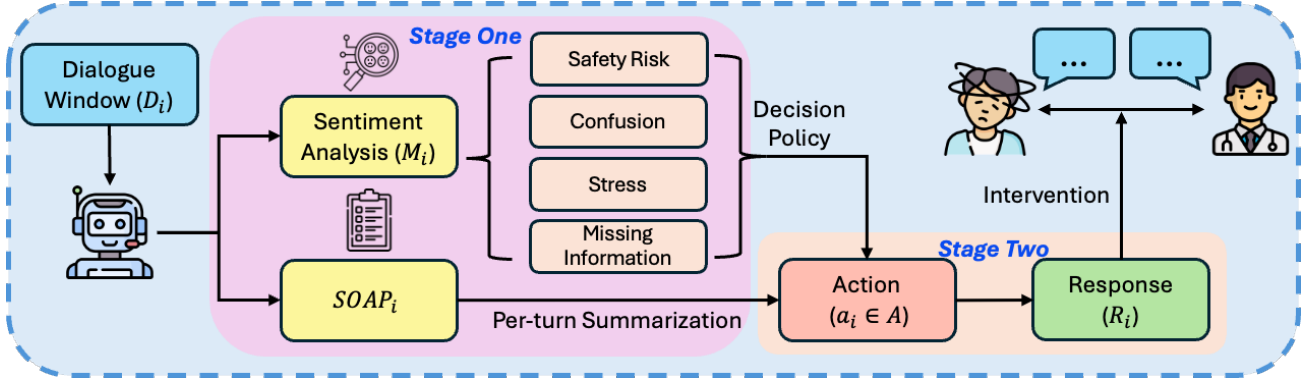


Fig. 2. Med-CoT Two-Stage Pipeline with Sentiment Assessment Integration and Cognitive-Aware Action Design

TABLE I
DEFINED ACTIONS FOR EACH DIALOGUE TURN

Action	Description
<assess>	Required every turn. Detects confusion, stress, safety risk, and missing information.
<retrieve>	References prior medical records to support the current assessment.
<elaborate>	Provides supportive information to resolve confusion and stress.
<clarify>	Asks a concise question to fill a knowledge gap.
<redirect>	Corrects misconceptions or addresses potential safety risks.
<summarize>	Generates a brief summary at the end of the conversation.
<defer>	Refrains from intervening when no action is required.

safety. If a safety risk is detected, such as signs of self-harm or an emerging medical emergency, the model bypasses all other logic except the <retrieve> and immediately selects <redirect>. This “safety gate” ensures that the background information is fully retrieved to support the security alert.

If no safety risks are identified, the policy changes to support the adequacy of the available clinical information and the stress level of the patient. When the patient’s confusion or stress is detected i) with critical missing information, it motivates the model to take the action <clarify> and ask a question to fill the information gap, which prevents offering inappropriate or general advice when lacking the necessary clinical context; ii) without missing information flag, the model takes action to <elaborate> and provide supportive information to resolve the patient’s confusion or release the stress. Finally, if none of these specific criteria are met, indicating that the dialogue is safe and proceeding smoothly, the action <defer> is encouraged to maintain a natural conversation flow between the patient and the doctor.

In many AI-augmented clinical support systems, the models tend to over-intervene, including constantly offering advice even when the doctor-patient dialogue is flowing naturally [12]. Some researchers argue that excessive interventions from the third-party could be distracting and harmful to the natural conversation flow. Therefore, AI interventions should be context-aware and minimally intrusive to maintain the user’s

optimal cognitive state [13]. Our approach follows this theory of “Cognitive Flow”, and prefers the <defer> action when no confusion, stress, or safety risk detected, thereby preserving the cognitive resources of the human participants.

D. Conversation State Management

The long context management in extended medical dialogues is a significant challenge for LLMs, primarily due to token limitations and the degradation of reasoning over long sequences [14], [15]. To handle this, Med-CoT adopts a stateful approach that maintain local accuracy with global coherence. To narrow the focus during assessment stage, the model works within a constricted dialogue window containing only the most recent doctor-patient exchange. It prevents the model from being distracted by historical context when assessing, ensuring that its immediate reaction is grounded in the current state. However, to maintain continuity, the system implicitly tracks the global conversation state via a unique session identifier, allowing access to the historical turns without overloading the active context window. With this conversation state management, for example, when the model selects <retrieve> at turn i , the system proactively injects the retrieved patient’s background information into the input schema for the subsequent turn $i + 1$. Med-CoT optimizes token usage and aligns its information retrieval process with the realistic workflow of a clinician.

Med-CoT Output Example

Dialogue Window:

- Patient: I do it at home. Mainly two to three times a day. I increase or decrease the amount of insulin that I take based on my numbers.
- Doctor: Don’t you consult anyone before doing that?

Sentiment Assessment:

- Confusion: False
- Stress: False
- Safety Risk: True
- Missing Information: True
- Notes: Patient self-reports adjusting insulin at home; clinician asks if they consult anyone—unsupervised insulin changes pose a medication safety risk.

SOAP:

- S: Patient self-monitors glucose 2–3 times daily and adjusts insulin based on readings.
- O: No new objective values in this turn; prior records show diabetes with HbA1c 6.8
- A: Potentially unsafe unsupervised insulin titration with risk of hypoglycemia or hyperglycemia; need clear dosing plan.
- P: Advise immediate safety measures, stop unsupervised dose changes, obtain specifics of how doses are chosen, and arrange diabetes-team review of insulin regimen.

Action: redirect

Response: Adjusting insulin without consulting your diabetes provider can be risky—please stop unsupervised dose changes and we need you to tell us exactly how you decide doses so we can make a safe plan; if you have signs of low blood sugar, treat with quick-acting carbs and seek urgent help if severe.

III. EVALUATION

We perform a comprehensive evaluation to compare the proactive decision-making capabilities of Med-CoT against a direct prompting baseline using GPT-5.

A. Dataset Preparation

Our experiments are conducted using the MTS-Dialogue dataset [16], which contains real-world multi-turn clinical dialogues between doctors and patients. We split each multi-turn dialogue into one-turn question-answer pairs and label them with a conversation ID, turn ID, and a last-turn indicator to construct the desired input dialogue window. Additionally, because the MTS-Dialogue dataset does not provide supportive materials for retrieval, such as patient profiles or medical records, we use GPT-4o to synthetically generate background information that is consistent with each dialogue but not explicitly stated in it. Due to time limitations, we use only the first 500 dialogues in the MTS-Dialogue dataset (see Table II), after filtering out dialogues with fewer than two turns to ensure that our Med-CoT design is fully tested.

TABLE II
DATASET STATISTICS

Statistic	# Conversation	# Turn	Avg. len(Dialogue Window)	Avg. len(Patient Profile)
MTS-Dialogue	500	4764	20.55	111.25

B. Action Distribution

From Figure 3, we observe that with Med-CoT, the model takes more balanced and diverse actions, effectively utilizing the full spectrum of available interventions: “clarify” (9.4%), “elaborate” (4.3%), “redirect” (3.0%), “retrieve” (10.5%), “summarize” (10.5%). This diversity indicates that the model is dynamically adapting strategy based on the specific nuances of each turn rather than relying on a default and general behavior. On the contrary, the baseline GPT-5 model shows a monolithic behavior pattern skewing towards the “clarify”, which accounts for about 86.3% of all turns. This suggests that without Med-CoT reasoning framework, the model defaults to asking follow-up questions as a safe but redundant strategy. Furthermore, it is notable that Med-CoT enables a strategic

choice to “defer” (62.3%) when the doctor-patient interaction is proceeding smoothly. This mechanism is essential to prevent the cognitive overload of human participants by avoiding unnecessary interruptions, thus ensuring the proactive contributions are both timely and valuable. In sharp contrast, the baseline’s inability to defer (0.0%) results in consistent interruptions that disrupts the natural flow of clinical dialogue.

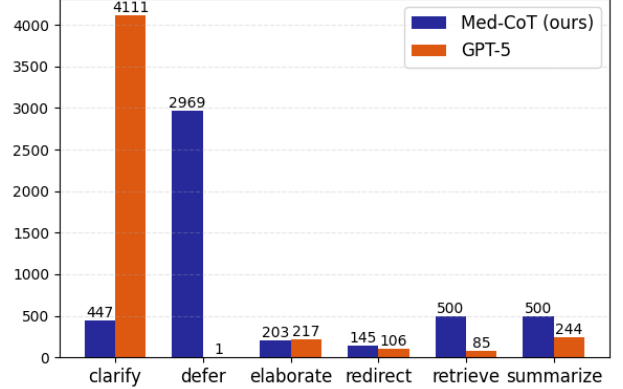


Fig. 3. Action Distribution

To evaluate the models’ responsiveness to critical situations, we invited two human experts to manually annotate the stress and safety risks based on the dialogues. For conversations flagged with a verified safety risk, we track whether the models trigger at least one “redirect” action to connect the patient with immediate help. As presented in Table III, Med-CoT achieves a significantly higher redirect accuracy compared to the baseline. While the baseline often continues to ask clarifying questions even in the face of clear safety risk signals, Med-CoT correctly identifies the urgency and executes a “redirect” action, demonstrating its superior capability in emergent clinical situations.

TABLE III
REDIRECT ACCURACY

Model	Med-CoT (ours)	GPT-5
Accuracy	68.8	18.8

C. Sentiment Detection Analysis

We analyze the distribution of sentiment assessment flags generated by Med-CoT, as shown in Table IV. It presents the “missing_info” flag is significantly more prevalent than other categories. This underscores the critical need for an AI-augmented proactive system to fill the substantial knowledge gaps between patients and doctors. Especially when the patient is unaware of the specific information they lack, a system capable of identifying these gaps can provide necessary elaboration. In addition, we utilized our manually labeled ground truth to rigorously evaluate the performance of the Med-CoT assessment module. Our model demonstrates strong performance in identifying critical patient conditions, achieving 71.4% accuracy in detecting safety risks and 88.4% accuracy in detecting patient stress. This high alignment between the model’s internal assessments and human labels validates the robustness of our sentiment analysis framework.

TABLE IV
SENTIMENT FLAG DISTRIBUTION

Sentiment Flag	Confusion	Stress	Safety Risk	Missing Information
%	5.9	13.4	3.6	45.5

D. Repetition Analysis

Additionally, we evaluate the repetitiveness of the models' behavior to further assess the quality of the generated interactions. A common issue in dialogue systems is the tendency to get stuck in loops, repeating the same actions or responses [17]. The results, presented in Table V, indicate that the GPT-5 baseline repeats non-defer actions in nearly every conversation (99.2%) with similar clarifying questions, while Med-CoT exhibits a remarkably low repetition rate, with only 28.4% of conversations containing repeated non-defer actions. Even when Med-CoT does intervene multiple times, the low cosine similarity of its final responses, 0.128 vs. 0.206 for the baseline, demonstrates that each intervention is distinct and contextually tailored. Therefore, Med-CoT is not merely reacting to keywords, but is maintaining a coherent and evolving understanding of the dialogue state.

TABLE V
ACTION AND RESPONSE DIVERSITY ANALYSIS

Statistic	Med-CoT	GPT-5
% Conversations with Repeated Non-defer Actions	28.4	99.2
% Repeated Non-defer Actions	25.7	81.7
Avg. Non-defer Final Responses Similarity	0.128	0.206

IV. CONCLUSION

In this work, we present Med-CoT, a novel framework that empowers clinical support agents to transition from passively responding to proactively intervening in healthcare delivery. By integrating sentiment assessment with structured reasoning trajectories, our approach enables Large Language Models (LLMs) to dynamically take actions to clarify, elaborate, redirect, retrieve, and summarize based on the evolving context of the dialogue. Our extensive evaluation on the MTS-Dialogue dataset demonstrates that Med-CoT significantly outperforms standard direct prompting methods, achieving a more balanced action distribution, superior safety risk detection, and better protection against repetitive behaviors. These findings underscore the potential of chain-of-thought reasoning to enhance the safety and utility of AI systems in potential or emerging clinical emergencies, paving the way for more autonomous and context-aware medical assistants.

ACKNOWLEDGMENTS

This research is supported by the National Science Foundation under Grant No. 2449011. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either

expressed or implied, of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

REFERENCES

- [1] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, C. Cui, G. Corrado, S. Thrun, and J. Dean, "A guide to deep learning in healthcare," *Nature medicine*, vol. 25, no. 1, pp. 24–29, 2019.
- [2] E. J. Topol, "High-performance medicine: the convergence of human and artificial intelligence," *Nature medicine*, vol. 25, no. 1, pp. 44–56, 2019.
- [3] C. J. Kelly, A. Karthikesalingam, M. Suleyman, G. Corrado, and D. King, "Key challenges for delivering clinical impact with artificial intelligence," *BMC medicine*, vol. 17, no. 1, p. 195, 2019.
- [4] A. C. Raldow, N. Raja, C. W. Villaflores, S. A. Skootsky, E. A. Jauregui, H. L. Rosenstein, S. D. Meshkat, S. S. Vangala, and C. A. Sarkisian, "Proactive care management of ai-identified at-risk patients decreases preventable admissions," *The American journal of managed care*, vol. 30, no. 11, p. 548, 2024.
- [5] D. W. Bates, D. Levine, A. Syrowatka, M. Kuznetsova, K. J. T. Craig, A. Rui, G. P. Jackson, and K. Rhee, "The potential of artificial intelligence to improve patient safety: a scoping review," *NPJ digital medicine*, vol. 4, no. 1, p. 54, 2021.
- [6] C. Dreisbach, T. A. Koleck, P. E. Bourne, and S. Bakken, "A systematic review of natural language processing and text mining of symptoms from electronic patient-authored text data," *International journal of medical informatics*, vol. 125, pp. 37–46, 2019.
- [7] T. A. Koleck, C. Dreisbach, P. E. Bourne, and S. Bakken, "Natural language processing of symptoms documented in free-text narratives of electronic health records: a systematic review," *Journal of the American Medical Association*, vol. 26, no. 4, pp. 364–379, 2019.
- [8] H. Bastani, K. Drakopoulos, V. Gupta, I. Vlachogiannis, C. Hadjichristodoulou, P. Lagiou, G. Magiorkinis, D. Paraskevis, and S. Tsiodras, "Efficient and targeted covid-19 border testing via reinforcement learning," *Nature*, vol. 599, no. 7883, pp. 108–113, 2021.
- [9] R. Adams, K. E. Henry, A. Sridharan, H. Soleimani, A. Zhan, N. Rawat, L. Johnson, D. N. Hager, S. E. Cosgrove, A. Markowski, *et al.*, "Prospective, multi-site study of patient outcomes after implementation of the trews machine learning-based early warning system for sepsis," *Nature medicine*, vol. 28, no. 7, pp. 1455–1460, 2022.
- [10] M. Laymouna, Y. Ma, D. Lessard, T. Schuster, K. Engler, and B. Lebouché, "Roles, users, benefits, and limitations of chatbots in health care: rapid review," *Journal of medical Internet research*, vol. 26, p. e56930, 2024.
- [11] V. Podder, V. Lew, and S. Ghassemzadeh, *SOAP Notes*. Treasure Island (FL): StatPearls Publishing, 2024.
- [12] Y. Kim, Z. Hu, H. Jeong, E. W. Park, S. S. Li, C. Park, S. Xiong, M. Lu, H. Lee, X. Liu, *et al.*, "Behaviorsft: Behavioral token conditioning for health agents across the proactivity spectrum," in *Findings of the Association for Computational Linguistics: EMNLP 2025*, pp. 9789–9817, 2025.
- [13] D. Dissanayake and S. Nanayakkara, "Navigating the state of cognitive flow: Context-aware ai interventions for effective reasoning support," *arXiv preprint arXiv:2504.16021*, 2025.
- [14] W. Xiong, J. Liu, I. Molybog, H. Zhang, P. Bhargava, R. Hou, L. Martin, R. Rungta, K. A. Sankararaman, B. Oguz, *et al.*, "Effective long-context scaling of foundation models," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 4643–4663, 2024.
- [15] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang, "Lost in the middle: How language models use long contexts," *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 157–173, 2024.
- [16] A. B. Abacha, W.-w. Yim, Y. Fan, and T. Lin, "An empirical study of clinical note generation from doctor-patient encounters," in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 2291–2302, 2023.
- [17] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi, "The curious case of neural text degeneration," *arXiv preprint arXiv:1904.09751*, 2019.