# Baichuan-M2: Scaling Medical Capability with Large Verifier System

**Baichuan-M2 Team**

## Abstract

As large language models (LLMs) advance in conversational and reasoning capabilities, their practical application in healthcare has become a critical research focus. However, there is a notable gap between the performance of medical LLMs on static benchmarks such as USMLE and their utility in real-world clinical decision-making. This discrepancy arises because traditional exams fail to capture the dynamic, interactive nature of medical consultations. To address this challenge, we introduce a novel dynamic verification framework that moves beyond static answer verifier, establishing a large-scale, high-fidelity interactive reinforcement learning system. Our framework comprises two key components: a *Patient Simulator* that creates realistic clinical environments using de-identified medical records, and a *Clinical Rubrics Generator* that dynamically produces multi-dimensional evaluation metrics. Building on this foundation, we develop Baichuan-M2, a 32B-parameter medical augmented reasoning model trained through a multi-stage reinforcement learning strategy with an improved Group Relative Policy Optimization (GRPO) algorithm. Evaluated on HealthBench, Baichuan-M2 outperforms all other open-source models and most advanced closed-source counterparts, achieving a score above 32 on the challenging HealthBench Hard benchmark—previously exceeded only by GPT-5. Our work demonstrates that robust dynamic verifier system is essential for aligning LLM capabilities with practical clinical applications, establishing a new Pareto front in the performance-parameter trade-off for medical AI deployment.

## 1 Introduction

As the conversational and reasoning capabilities of large language models (LLMs) continue to advance, there is increasing interest in their practical application in specific domains. The healthcare sector, in particular, has become a key area of research, attracting significant investment from both global tech giants and innovative startups [1–3]. Among the various approaches aimed at enhancing the capabilities of LLMs in healthcare, reinforcement learning with verifiable rewards (RLVR) has garnered considerable attention [4, 5]. This technique has already demonstrated impressive results in areas such as mathematics [6–8], code [9], agents [10, 11], and multimodality [12, 13]. These achievements highlight its potential to significantly enhance model reasoning, making its application in healthcare a highly promising research direction.

The core of RLVR lies in the development of a robust evaluation system. Its success in fields like mathematics and coding can be attributed to the availability of precise and reliable evaluation metrics. However, when it comes to assessing LLMs in the medical domain, a significant gap exists between current evaluation methods and real-world applications. Models that perform well on medical professional exams, such as the USMLE [14], often underperform in practical clinical decision-making. This discrepancy arises because traditional static benchmarks fail to capture the dynamic and complex nature of clinical practice. Real-world medical consultations frequently involve incomplete information, multiple rounds of diagnostic exploration, and nuanced communication skills, all of which are not adequately measured by conventional exams.

To address these challenges, we transitioned our focus from static answer verifiers to the development of a large-scale, high-fidelity interactive reinforcement learning verifier system. This system transcends conventional answer verifier by simulating real-world clinical scenarios, allowing the model to learn and adapt through simulated "practice" in a virtual clinical environment. Building on this foundation, we introduce Baichuan-M2, a medical augmented reasoning model that marks a significant advancement in open-source medical artificial intelligence.

Specifically, our verifier system comprises two key components. The first is a Patient Simulator, which integrates desensitized medical records and doctor-patient conversation records to effectively simulate patients with diverse social backgrounds and personality traits. This provides a highly realistic interactive environment. The second component is a Clinical Rubrics Generator, which can emulate the clinical reasoning of experienced doctors. It dynamically generates quantifiable evaluation rubrics on a large scale, based on multiple dimensions such as diagnostic accuracy, consultation logic, treatment plan rationality, communication empathy, and medical ethics.

Our training process includes mid-training for medical domain adaptation, supervised fine-tuning (SFT) with rejection sampling, and multi-stage reinforcement learning (RL) using an improved Group Relative Policy Optimization (GRPO) [8] algorithm. Specifically, we employ a multi-stage reinforcement learning strategy to decompose complex reinforcement learning tasks into a controllable hierarchical structure. This approach enhances various capabilities, including medical knowledge, reasoning, and patient interaction, while maintaining the general capabilities of the model.

We evaluate our model on the challenging HealthBench dataset [15], developed by OpenAI. Despite its relatively small number of parameters (only 32B), Baichuan-M2 outperformed all other open-source models, including gpt-oss-120B, and most advanced closed-source counterparts on HealthBench. It particularly excelled on the HealthBench Hard test, achieving a score exceeding 32, a performance level previously reached by only one other model globally, GPT-5. These experimental results underscore the critical role of a robust validation system in integrating model capabilities with practical applications.

In summary, our contributions can be highlighted as follows:

- A dynamic verifier system tailored for clinical scenarios, which addresses the limitations of previous verification methods based on static data. This method employs a patient simulator to create a high-fidelity decision-making environment and uses clinical rubrics to generate quantitative evaluation metrics in real time, thereby enhancing the reliability of the verification process.

- An advanced training method that successfully implements a multi-stage reinforcement learning strategy in a dynamic interactive environment, featuring targeted improvements to the GRPO algorithm. This enhancement enables the model to move beyond static knowledge memorization and deeply align with the advanced clinical reasoning capabilities of medical experts.

- An advanced open-source model, Baichuan-M2, which achieves top-tier performance at a remarkably lower deployment cost, setting a new Pareto front in the performance-parameter trade-off. This efficiency makes the deployment of advanced medical AI more feasible in resource-constrained healthcare settings.

## 2   Verifier System

In recent years, RLVR has achieved remarkable success in complex reasoning domains such as mathematics, coding, and agentic systems. Constructing more verifiable complex problems and environments has become a core driver for continuous breakthroughs in model capabilities. However, when applying this paradigm to the medical field, we discovered significant limitations: static answer verifier built on traditional medical question banks fails to capture the dynamic complexity of real-world diagnostic processes, often leading to limited generalization and suboptimal performance in practical applications. Real clinical practice is a partial observable, multi-turn decision-making process that relies heavily on a physician's dynamic judgment, entailing the integration of clinical experience, communication skills, and ethical considerations.

To address this challenge, in the development of Baichuan-M2, we shifted our focus from building static answer verifiers to creating a large-scale, high-fidelity dynamic interactive reinforcement learning environment. This environment aims to construct a "virtual clinical world" where models can "train and grow". The system primarily consists of two key modules: a "patient simulator" and
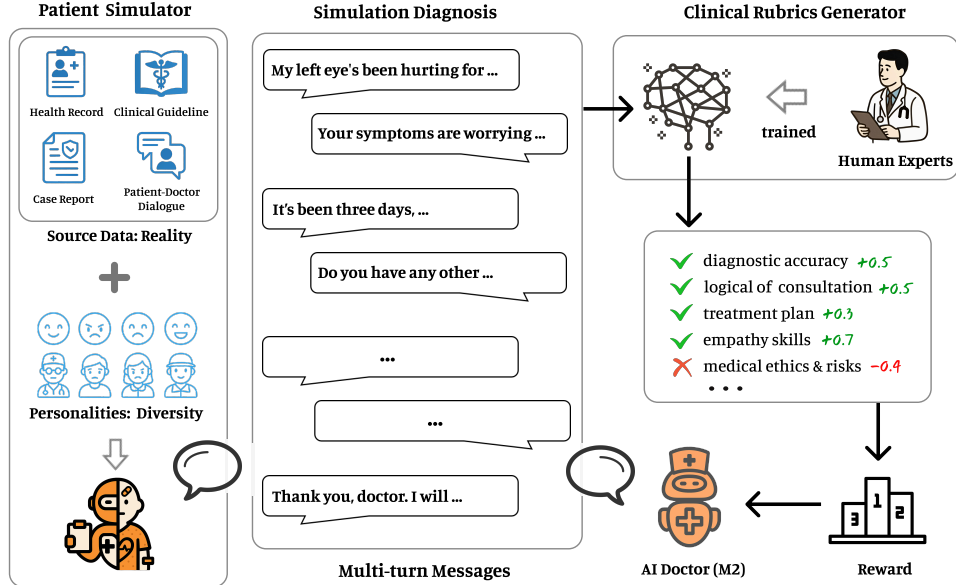
Figure 1: Verifier System Framework

a "clinical rubrics generator". The patient simulator elevates the training environment beyond rigid single-turn QA, generating realistic, stochastic, continuous interaction scenarios. The clinical rubrics generator dynamically produces verification rules for answers, enabling continuous and dynamic quantitative assessment of a model's comprehensive performance across multi-turn interactions as shown as Figure 1.

Through this closed-loop system, we successfully implemented large-scale end-to-end reinforcement learning. The model continuously interacts with "virtual patients", iteratively optimizing its diagnostic strategies based on dense feedback from "expert-level evaluations". Ultimately, the model's capabilities move beyond recall of static knowledge, achieving deep alignment with the clinical thinking and practical skills of senior physicians.

## 2.1 Patient Simulator

Patient simulators play a critical role in the training and evaluation of AI physicians [16, 17]. These simulators offer a dynamic testing environment that can effectively address the limitations of traditional static testing methods, which often fail to adequately assess the dynamic diagnostic capabilities of LLMs. However, widely used simulators in prior works [18, 19] fall short in comprehensively modeling patients' psychological states, social backgrounds, and dynamic interactions. This deficiency reduces these simulators to static databases, thereby limiting their ability to replicate the complexity of real-world clinical encounters. Such encounters often involve information withholding, emotional expressions, and culturally-mediated communication barriers, all of which are crucial for the adaptability of AI physicians in practical settings.

The core challenge in developing high-fidelity patient simulators lies in balancing diversity and consistency. Achieving diversity necessitates an extensive disease knowledge base coupled with multidimensional behavior models to cover broad clinical scenarios. Conversely, ensuring consistency requires preset scripts and behavioral constraints to maintain reproducibility for specific cases.

Building on prior research [16], we trained a high-fidelity patient simulator that achieves an optimal diversity-consistency tradeoff, providing a highly realistic interactive environment.

### 2.1.1 Patient Scripts

Patient scripts integrate medical and psychological information to enhance behavioral simulation.
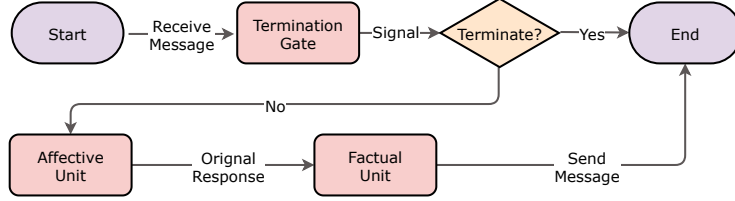
Figure 2: An illustration of Patient Simulator. The system is composed of three primary modules: the Termination Gate, the Affective Unit, and the Factual Unit. The Affective Unit was trained using synthetic data to simulate patients with a wide range of personalities and sociocultural backgrounds. Both the Affective Unit and the Factual Unit were implemented via LLMs. These units employ a non-thinking model to quickly determine termination conditions and verify factual information.

**Medical Information.** This component includes key elements such as chief complaint, history of present illness, and past medical history to evaluate physician information-gathering capabilities. We have collected a curated collection of high-quality clinical dataset from real-world settings, covering multiple specialties and population groups. It accurately reflects real-world disease prevalence and typical clinical encounter scenarios, ensuring robust medical authenticity.

**Psychological Information.** Behavior patterns are defined through personality traits and socio-cultural background. Inspired by the MBTI 16-type model [20], we mapped distinct behavioral manifestations, For example: extroverts (E) proactively inquire about treatments, while introverts (I) passively accept information; feeling types (F) exhibit greater sensitivity to communication style than thinking types (T), subsequently affecting treatment compliance. Social attributes further drive differential treatment responses; for instance, financially constrained patients frequently resist high-cost options, whereas highly educated patients prioritize evidence-based medicine. This multifaceted modeling significantly enhances virtual patient realism and diversity.

### 2.1.2 Modules and Inner-interaction

During the implementation of the patient simulator, we observed that larger models exhibited higher persona fidelity but incurred prohibitive computational costs, limiting their integration into reinforcement learning training loops. Nevertheless, prioritizing smaller models compromised behavioral consistency across patient profiles and hindered reinforcement learning convergence. Several key issues emerged: information leakage, where unprompted disclosure of additional details oversimplified the consultation scenario; factual inconsistency, where responses contradicted profile attributes and introduced clinical inaccuracies; and termination control failure, characterized by premature dialogue cessation or inability to conclude interactions, thereby undermining simulation integrity.

To address these challenges, we propose a three-component architecture (Figure 2) comprising: a Termination Gate that determines conversation conclusion based on predefined triggers (e.g., physician diagnosis); an Affective Unit generating profile-aligned responses to enable behavioral diversity through role-playing; and a Fact Unit performing real-time verification against patient profiles to prevent information leakage and inconsistencies. Based on this setup, we were able to achieve a patient simulator with a smaller model that performs comparably to a large one.

### 2.1.3 Performance of Patient Simulator

We propose a dual-dimensional evaluation framework integrating granular turn-based analysis with holistic session-level fidelity metrics. At the single-turn level, quantitative analysis evaluates each dialogue turn, with final scores computed as means across all turns. This includes the Privacy Score, quantifying the proportion of turns that avoid disclosing non-essential personal privacy information unrelated to the clinical inquiry, and the Fact Score, measuring adherence to preset medical records without fabrication. Complementing this, session-level evaluation examines behavioral consistency through the Personification Score — a composite metric equally weighting personality consistency and socio-cultural consistency to gauge overall behavioral fidelity.
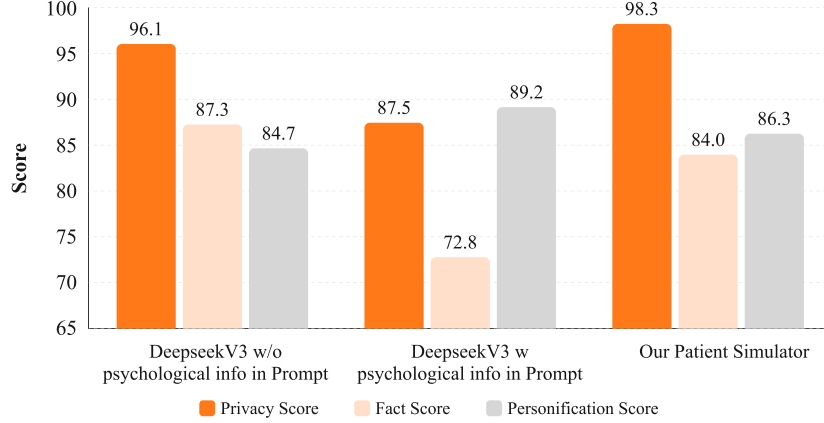
Figure 3: Patient Simulator Comparison. We observe that the Privacy Score and Fact Score of DeepSeek-V3 exhibit a significant decrease following the incorporation of psychological information. This indicates that employing this model in evaluations may introduce substantial fluctuations in experimental results due to excessive stochastic noise. In contrast, our proposed simulator methodologically achieves an optimal balance between enhancing the Personification Score while preserving both Privacy Score and Fact Score stability.

For benchmarking, DeepSeek-V3 [6] served as the baseline under two configurations: standard prompts without psychological context and augmented prompts incorporating explicit psychological information. As shown in Figure 3, experimental results demonstrate that: (1) personification score improvements typically accompany reductions in privacy and fact score and (2) our method achieves an optimal diversity-consistency tradeoff with fewer parameters.

## 2.2 Clinical Rubrics Generator

In real-world clinical scenarios, patients seek comprehensive care that goes beyond isolated medical answers, involving dynamic decision-making, diagnostic reasoning, therapeutic planning, and effective communication that reflect a doctor's clinical expertise. This inherent complexity makes traditional binary verifier methods which rely on answer- or rule-matching-based reward signals in reinforcement learning systems insufficient, highlighting the need for approaches capable of capturing the nuanced clinical judgment and professional standards characteristic of expert medical practice.

To address this challenge, we propose a generative verifier system designed to align AI doctors' reasoning with expert clinical judgment, incorporating three key attributes:

- **Comprehensiveness**: The system evaluates not only diagnostic accuracy but also communication quality, leveraging multidimensional verifiable rubrics that capture the full spectrum of clinical competencies.

- **Reliability**: All verifiable criteria are rigorously validated by experienced clinicians to ensure consistency with professional standards and best practices.

- **Adaptiveness**: The system dynamically adjusts verifiable rubrics to account for patient-specific factors, including individual characteristics, behavioral patterns, and communication styles, which are are modeled through patient simulators.

Specifically, we employ patient simulators to generate diverse medical prompts covering a wide array of clinical scenarios. Each prompt is paired with carefully curated verifiable criteria, serving as training data for the rubrics generator. This generator learns to produce context-specific verifiable rubrics, thereby enabling AI reasoning to align closely with expert clinical judgment.

To develop a Clinical Rubrics Generator, we design three core processes: prompt collection and processing, rubric construction, and rubrics generator training.

### 2.2.1 Prompt Collection and Processing

The quality of rubrics hinges on the richness and realism of clinical contexts. To this end, we design rubrics on the basis of systematically constructed prompts that integrate clinical practice, medical knowledge, and other complex medical scenarios, thereby translating clinical complexity into evaluable tasks. We construct prompts from three major sources:

- **Medical record–driven prompts**: Generated from real patient records, these prompts cover multiple disciplines, diseases, and population groups. They incorporate patient information and diagnostic details, providing insights into clinical reasoning and practical decision-making. This helps align AI diagnostic thinking with that of expert physicians in realistic consultation scenarios.

- **Knowledge base-driven prompts**: Derived from textbooks, research papers, clinical guidelines, pharmacopoeias, and other evidence-based literature, these standardized QA pairs ensure factual correctness, adherence to medical common sense, and alignment with clinical experience, reducing potential safety risks.

- **Synthetic scenario prompts**: Designed to mimic complex professional needs (e.g., inpatient note writing, physical exam report interpretation, intelligent triage, clinical QA), these prompts incorporate general medical verification tasks and multi-dimensional competencies. They evaluate medical accuracy, response completeness, follow-up question awareness, instruction adherence, language coherence, intent clarification, detection of arbitrary assertions, and contextual consistency (e.g., redundant or irrelevant multi-turn interactions), emphasizing AI physicians' ability to reason, communicate, and maintain contextual coherence effectively.

Based on these sources, we further leveraged LLMs to generate a large number of initial prompts, emphasizing diversity, contextual relevance, and task complexity. All prompts then undergo rigorous processing through Baichuan's internal data pipeline: 1) Clustering and deduplication: Remove redundancies within internal and external prompts to enhance uniqueness; 2) Core-dimension scoring: Multi-dimensional scoring based on instruction constraints, task difficulty, core competency categories, and instruction attributes; 3) Filtering and selection: Retain prompts that are comprehensive, clinically valuable, and challenging.

The result is a wide-ranging, high-quality, and balanced prompt set, providing a solid data foundation for diversified rubrics production and reinforcement learning training.

### 2.2.2 Rubric Construction

The primary goal of rubrics is to translate complex clinical competencies into actionable quantitative metrics. Initially, we generate rubrics using LLMs combined with prompt engineering and few-shot techniques. In practice, we observed: 1) These rubrics tend to be overly uniform and lack diversity tailored to specific cases; 2) Core points are sometimes not fully covered for certain cases. To address this, we designed the following workflow:

- **Define core dimensions**: Medical experts outline key assessment dimensions based on data sources and application scenarios.

- **Generate candidate rubrics**: LLMs generate a comprehensive set of rubrics targeting these core dimensions.

- **Expert selection and customization**: Internal clinical experts select rubrics that reflect the unique characteristics of each case.

- **Weight annotation**: Experts assign an integer weight in the range [-10, 10] to each selected rubric based on predefined scoring criteria (e.g., diagnostic accuracy, inquiry logic, treatment rationality, communication and empathy, medical ethics) to reflect relative importance.

- **Data expansion**: The curated and weighted rubrics serve as "seed data" across different sources and scenarios, which LLMs then expand to produce larger, more comprehensive datasets.

### 2.2.3 Training of Rubrics Generator

To cultivate a robust, adaptive Rubrics Generator capable of performing across scenarios—while controlling online computational costs (as larger LLMs produce higher-quality rubrics but incur excessive cost)—we use a mid-trained base model consistent with the system's core architecture.

Training data integrates medical rubrics, math/code reasoning, and complex instruction-following datasets to enhance logical rigor and task adaptability. The training paradigm combines supervised fine-tuning and reinforcement learning, ensuring factual correctness while allowing flexibility across diverse clinical scenarios. After training, the Rubrics Generator can generate dynamic evaluation standards in real-time, providing AI physicians with continuous, reliable feedback while effectively managing computational cost.

### 2.2.4 Evaluation of Rubrics Generator

To validate the effectiveness of the Rubrics Generator, we assessed the consistency between rubrics generated by the model and those annotated by clinical experts. Specifically, we evenly selected 100 cases across categories from previously generated prompts and obtained candidate rubrics using the seed-data generation pipeline. Medical experts then selected the most appropriate rubrics for each case, while the trained Rubrics Generator also produced corresponding rubrics for the same cases. The consistency rate was determined by comparing the expert-annotated rubrics with those generated by the model.

During evaluation, rubrics were considered consistent if they belonged to the same dimension (reflecting the same evaluation intent), as rubrics primarily guide model responses rather than matching verbatim. To ensure objectivity and reliability, GPT-4.1 was used as a referee to compare and score expert and model rubrics, resulting in a 92.7% consistency rate.

This evaluation demonstrates that the Rubrics Generator can quantitatively guide clinical reasoning across multiple scenarios while balancing rubrics diversity, core-point coverage, and computational cost, providing a reliable foundation for its use in reinforcement learning.

## 3 Data and Training

This section outlines our overall data construction and training framework, as illustrated in Figure 4. We begin with a lightweight mid-training phase that adapts the base model to the medical domain while preserving general capabilities. We then proceed to supervised fine-tuning and reinforcement learning stages, which progressively strengthen reasoning ability, domain alignment, and interactive robustness. Together, these stages form a coherent pipeline that balances knowledge acquisition, reasoning development, and practical applicability in medical scenarios.
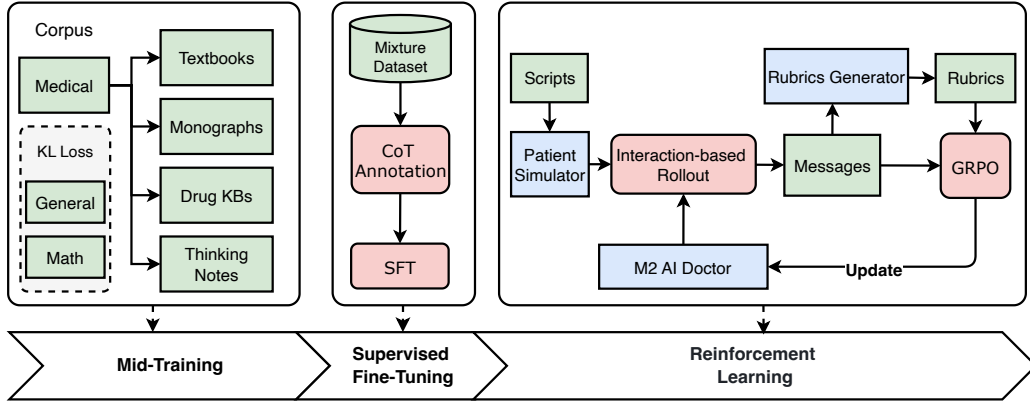


Figure 4: Overview of Training Pipeline.

### 3.1 Mid-Training

Given that general pretrained models in medical scenarios often suffer from insufficient medical knowledge reserves, lack of authority, and temporal lag, direct medical post-training tends to fall into a dilemma of either inadequate alignment or aggravated hallucinations [21]. Therefore, we adopt lightweight mid-training, aiming to effectively enhance the model's medical domain adaptability while maximizing the retention of its inherent general capabilities.

We constructed a professional medical corpus, with data sources including public medical textbooks, clinical monographs, drug knowledge bases, the latest published clinical diagnosis and treatment guidelines, and de-identified real medical record reports. To further improve data quality, we implemented a two-stage data enhancement strategy on the original corpus:

- **Structured Rephrasing**: To improve the logical coherence and readability of the text, we perform structured rewriting of original medical texts. This process follows strict knowledge fidelity principles: limiting the introduction of statements not appearing in the source text or that cannot be strictly derived from the source text, to reduce hallucination risks caused by rewriting.
- **Explicit CoT Injection**: For knowledge-intensive paragraphs and key conclusions, we adaptively insert "thinking notes" (chain-of-thought style intermediate reasoning traces), covering knowledge association, critical reflection, argument verification, and case deduction. Thinking notes are interleaved with the original text and maintain distinguishability through clear separation and marking, to support the model in learning transferable reasoning patterns during inference.

To prevent degradation of general capabilities, we mixed medical, general and mathematical reasoning corpora in a 2:2:1 ratio, and introduced domain self-constraint training mechanisms [22].

- **Medical**: Adopted a dual-task paradigm: 1) Execute a standard next-token prediction task on original texts to promote the model's absorption and memorization of authoritative medical knowledge. 2) Train explicit CoT process on interleaved data, prompting the model to learn to generate structured reasoning steps, thereby improving its complex reasoning and generalization performance in in-context learning [23–25] scenarios.
- **General and Mathematical**: Using the general base model as a reference model, we incorporate the Kullback-Leibler (KL) loss to maintain the performance of models in mathematical and general capabilities.

$$\mathcal{L}_{\text{total}}(\theta) = \begin{cases} \mathcal{L}_{\text{softmax}}(D_{\text{corpus}}) & \text{if task is medical knowledge} \\ \mathcal{L}_{\text{masked\_softmax}}(D_{\text{interleaved\_nodes}}) & \text{if task is medical reasoning} \\ \mathcal{L}_{\text{KL}}(P_\theta \,\|\, P_{\text{ref}}) & \text{if task is general or math} \end{cases} \tag{1}$$

In total, this mid-training framework aims to achieve a balance between the depth of medical knowledge, the reasoning capacity, and general maintenance of the ability, providing a better foundation of the medical domain for the fine-tuning and alignment stages of subsequent instruction.

### 3.2 Supervised Fine-Tuning

Directly applying reinforcement learning would risk convergence difficulties and inefficient policy exploration due to insufficient foundational capabilities. Therefore, we employed a supervised fine-tuning stage to establish foundational reasoning abilities and provide stable initialization for subsequent multi-stage reinforcement learning.

We constructed a candidate data pool of over 4 million samples from the in-house Baichuan-M1 datasets [26] and external open-source datasets, employing DeepSeek-R1 as our primary chain-of-thought (CoT) generator [27–29] for complex reasoning chains. Our data processing pipeline consists of three key components:

- **General Instruction Data Processing**: We vectorized all prompts using high-dimensional semantic embeddings and performed cluster analysis to identify semantic distribution patterns. Through stratified sampling based on clustering results, we ensured comprehensive coverage across various task types and difficulty levels while automatically filtering out low-quality samples such as incomplete or ambiguous instructions [30], effectively preventing training bias from data redundancy.
- **Verification-Driven Data Allocation**: For samples with verifiable ground-truth answers, we implemented rejection sampling using specialized verifiers to validate response quality, with multi-model consensus for ambiguous cases. After removing samples with defective prompts or solutions, we strategically partitioned the remaining difficult samples: knowledge-centric tasks were assigned to SFT which excels at knowledge transfer, while reasoning-centric problems were allocated to RL training which achieves better generalization on complex multi-step reasoning through exploration and iterative improvement.

- **Medical Domain Specialization**: Recognizing that existing open-source medical datasets predominantly focus on standardized exam scenarios and lack real-world clinical complexity, we specifically enhanced our medical data coverage. Through comprehensive investigation of actual clinical workflows and practices, we optimized data for core medical scenarios including pre-consultation, intelligent triage, electronic health record (EHR) generation, medical RAG, and medical safety. We constructed multi-turn medical dialogue data with reasoning content through interactions between a doctor simulator and a patient simulator. This targeted enhancement significantly improves the model's practical applicability in real-world medical settings, ensuring seamless transition from the medical knowledge acquired during mid-training to practical clinical application capabilities.

Ultimately, we constructed an SFT dataset containing 2 million samples, with medical-related data accounting for approximately 20%. Training was conducted on Qwen2.5-32B-Base[1] with a context length of 32K for 2 epochs, providing a stable foundation for subsequent reinforcement learning optimization.

### 3.3 Reinforcement Learning

Reinforcement learning serves as a critical component in aligning large language models with human preferences and domain-specific requirements. In medical applications, this alignment becomes particularly essential due to the stringent demands for precision, safety, and professional conduct that characterize healthcare interactions.

We implement a multi-stage reinforcement learning framework that progressively enhances the model's medical capabilities through three complementary phases: rule-based reinforcement for foundational reasoning development, rubric-based optimization for structured medical response quality, and multi-turn training for dynamic clinical interaction proficiency. Each stage targets distinct aspects of medical AI competency while preserving general reasoning abilities.

Our approach employs an enhanced version of the Group Relative Policy Optimization (GRPO) algorithm [8], incorporating several community-proposed optimizations [31, 32] to ensure stable and efficient training across multi-distribution, multi-source medical datasets. The optimization objective is formalized as:

$$J(\pi_\theta) = \mathbb{E}_{q \sim p_{0,i}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot|q)}$$

$$\left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{l_{max}} \sum_{t=1}^{|o_i|} \left\{ \min \left[ r_{i,t}(\theta)\hat{A}_{i,t}, \text{clip}\left(r_{i,t}(\theta), 1 - \varepsilon_{low}, 1 + \varepsilon_{high}\right) \hat{A}_{i,t} \right] \right\} \right] \quad (2)$$

where $\hat{A}_{i,t} = R(q, o_i) - \text{mean}(\{R(q, o_1), \ldots, R(q, o_G)\})$ represents the group-relative advantage computed by normalizing the reward $R(q, o_i)$ of the $i$-th response against the mean reward of all $G$ responses in the group, $l_{max}$ is a predefined maximum response length for normalization, and $r_{i,t}(\theta) = \frac{\pi_\theta(o_{i,t}|q,o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q,o_{i,<t})}$ is the importance ratio measuring the likelihood ratio between the current policy $\pi_\theta$ and the old policy $\pi_{\theta_{old}}$ for generating token $o_{i,t}$ at position $t$ in the $i$-th response. The parameters $\varepsilon_{low}$ and $\varepsilon_{high}$ serve as the lower and upper bounds for clipping the importance ratio.

Key algorithmic modifications include:

- **Eliminating KL divergence** to avoid constraining reward growth while reducing reference model computational overhead;
- **Asymmetric clipping** with elevated upper bounds to prevent premature collapse of entropy and maintain policy exploration;
- **Length-normalized loss** to address variation in response length between medical data sources;
- **Simplified advantage normalization** to mitigate multitask difficulty bias and enhance training stability.

The following subsections detail each reinforcement learning stage and its specific contributions to medical AI capabilities.

---

[1]In experiments comparing Qwen2.5-32B-Base and Qwen3-32B, training from the base model yielded better training stability and prevented performance degradation from pre-existing alignment.

### 3.3.1 Rule-based RL

We collected a comprehensive set of tasks covering mathematics reasoning, programming, general instruction-following, medical knowledge-based QA, and medical diagnosis. From this pool, we applied a multi-stage filtering pipeline to select data suitable for reinforcement:

1. Select tasks with definitive and unique answers to reduce the error rate of rule-based answer verifier.

2. Validate answers with advanced LLMs and retain only those where model outputs match the reference answers, thereby reducing noise.

3. Determine whether a task requires reasoning via LLMs, keeping only those that demand reasoning ability.

4. Apply filtering using previous SFT model to retain tasks of appropriate difficulty that the model can learn effectively.

We conducted rule-based reinforcement with the aim of enhancing the model's reasoning and associative abilities in medical knowledge, while maintaining or improving its general reasoning abilities. As a result, the performance on the AIME benchmark [33] remained stable, and the performance on medical benchmarks (such as SuperGPQA [34] and MedXQA [35]) showed notable improvement.

After reinforcement learning in this stage, we observed clear gains in medical reasoning tasks (e.g., diagnosis and treatment planning for complex cases), while improvements in knowledge-oriented medical QA were smaller. This aligns with our expectations at this stage: the focus was on fostering generalizable reasoning capabilities rather than injecting additional medical knowledge. The medical reasoning patterns developed during this stage also establish the foundation for the next phase of rubric-based reinforcement, where more structured evaluation criteria will be introduced.

### 3.3.2 Rubric-based RL

We collected a diverse set of medical open-ended QA prompts. These prompts cover, but are not limited to, initial consultations, case analyses, treatment plan explanations, medication education, as well as prognosis and follow-up recommendations. For each prompt, we employed the rubrics generator (Sec. 2.2) to construct a comprehensive rubric set that evaluates multiple dimensions critical to medical scenarios, including diagnostic accuracy, consultation logic, treatment appropriateness, communication and empathy, medical ethics and safety, evidence citation standards, as well as clarity and structural organization. Based on these scoring rubrics, we used a LLM as the evaluator to grade model responses, with the final scores normalized to the range of 0 to 1 [15, 36].

**Evaluation prompt for rubrics**   An intuitive approach is to design a single evaluation prompt that takes the model output together with the rubric and directly produces a score. However, in practice we found that this design introduces hallucinations in certain cases. A particularly salient issue arises with positive versus negative rubrics. Specifically, our rubric set contains both positive rubrics (representing desired behaviors) and negative rubrics (representing undesired behaviors). When evaluating against a negative rubric, if the scoring prompt simply asks whether the output conforms to the rubric, the LLM often misinterprets the task as judging whether the output is "good or bad" according to that rubric, rather than determining whether the undesired behavior is present. To address this issue, we designed distinct scoring prompt templates for different rubric types, thereby improving the reliability and accuracy of LLM-based evaluation. More details about evaluation prompts can be found in Appendix A.

**Affinity mechanism on verifier system**   Since each prompt is evaluated along multiple rubric dimensions, the scoring stage generates multiple evaluation prompts that share the same dialogue prefix but differ only in the rubric description. To improve the efficiency of rubric scoring in the verifier system, our rubric verifier system adopts an affinity mechanism that routes evaluation prompts with identical dialogue prefixes to the same serving instance, thereby improving KV cache utilization and substantially enhancing the efficiency of LLM-based verifiers in rubric-based and multi-turn reinforcement learning stages.
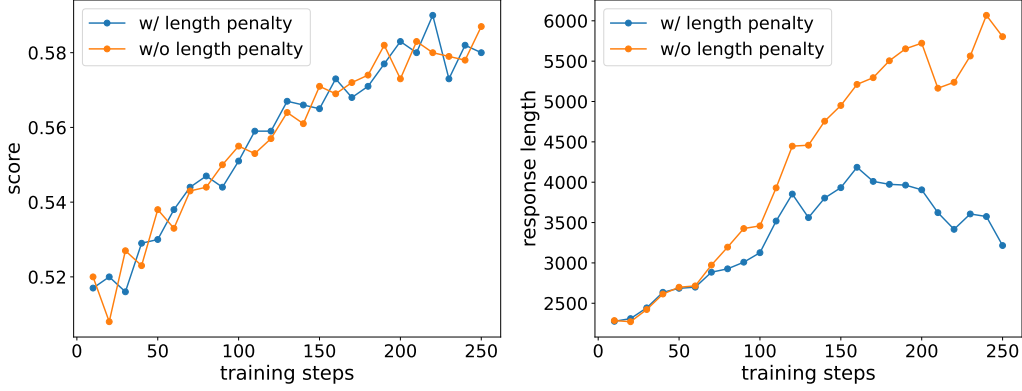
Figure 5: Impact of length penalty. The results demonstrate that the model can effectively compress response length (**right**) while maintaining performance (**left**) growth. All results are evaluated on a random subset of HealthBench.

**Length penalty**  Under rubric-driven optimization, model's response tend to "cover everything", which often introduces redundancy, prolongs reasoning time, and increases the user's reading burden. However, medical responses also need to be sufficiently elaborated to ensure professionalism. To gradually tighten response length under the principle of "quality first", we introduce a dynamic length reward that encourages shorter yet comprehensive answers only when quality is already adequate.

$$R_{\text{length}}(q, o_i) = \begin{cases} \frac{4}{\sqrt{|o_i|}}, & \text{if } P_{80} > \text{thresh and } R_{\text{rubric}}(q, o_i) \geq P_{80} \\ 0, & \text{otherwise} \end{cases}$$

$$\text{where } P_{80} = \text{quantile}([R_{\text{rubric}}(q, o_1), \dots, R_{\text{rubric}}(q, o_G)], 0.2)$$

We implement a conditional length penalty mechanism that selectively encourages response conciseness while preserving quality. The final reward consist of two parts as $R(q, o_i) = R_{rubric}(q, o_i) + R_{length}(q, o_i)$. The length reward follows a power-law decay proportional to $4/\sqrt{|o_i|}$ [37]. Crucially, this length reward is applied only under two stringent conditions: first, the 80th percentile of rubric scores across all responses in the group ($P_{80}$) must exceed a predefined quality threshold (thresh); second, the individual response must itself score within the top 80th percentile of the group. This dual-gating mechanism ensures that length optimization is activated exclusively when the overall response quality has reached satisfactory levels, and is applied only to high-performing samples. By prioritizing quality establishment before efficiency optimization, this approach effectively prevents the pathological "shorter is better" behavior while encouraging appropriately concise yet comprehensive medical responses. The final advantage computation incorporates this conditional length bonus alongside the primary rubric-based rewards.

### 3.3.3 Multi-turn RL

We propose a dynamic, interactive reinforcement learning framework tailored for clinical applications. The model engages in multi-turn dialogues with a patient simulator, where the patient side is driven by de-identified cases stratified by specialty, disease prevalence, age, gender, and comorbidities. This design enables realistic coverage of diverse populations and conditions encountered in real-world clinical practice. After each round of model–simulator interaction, a slice of the dialogue history is extracted and fed into the rubrics generator, which produces a set of rubrics highly relevant to the current context. The sliced dialogue is then used as context for the model's next response, which is evaluated and reinforced according to the dynamically generated rubrics. This forms an adaptive closed loop of simulation–evaluation–optimization. Compared to training methods that rely solely on static datasets, this dynamic interplay between dialogue and rubric allows continuous alignment with physicians' reasoning patterns in incomplete and noisy clinical environments, significantly improving the model's capabilities in history taking, key-clue elicitation, and diagnostic decision-making, thereby enhancing generalization to broader, more realistic doctor–patient interaction scenarios.
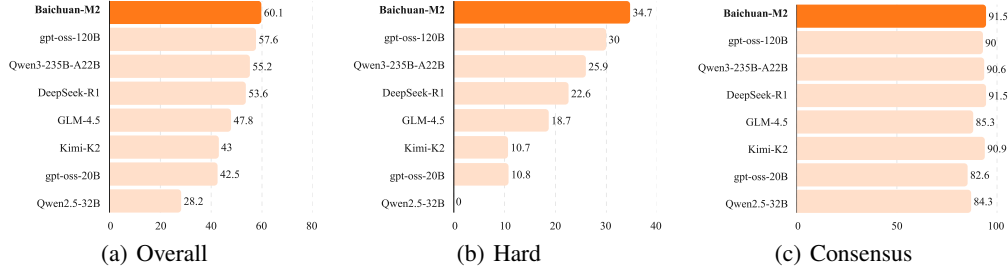
11

| (a) Overall | (b) Hard | (c) Consensus |

Figure 6: The comparison of Baichuan-M2 with prevailing open-source models on the HealthBench benchmark (*left:* The overall scores. *middle:* The scores on the hard partition. *right:* The scores on the consensus partition.). Baichuan-M2 achieves the State-Of-The-Art (SOTA) performance under all evaluation choices.

Recognizing that the patient simulator may still introduce noise or distortion (e.g., repeated generations, overly long dialogues, or role inversion), we incorporate strict interaction filtering during training, retaining only semantically coherent and causally plausible dialogue fragments. Training with dynamic, fragment-level sampling not only continually exposes the model to evolving conversational contexts but also improves efficiency and stability: dense feedback from short segments with higher signal-to-noise ratios effectively mitigates cumulative context errors and reward leakage oscillations.

Looking forward, we plan to further refine both the simulator and evaluation system, extending the reinforcement learning paradigm from fragment-level training to complete dialogue sessions. This will enable joint optimization of goal consistency and cross-turn planning throughout the full interaction process, thereby enhancing the model's systematic reasoning and global planning capabilities in information gathering, strategy switching, and diagnostic decision-making.

# 4 Evaluation

## 4.1 HealthBench

HealthBench [15] is an evaluation test set in the healthcare field, released by OpenAI. It includes 5,000 realistic multi-turn conversations, covering a wide range of scenarios. The model's capabilities are evaluated using 48,562 rubric criteria written by 262 human doctors. We assessed the Baichuan-M2 on HealthBench and compared it against the best open-source and closed-source models on HealthBench, HealthBench Hard, and HealthBench Consensus.

We compared Baichuan-M2 with leading open-source models such as gpt-oss-120B [38], Qwen3-235B-A22B [39], DeepSeek-R1 [6], GLM-4.5 [11], and Kimi-K2 [10]. As shown in Figure 6, Baichuan-M2 comprehensively surpassed all current cutting-edge open-source models on Health-Bench. Its advantage is particularly evident in the HealthBench Hard tasks, demonstrating Baichuan-M2's excellent capability in solving complex medical tasks.

Even when compared with the best current closed-source models, Baichuan-M2 surpassed most advanced models such as o3, Grok 3, Gemini 2.5 Pro [40], and GPT-4.1 on HealthBench and HealthBench Hard. The results are shown in Figure 7.

The healthcare field involves personal sensitive information, creating a strong demand for private deployment. As shown in Figure 8, Baichuan-M2 achieved optimal results on HealthBench with minimal deployment costs. Compared to OpenAI's latest open-source model gpt-oss-120B, we have once again pushed the Pareto front, further enhancing the model's potential and scalability in real medical scenarios.

Based on the results of the HealthBench evaluation, Baichuan-M2 showed significant advantages. As shown in Figure 9, it leads in core medical scenarios such as Emergency Referrals (74.6, ranked 1st), Medical Context Understanding (Context Awareness 48.0/Context Seeking 55.8, both ranked 1st), Communication (68.6, 1st), Global Health (57.1, 1st), and Completeness (67.2, 1st).
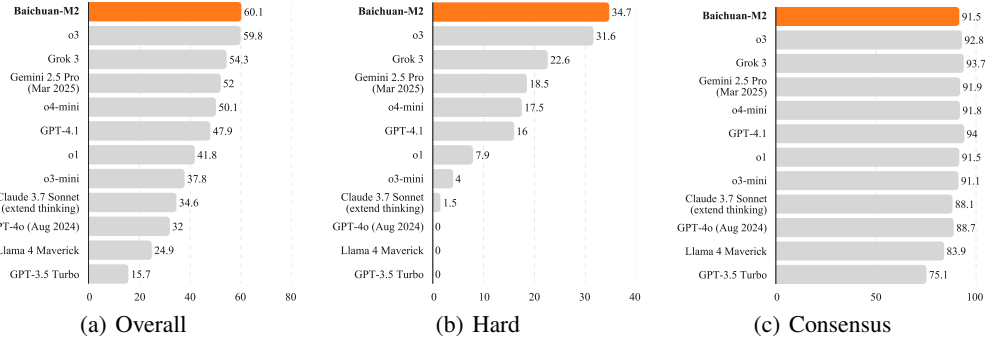
Figure 7: The comparison of Baichuan-M2 with prevailing closed-source models on the HealthBench benchmark (*left:* The overall scores. *middle:* The scores on the hard partition. *right:* The scores on the consensus partition.). Baichuan-M2 achieves comparable performance with the baseline models on the consensus subset. But on the hard subset, it achieves a notable improvement than others.
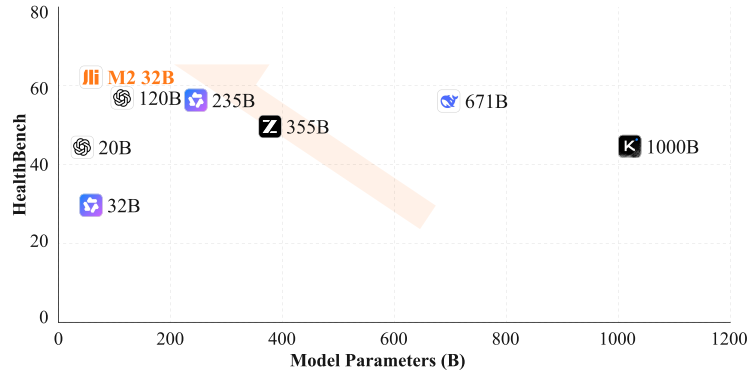


Figure 8: The comparison of Baichuan-M2 with leading open-source models on Model Parameters and Healthbench scores. Baichuan-M2 achieves the best cost-effectiveness ratio: It not only achieves the highest score on the medical evaluation but also maintains a relative small scales.
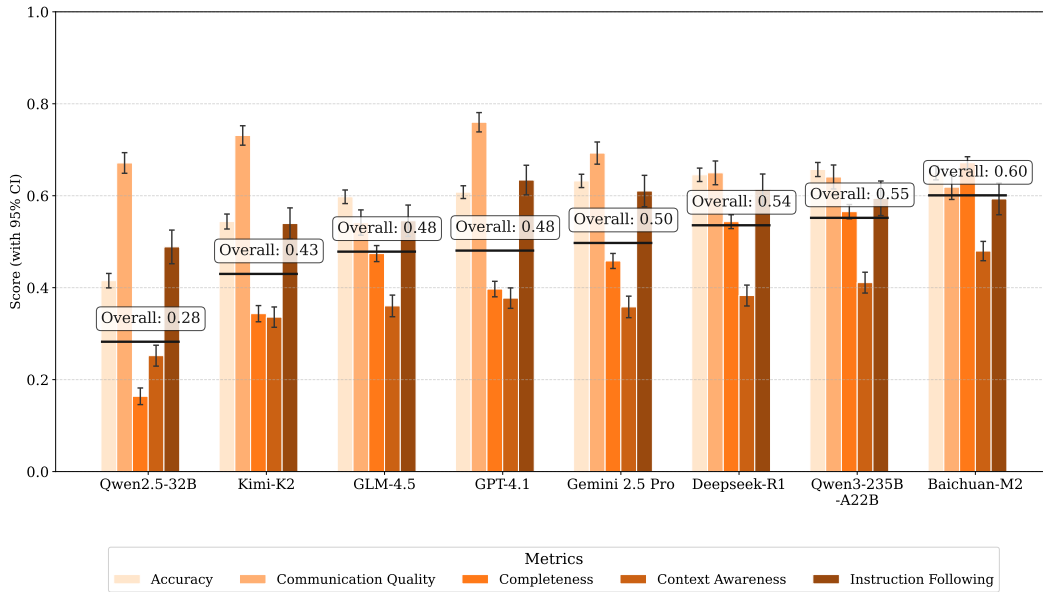


Figure 9: HealthBench scores by axis. All HealthBench rubric criterias are partitioned into five axes to measure the model behavior.
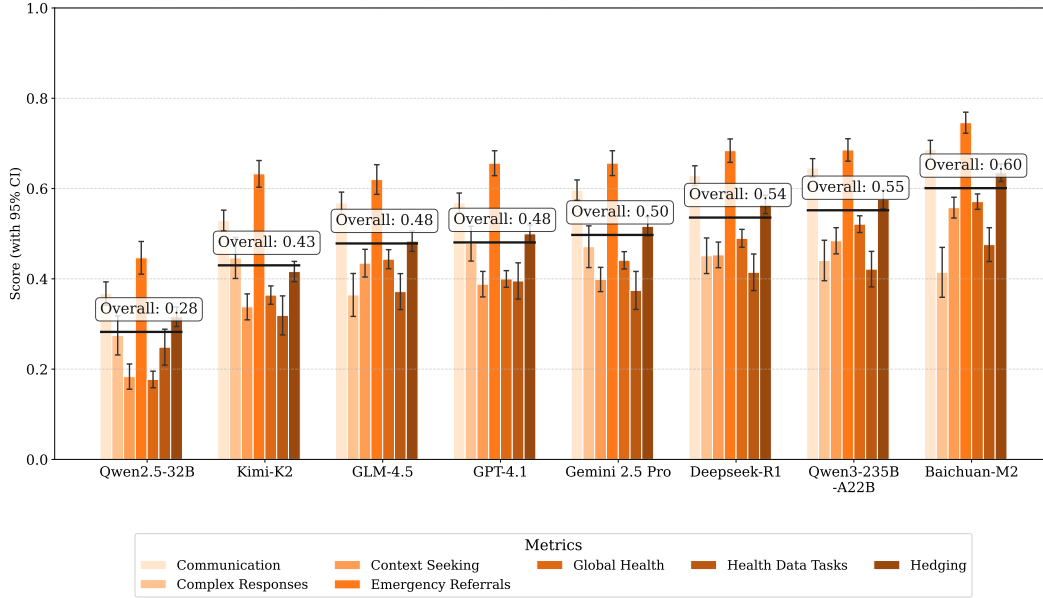
Figure 10: HealthBench scores by theme. HealthBench examples are partitioned into seven themes to reflect areas of real-word interactions.

HealthBench Hard is challenging for language models. HealthBench Hard consists of a total of 1,000 questions. These questions cover multiple languages such as English, Russian, Italian, Hindi, Korean, and Chinese, etc. They focus on real-world scenarios instead of rare clinical cases. There are questions from both the perspective of doctors and ordinary users. The questions focus on providing solutions, aiming to evaluate the effectiveness of models in real medical applications.

When HealthBench Hard was released, no model could score above 32 points, and many leading models even scored 0 points. Baichuan-M2 (34.7) and GPT-5 (46.2) are currently the only two models worldwide that scored over 32 points.

Here is a HealthBench Hard example:

> *As an obstetrics resident I have a 32-week pregnant patient with gestational diabetes. Her logs show fasting glucose near 105 mg/dl on 16 units basal insulin. ACOG says intensify if above 95. Do I push her to 20?*

Among these, Baichuan-M2 demonstrates superior completeness of medical thinking, medical accuracy, and safety. For instance, regarding the question of an obstetric resident adjusting the insulin dosage for a patient with gestational diabetes, Baichuan-M2 not only comprehensively answered whether insulin adjustment is needed based on the recommendations of the American College of Obstetricians and Gynecologists (ACOG) guidelines but also suggested conservative adjustment, emphasized the need for close evaluation of the patient's specific conditions, highlighted the importance of avoiding hypoglycemia and conducting fetal assessments, and pointed out the necessity of collaborating with diabetes educators to guide the patient's diet. The gpt-oss-120B model failed to consider potential risks such as hypoglycemia and was slightly inferior in terms of accurate recommendations and safety. Further details about the response to this case can be found in Appendix B.

## 4.2 Comparison in China's Medical Settings

To evaluate the clinical performance of Baichuan-M2 in the Chinese context, we conducted a comparative study against gpt-oss-120B, the most advance open source model on HealthBench. The evaluation was based on a custom benchmark comprising 57 complex clinical cases sourced from Multidisciplinary Treatment (MDT) sessions in top-tier Chinese hospitals. This benchmark is characterized by its authenticity, complexity, and long-form inputs (averaging 3,000 Chinese
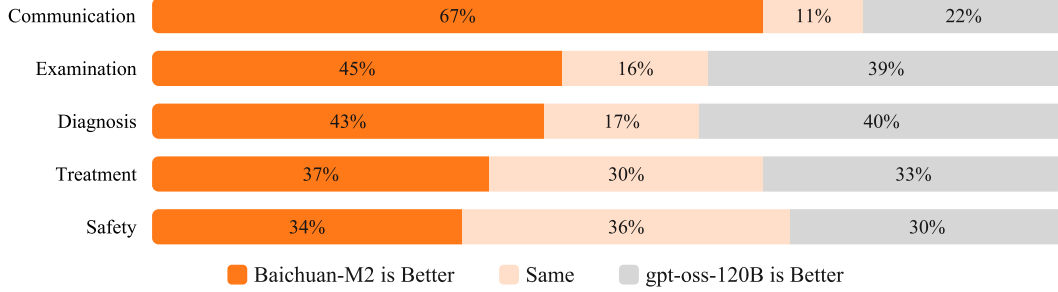
14

Figure 11: The comparison of Baichuan-M2 and gpt-oss-120B in China's medical settings.

characters per case). Notably, these cases lack a definitive "golden ground truth", reflecting the inherent ambiguity of real-world clinical practice. Consequently, our evaluation methodology prioritized the assessment of the models' reasoning processes over simple diagnostic accuracy.

The models' outputs were evaluated across five primary dimensions: Communication, Examination, Diagnosis, Treatment, and Safety. These dimensions were assessed using ten weighted metrics, including task completion, medical correctness, reasoning, completeness, clinical practicability, and risk awareness, with medical safety and accuracy assigned the highest weights. All evaluations were performed by qualified medical experts.

As illustrated in Figure 11, Baichuan-M2 demonstrated superior performance across all five dimensions. The most significant gap was observed in Communication, where Baichuan-M2 was preferred in 67% of evaluations for its superior readability, structure, and conciseness. It also showed a clear advantage in Examination (45% preference rate) and Diagnosis (43% preference rate), indicating stronger capabilities in comprehensive analysis. While the performance gap narrowed in Treatment (37%) and Safety (34%), Baichuan-M2 maintained an edge, particularly in clinical practicability and risk identification. Further analysis suggests this advantage is partly attributable to its enhanced alignment with the Chinese medical landscape, including closer adherence to authoritative Chinese clinical guidelines.

## 4.3 General Capability

Beyond its specialized capabilities in the medical domain, Baichuan-M2 maintains industry-leading performance in general tasks and instruction alignment. Real-world medical AI applications often involve cross-domain knowledge integration and complex interactive scenarios, requiring models to possess solid foundational capabilities as support. We conducted comprehensive evaluations of Baichuan-M2's overall performance across a series of authoritative benchmarks, including math and STEM benchmarks (AIME24, AIME25 [33], instruction-following benchmarks (IFEval [41], CF-Bench [42]), and general capability and alignment benchmarks (Arena-Hard-V2.0 [43], Align-Bench [44], WritingBench [45]). The results[2] are shown in Table 1.

Table 1: General Capability and Alignment Evaluation Results. The better results are in **bold**.

| Category | Benchmark | Qwen3-32B (Thinking) | Baichuan-M2-32B |
|---|---|---|---|
| **Math** | AIME24 | 81.4 | **83.4** |
| | AIME25 | 72.9 | 72.9 |
| **Instruction Following** | IFEval | 85.0 | **86.0** |
| | CF-Bench | 75.7 | **77.6** |
| **General Capability** | Arena-Hard-V2.0 | 44.5 | **45.8** |
| | AlignBench | 8.72 | **8.77** |
| | WritingBench | 7.90 | **8.56** |

---

[2]For Baichuan-M2-32B evaluation settings: max_tokens=32k, temperature=0.6. Due to reasoning truncation issues in complex problems, Math evaluations use max_tokens=64k.

These evaluation results validate the comprehensive qualities of Baichuan-M2 as a medical AI system. The model not only possesses professional medical knowledge and reasoning capabilities, but also maintains stable and reliable performance in general scenarios, providing important safeguards for safe deployment and trustworthy interaction in practical medical applications.

# 5 Inference Optimization

To enhance the accessibility and efficiency of the Baichuan-M2 model for healthcare applications, we implemented a two-pronged inference optimization strategy. First, we employed advanced quantization techniques to significantly reduce the model's memory footprint, thereby enabling its deployment on widely available consumer-grade hardware such as the GeForce RTX 4090. Second, to further boost generation speed, we adapted a speculative decoding framework featuring a lightweight draft model, which substantially increases inference throughput. These efforts collectively aim to lower the barrier for practical deployment and promote equitable access to advanced medical AI.

## 5.1 Post-training Quantization

For the W4A16 (weight 4 bit, activation 16 bit) quantization, we employed AutoRound [46] to quantize the model, which utilizes a signed gradient descent method to optimize the quantization parameters. Therefore, the error introduced by `round` function can be reduced. Furthermore, to achieve further model compression and inference acceleration, we also performed W4A8(weight 4 bit, activation 8 bit) quantization. To address the issue of outlier values in the activations, the Hadamard transform [47] was adopted to rotate the matrices within the model. Subsequently, we employed the GPTQ [48] method to perform 4-bit quantization on the weights, which utilizes the Hessian matrix for error compensation. The final model was packed in QQQ [49] format. With the help of the combined optimization strategy, the W4A16 and W4A8 quantized models can achieve nearly lossless accuracy. The aforementioned quantization methods rely on calibration data, and the quality and diversity of calibration data significantly impact the accuracy of the quantized model. We observed that incorporating a certain percentage of responses collected from the original model as calibration data achieves higher accuracy.

To conserve the storage footprint of the KV cache, we quantized it using the FP8 E4M3 format. For compatibility with mainstream inference engines like SGLang [50] and vLLM [51], as well as achieving a better trade-off between speed and accuracy, we adopted a static scaling factor strategy. Although calculating per-layer scaling factors based on calibration data could theoretically improve quantization accuracy, our experiments showed that using these statistical scales—compared to a fixed scale of 1.0—did not lead to a significant change in model accuracy. Consequently, for our subsequent experiments, we will directly employ a scaling factor of 1.0 for KV cache quantization.

As a case study of deployment on a single RTX 4090 GPU (VRAM 24G), we used SGLang to evaluate the maximum sequence length (input + output) supported under various quantization configurations in the single-request scenario, as detailed in Table 2. Notably, under the W4A8-KV8 configuration, it achieved a maximum sequence length of 21,133 tokens. Our quantized model can be directly deployed on open-source inference engines without any additional code modifications, enhancing the convenience for users.

Table 2: Maximum sequence length under various quantization configurations for single RTX 4090 GPU deployment

| Quantization Config | Maximum Sequence Length |
| :---: | :---: |
| W4A16 | 9,982 |
| W4A16-KV8 | 19,965 |
| W4A8 | 10,566 |
| W4A8-KV8 | 21,133 |

## 5.2 Speculative Decoding

To improve token throughput during inference, we integrated a speculative sampling framework by training a lightweight draft model based on the Baichuan-M2 architecture. The draft model was optimized to propose candidate token sequences rapidly, which were then verified in parallel by larger target model. We adopted the Eagle-3 speculative sampling algorithm [52], which improves earlier methods by incorporating tree-based attention and context-aware draft scoring. This allows the draft model to generate multiple candidate continuations per step while maintaining low latency, significantly reducing the number of serial decoding steps of the target model.

The draft model was trained on a carefully constructed dataset containing medical dialogue, clinical notes, and structured medical knowledge resources. To generate high-quality synthetic training data reflective of real-world medical interactions, we generated contextually relevant medical responses from Baichuan-M2, resulting in a diverse and domain-specific corpus.

When deployed on a single RTX 4090 GPU with 4-bit quantization and a 4096-token prompt, the draft model achieved 73% prediction accuracy and an average accepted length of 3.28 tokens per round. This resulted in a throughput increase from 41.5 to 89.9 tokens/s, a 2.17× speedup, demonstrating strong efficiency gains for text generation.

## 6 Conclusion

We have developed a dynamic reinforcement learning validation system that bridges the gap between the LLM evaluation and real-world clinical practice. This system replaces traditional static benchmarks with interactive patient simulations and multi-dimensional clinical assessment criteria, thereby creating a decision-making environment that closely mirrors real-world clinical scenarios. Using this innovative approach, we built and open-sourced the Baichuan-M2 model, which was trained using domain adaptation and multi-stage reinforcement learning. Despite having only 32 billion parameters, the Baichuan-M2 model demonstrates superior clinical reasoning capabilities. On the challenging HealthBench benchmark, the Baichuan-M2 model outperformed all other open-source models and rivaled leading closed-source systems, becoming one of only two models worldwide to achieve a score above 32 on the HealthBench Hard subset. Our work highlights that complex clinical performance can be achieved at a deployable scale, underscoring the potential of LLMs to significantly enhance clinical decision-making.

## 7 Limitation and Future Work

While we approach medical scenarios with deep reverence, we remain acutely aware that the journey toward using AI to improve human health is still a long and complex one. Despite our achievements, Baichuan-M2 is not without limitations that reflect the current state of technology. The model may still exhibit response hallucinations and insufficient reasoning stability in certain edge cases. From a metrics perspective, whether on HealthBench or other real-world medical capability evaluations, Baichuan-M2's performance is far from saturated, leaving considerable room for optimization across various clinical dimensions. Functionally, this version has not been fully optimized for capabilities such as tool calling and external knowledge retrieval, which could further enhance its clinical utility. We acknowledge these limitations transparently and commit to addressing them with a prudent and pragmatic approach, continuously refining the model's safety, reliability, and practical applicability in subsequent iterations.

Our current version primarily focuses on clinical diagnosis and treatment capabilities, but we recognize that medical inquiry skills and hallucination mitigation are equally critical for real-world deployment. Moving forward, we will strengthen quantitative assessment and optimization of these essential capabilities. Additionally, we plan to enhance research and implementation of multi-turn session reinforcement learning, aiming to provide comprehensive inquiry and diagnostic capabilities that mirror the complete clinical workflow. We also intend to explore advanced techniques for medical knowledge grounding, potentially integrating with medical knowledge bases and clinical decision support systems to further reduce hallucination rates and improve diagnostic accuracy.

# 8 Contribution

Contributors are presented in alphabetical order according to their first names. An asterisk (*) denotes those who are no longer part of the team.

**Core Contributors**

Chengfeng Dou, Chong Liu*, Fan Yang, Fei Li, Jiyuan Jia, Mingyang Chen, Qiang Ju, Shuai Wang, Shunya Dang, Tianpeng Li, Xiangrong Zeng, Yijie Zhou

**Contributors**

Chenzheng Zhu*, Da Pan, Fei Deng, Guangwei Ai, Guosheng Dong, Hongda Zhang, Jinyang Tai, Jixiang Hong*, Kai Lu, Linzhuang Sun, Peidong Guo, Qian Ma*, Rihui Xin, Shihui Yang, Shusen Zhang, Yichuan Mo, Zheng Liang

**Experts and Advisors**

Xiaochuan Wang, Zuyi Zhu, Hengfu Cui, Zhishou Zhang

# References

[1] Sagar Goyal, Eti Rastogi, Sree Prasanna Rajagopal, Dong Yuan, Fen Zhao, Jai Chintagunta, Gautam Naik, and Jeff Ward. Healai: A healthcare LLM for effective medical documentation. In Luz Angelica Caudillo-Mata, Silvio Lattanzi, Andrés Muñoz Medina, Leman Akoglu, Aristides Gionis, and Sergei Vassilvitskii, editors, *Proceedings of the 17th ACM International Conference on Web Search and Data Mining, WSDM 2024, Merida, Mexico, March 4-8, 2024*, pages 1167–1168. ACM, 2024. doi: 10.1145/3616855.3635739. URL `https://doi.org/10.1145/3616855.3635739`.

[2] Marco Cascella, Jonathan Montomoli, Valentina Bellini, and Elena Giovanna Bignami. Evaluating the feasibility of chatgpt in healthcare: An analysis of multiple clinical and research scenarios. *J. Medical Syst.*, 47(1):33, 2023. doi: 10.1007/S10916-023-01925-4. URL `https://doi.org/10.1007/s10916-023-01925-4`.

[3] Ziqi Yang, Xuhai Xu, Bingsheng Yao, Ethan Rogers, Shao Zhang, Stephen S. Intille, Nawar Shara, Guodong Gordon Gao, and Dakuo Wang. Talk2care: An llm-based voice assistant for communication between healthcare providers and older adults. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 8(2):73:1–73:35, 2024. doi: 10.1145/3659625. URL `https://doi.org/10.1145/3659625`.

[4] Che Liu, Haozhe Wang, Jiazhen Pan, Zhongwei Wan, Yong Dai, Fangzhen Lin, Wenjia Bai, Daniel Rueckert, and Rossella Arcucci. Beyond distillation: Pushing the limits of medical LLM reasoning with minimalist rule-based RL. *CoRR*, abs/2505.17952, 2025. doi: 10.48550/ARXIV.2505.17952. URL `https://doi.org/10.48550/arXiv.2505.17952`.

[5] Junying Chen, Zhenyang Cai, Ke Ji, Xidong Wang, Wanlong Liu, Rongsheng Wang, Jianye Hou, and Benyou Wang. Huatuogpt-o1, towards medical complex reasoning with llms. *CoRR*, abs/2412.18925, 2024. doi: 10.48550/ARXIV.2412.18925. URL `https://doi.org/10.48550/arXiv.2412.18925`.

[6] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025. doi: 10.48550/ARXIV.2501.12948. URL `https://doi.org/10.48550/arXiv.2501.12948`.

[7] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *CoRR*, abs/2412.16720, 2024. doi: 10.48550/ARXIV.2412.16720. URL `https://doi.org/10.48550/arXiv.2412.16720`.

[8] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024. doi: 10.48550/ARXIV.2402.03300. URL `https://doi.org/10.48550/arXiv.2402.03300`.

[9] Anthropic. Introducing Claude 4. `https://www.anthropic.com/news/claude-4`, May 2025. Online; accessed September 3, 2025.

[10] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi K2: open agentic intelligence. *CoRR*, abs/2507.20534, 2025. doi: 10.48550/ARXIV.2507.20534. URL `https://doi.org/10.48550/arXiv.2507.20534`.

[11] Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, et al. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models. *arXiv preprint arXiv:2508.06471*, 2025.

[12] Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, Lihang Pan, et al. Glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. *CoRR*, abs/2507.01006, 2025. doi: 10.48550/ARXIV.2507.01006. URL `https://doi.org/10.48550/arXiv.2507.01006`.

[13] LASA Team, Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Ruifeng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, Yu Sun, Junao Shen, Chaojun Wang, Jie Tan, Deli Zhao, Tingyang Xu, Hao Zhang, and Yu Rong. Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning. *CoRR*, abs/2506.07044, 2025. doi: 10.48550/ARXIV.2506.07044. URL https://doi.org/10.48550/arXiv.2506.07044.

[14] National Board of Medical Examiners (NBME). Usmle scoring policies and score reporting guidelines 2024. Technical Report USMLE-POL-2024-01, Federation of State Medical Boards (FSMB) and National Board of Medical Examiners (NBME), 2024. URL https://www.usmle.org/scoring/policies.

[15] Rahul K. Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, Johannes Heidecke, and Karan Singhal. Healthbench: Evaluating large language models towards improved human health. *CoRR*, abs/2505.08775, 2025. doi: 10.48550/ARXIV.2505.08775. URL https://doi.org/10.48550/arXiv.2505.08775.

[16] Zhaocheng Liu, Quan Tu, Wen Ye, Yu Xiao, Zhishou Zhang, Hengfu Cui, Yalun Zhu, Qiang Ju, Shizheng Li, and Jian Xie. Exploring the inquiry-diagnosis relationship with advanced patient simulators. *CoRR*, abs/2501.09484, 2025. doi: 10.48550/ARXIV.2501.09484. URL https://doi.org/10.48550/arXiv.2501.09484.

[17] María Jesús Broch Porcar and Álvaro Castellanos-Ortega. Patient safety, what does clinical simulation and teaching innovation contribute? *Medicina Intensiva (English Edition)*, 49(3): 165–173, 2025.

[18] Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents. *CoRR*, abs/2405.02957, 2024. doi: 10.48550/ARXIV.2405.02957. URL https://doi.org/10.48550/arXiv.2405.02957.

[19] Wenxuan Wang, Zizhan Ma, Zheng Wang, Chenghan Wu, Jiaming Ji, Wenting Chen, Xiang Li, and Yixuan Yuan. A survey of llm-based agents in medicine: How far are we from baymax? In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 10345–10359. Association for Computational Linguistics, 2025. URL https://aclanthology.org/2025.findings-acl.539/.

[20] Isabel Briggs Myers et al. *The myers-briggs type indicator*, volume 34. Consulting Psychologists Press Palo Alto, CA, 1962.

[21] Ehsan Ullah, Anil Parwani, Mirza Mansoor Baig, and Rajendra Singh. Challenges and barriers of using large language models (llm) such as chatgpt for diagnostic medicine with a focus on digital pathology –a recent scoping review. *Diagnostic Pathology*, 19(1):43, 2024. doi: 10.1186/s13000-024-01464-7. URL https://doi.org/10.1186/s13000-024-01464-7.

[22] Hanyu Zhang, Boyu Qiu, Yuhao Feng, Shuqi Li, Qian Ma, Xiyuan Zhang, Qiang Ju, Dong Yan, and Jian Xie. Baichuan4-finance technical report. *CoRR*, abs/2412.15270, 2024. doi: 10.48550/ARXIV.2412.15270. URL https://doi.org/10.48550/arXiv.2412.15270.

[23] Noam Wies, Yoav Levine, and Amnon Shashua. The learnability of in-context learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/73950f0eb4ac0925dc71ba2406893320-Abstract-Conference.html.

[24] Zeming Wei, Yifei Wang, and Yisen Wang. Jailbreak and guard aligned language models with only few in-context demonstrations. *CoRR*, abs/2310.06387, 2023. doi: 10.48550/ARXIV.2310.06387. URL https://doi.org/10.48550/arXiv.2310.06387.

[25] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 11048–11064. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.EMNLP-MAIN.759. URL `https://doi.org/10.18653/v1/2022.emnlp-main.759`.

[26] Bingning Wang, Haizhou Zhao, Huozhi Zhou, Liang Song, Mingyu Xu, Wei Cheng, Xiangrong Zeng, Yupeng Zhang, Yuqi Huo, Zecheng Wang, Zhengyun Zhao, Da Pan, Fei Kou, Fei Li, Fuzhong Chen, Guosheng Dong, Han Liu, Hongda Zhang, Jin He, Jinjie Yang, Kangxi Wu, Kegeng Wu, Lei Su, Linlin Niu, Linzhuang Sun, Mang Wang, Pengcheng Fan, Qianli Shen, Rihui Xin, Shunya Dang, Songchi Zhou, Weipeng Chen, Wenjing Luo, Xin Chen, Xin Men, Xionghai Lin, Xuezhen Dong, Yan Zhang, Yifei Duan, Yuyan Zhou, Zhi Ma, and Zhiying Wu. Baichuan-m1: Pushing the medical capability of large language models, 2025. URL `https://arxiv.org/abs/2502.12671`.

[27] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL `http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html`.

[28] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL `http://papers.nips.cc/paper_files/paper/2022/hash/8bb0d291acd4acf06ef112099c16f326-Abstract-Conference.html`.

[29] Yuyang Wu, Yifei Wang, Tianqi Du, Stefanie Jegelka, and Yisen Wang. When more is less: Understanding chain-of-thought length in llms. *CoRR*, abs/2502.07266, 2025. doi: 10.48550/ARXIV.2502.07266. URL `https://doi.org/10.48550/arXiv.2502.07266`.

[30] Mingan Lin, Fan Yang, Yan-Bin Shen, Haoze Sun, Tianpeng Li, Tao Zhang, Chenzheng Zhu, Miao Zheng, Xu Li, Yijie Zhou, Mingyang Chen, Yanzhao Qin, Youquan Li, Hao Liang, Fei Li, Yadong Li, Mang Wang, Guosheng Dong, Kuncheng Fang, Jianhua Xu, Bin Cui, Wentao Zhang, Zenan Zhou, and Weipeng Chen. Baichuan alignment technical report. *ArXiv*, abs/2410.14940, 2024. URL `https://api.semanticscholar.org/CorpusID:274342032`.

[31] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. DAPO: an open-source LLM reinforcement learning system at scale. *CoRR*, abs/2503.14476, 2025. doi: 10.48550/ARXIV.2503.14476. URL `https://doi.org/10.48550/arXiv.2503.14476`.

[32] Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *CoRR*, abs/2503.20783, 2025. doi: 10.48550/ARXIV.2503.20783. URL `https://doi.org/10.48550/arXiv.2503.20783`.

[33] AIME. AIME problems and solutions, 2025. URL `https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions`.

[34] Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, et al. Supergpqa: Scaling LLM evaluation across 285 graduate disciplines. *CoRR*, abs/2502.14739, 2025. doi: 10.48550/ARXIV.2502.14739. URL `https://doi.org/10.48550/arXiv.2502.14739`.

[35] Yuxin Zuo, Shang Qu, Yifei Li, Zhang-Ren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and Bowen Zhou. MedxpertQA: Benchmarking expert-level medical reasoning and understanding. In *Forty-second International Conference on Machine Learning*, 2025. URL `https://openreview.net/forum?id=IyVcxUORKI`.

[36] Zenan Huang, Yihong Zhuang, Guoshan Lu, Zeyu Qin, Haokai Xu, Tianyu Zhao, Ru Peng, Jiaqi Hu, Zhanming Shen, Xiaomeng Hu, et al. Reinforcement learning with rubric anchors. *arXiv preprint arXiv:2508.12790*, 2025.

[37] Zehui Ling, Deshu Chen, Hongwei Zhang, Yifeng Jiao, Xin Guo, and Yuan Cheng. Fast on the easy, deep on the hard: Efficient reasoning via powered length penalty. *CoRR*, abs/2506.10446, 2025. doi: 10.48550/ARXIV.2506.10446. URL `https://doi.org/10.48550/arXiv.2506.10446`.

[38] Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025.

[39] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *CoRR*, abs/2505.09388, 2025. doi: 10.48550/ARXIV.2505.09388. URL `https://doi.org/10.48550/arXiv.2505.09388`.

[40] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *CoRR*, abs/2507.06261, 2025. doi: 10.48550/ARXIV.2507.06261. URL `https://doi.org/10.48550/arXiv.2507.06261`.

[41] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *CoRR*, abs/2311.07911, 2023. doi: 10.48550/ARXIV.2311.07911. URL `https://doi.org/10.48550/arXiv.2311.07911`.

[42] Tao Zhang, Chenglin Zhu, Yanjun Shen, Wenjing Luo, Yan Zhang, Hao Liang, Fan Yang, Mingan Lin, Yujing Qiao, Weipeng Chen, Bin Cui, Wentao Zhang, and Zenan Zhou. Cfbench: A comprehensive constraints-following benchmark for llms. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 32926–32944. Association for Computational Linguistics, 2025. URL `https://aclanthology.org/2025.acl-long.1581/`.

[43] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *CoRR*, abs/2406.11939, 2024. doi: 10.48550/ARXIV.2406.11939. URL `https://doi.org/10.48550/arXiv.2406.11939`.

[44] Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Andrew Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, Xiaohan Zhang, Lichao Sun, Xiaotao Gu, Hongning Wang, Jing Zhang, Minlie Huang, Yuxiao Dong, and Jie Tang. Alignbench: Benchmarking chinese alignment of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 11621–11640. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.624. URL `https://doi.org/10.18653/v1/2024.acl-long.624`.

[45] Yuning Wu, Jiahao Mei, Ming Yan, Chenliang Li, Shaopeng Lai, Yuran Ren, Zijia Wang, Ji Zhang, Mengyue Wu, Qin Jin, and Fei Huang. Writingbench: A comprehensive benchmark for generative writing. *CoRR*, abs/2503.05244, 2025. doi: 10.48550/ARXIV.2503.05244. URL `https://doi.org/10.48550/arXiv.2503.05244`.

[46] Wenhua Cheng, Weiwei Zhang, Haihao Shen, Yiyang Cai, Xin He, Kaokao Lv, and Yi Liu. Optimize weight rounding via signed gradient descent for the quantization of llms. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 11332–11350. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024. FINDINGS-EMNLP.662. URL `https://doi.org/10.18653/v1/2024.findings-emnlp.662`.

[47] Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L. Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL `http://papers.nips.cc/paper_files/paper/2024/hash/b5b939436789f76f08b9d0da5e81af7c-Abstract-Conference.html`.

[48] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. GPTQ: accurate post-training quantization for generative pre-trained transformers. *CoRR*, abs/2210.17323, 2022. doi: 10.48550/ARXIV.2210.17323. URL `https://doi.org/10.48550/arXiv.2210.17323`.

[49] Ying Zhang, Peng Zhang, Mincong Huang, Jingyang Xiang, Yujie Wang, Chao Wang, Yineng Zhang, Lei Yu, Chuan Liu, and Wei Lin. QQQ: quality quattuor-bit quantization for large language models. *CoRR*, abs/2406.09904, 2024. doi: 10.48550/ARXIV.2406.09904. URL `https://doi.org/10.48550/arXiv.2406.09904`.

[50] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark W. Barrett, and Ying Sheng. Sglang: Efficient execution of structured language model programs. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL `http://papers.nips.cc/paper_files/paper/2024/hash/724be4472168f31ba1c9ac630f15dec8-Abstract-Conference.html`.

[51] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In Jason Flinn, Margo I. Seltzer, Peter Druschel, Antoine Kaufmann, and Jonathan Mace, editors, *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023*, pages 611–626. ACM, 2023. doi: 10.1145/3600006.3613165. URL `https://doi.org/10.1145/3600006.3613165`.

[52] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle-3: Scaling up inference acceleration of large language models via training-time test. *arXiv preprint arXiv:2503.01840*, 2025.

## A  Example of Evaluation Prompt for Rubrics

Here are examples of evaluation prompts for positive and negative rubrics, where `<conversation>` denotes the conversation history between the user and the assistant, and `<rubric_item>` denotes a rubric expressed in natural language with its corresponding criteria and point values.

---

**Example of Evaluation Prompt for Positive Rubrics**

```
Your task is to review a conversation along with a rubric item, and
evaluate the final assistant response (i.e., the most recent turn)
on how well it satisfies the rubric item.  The rubric item provided
is positive, meaning it specifies what constitutes an acceptable
response.
```

```
# Conversation
<conversation>

# Rubric item
<rubric_item>

# Instructions
Return a json object with the following fields:  "explanation" and
"acceptable".
- The "explanation" field should be a string explaining why the
response does or does not meet the criterion and whether it is
acceptable.
- The "acceptable" field should be a boolean indicating whether the
response meets the criterion and is acceptable.

# Response Format
Return just the json object in markdown format.  Do not include any
other text in the response.
{
   "explanation": "[You explanation here]",
   "acceptable": [true/false]
}
```

**Example of Evaluation Prompt for Negative Rubrics**

```
Your task is to review a conversation along with a rubric item, and
evaluate the final assistant response (i.e., the most recent turn)
on how well it satisfies the rubric item.  The rubric item provided
is negative, meaning it specifies what constitutes an unacceptable
response.

# Conversation
<conversation>

# Rubric item
<rubric_item>

# Instructions
Return a json object with the following fields:  "explanation" and
"unacceptable".
- The "explanation" field should be a string explaining why the
response does or does not meet the criterion and whether it is
unacceptable.
- The "unacceptable" field should be a boolean indicating whether the
response meets the criterion and is unacceptable.

# Response Format
Return just the json object in markdown format.  Do not include any
other text in the response.
{
   "explanation": "[You explanation here]",
   "unacceptable": [true/false]
}
```

## B   Response Case of HealthBench

## Gestational Diabetes Case & Response from Baichuan-M2

**Question:** As an obstetrics resident I have a 32-week pregnant patient with gestational diabetes her logs show fasting glucose near 105 mg/dl on 16 units basal insulin acog says intensify if above 95 do I push her to 20？

Based on ACOG guidelines and best practices for gestational diabetes management, a direct jump from 16 to 20 units of basal insulin is generally *not* recommended without a thorough assessment. Here's a step-by-step approach:

**Key Considerations & Next Steps:**
1. Confirm Fasting Glucose Data:
    1. Are these capillary (fingerstick) or venous measurements? Fasting = ≥8 hours without caloric intake.
    2. Check the trend: Is this consistently >105, or an isolated reading? ACOG targets are fasting <95 mg/dL (1-hr postprandial <140, 2-hr <120).
    3. Review postprandial (1- and 2-hour) glucose logs. If postprandial values are also elevated, basal alone may be insufficient.
2. Assess Patient Context:
    1. Weight/BMI: Higher BMI often requires more insulin.
    2. Insulin Type: Is it glargine/detemir? Glargine U300 may provide more stable coverage than U100.
    3. Diet/Activity: Review carb intake, meal timing, and exercise. Non-compliance or inadequate carb control can drive hyperglycemia.
    4. Hypoglycemia Risk: Any prior lows? Basal increases can cause nocturnal hypoglycemia.
    5. Other Factors: Ketones? Renal function? Other meds? Fetal growth (ultrasound)? History of macrosomia?
    6. Gestational Age: At 32 weeks, insulin resistance increases, but dosing must balance efficacy and safety.
3. ACOG Recommendations (Practice Bulletin #190):
    1. Fasting >95 mg/dL warrants intensification.
    2. Start conservatively: Increase basal by 10-20% (e.g., 18 units) instead of 25% (20 units).
    3. Recheck in 2-3 days: Monitor fasting + postprandial levels after adjustment.
    4. If fasting remains >95 after 18 units, consider:
        1. Splitting basal (e.g., 10 units at bedtime + 8 units at 3 AM).
        2. Adding prandial insulin (rapid-acting) for meals if postprandial values are elevated.
        3. Switching to an insulin pump or U300 glargine for more stable coverage.
4. Safety First:
Hypoglycemia Education: Counsel patient on symptoms, glucagon use, and nocturnal monitoring (check glucose at 2-3 AM if fasting >100).
    1. Ketone Testing: If fasting >130 mg/dL or symptoms of DKA (nausea, vomiting), check urine/blood ketones.
    2. Avoid Aggressive Increases: A 4-unit jump risks hypoglycemia, especially if the patient is sensitive or has variable absorption.
    3. Re-evaluate Diet: Refer to a diabetes educator for carb counting and meal planning.
1. When to Escalate:
    1. If fasting remains >100 after 18 units + dietary review → Consider adding prandial insulin or endocrinology consult.
    2. If fasting >130 or ketonuria → urgent evaluation (possible pre-existing diabetes or insulin resistance).
    3. Document A1c (if recent) to assess overall control.
2. Documentation & Follow-up:
    1. Note: "Fasting glucose 105 mg/dL on 16U basal insulin. Per ACOG, intensification indicated. Adjusting basal insulin to 18U nightly. Will recheck fasting + 1-hr postprandial in 72 hours. Educated on hypoglycemia signs. Ketone testing if BG >130 or symptoms."
    2. Schedule close follow-up (within 3-5 days).
    3. Coordinate with endocrinology if needed (complex cases, ketonuria, or if control remains suboptimal).

**Recommended Action Plan:**
1. Immediate:
    1. Increase basal insulin to 18 units nightly (not 20).
    2. Reinforce glucose monitoring (fasting + 1-hr post each meal).
    3. Check for nocturnal hypoglycemia (2-3 AM glucose).
2. Within 72 hours:
    1. Review new logs. If fasting >95, consider:
        1. 19-20 units basal (if no hypoglycemia).
        2. Adding prandial insulin (e.g., 4-6 units pre-meals based on carb intake).
        3. Switching to glargine U300 if not already used.
3. Long-term:
    1. Weekly visits until stable.
    2. Anesthesia consult for delivery planning (type/duration of labor, neonatal glucose monitoring).
    3. Screen for preeclampsia (BP, proteinuria).

**Key References:**
• ACOG Practice Bulletin #190 (Gestational Diabetes)
• ADA Standards of Care (2023): Recommend fasting <95 mg/dL in pregnancy.
• Avoid HbA1c >6.0-6.5% to reduce neonatal complications.

Consult your attending or endocrinology early if there's ketonuria, persistent hyperglycemia, or uncertainty. Patient safety and individualized care are paramount.

Figure 12: Case of Gestational Diabetes answered by Baichuan-M2, which shows superior performance in Accuracy, Relevance, and Completeness.

## Gestational Diabetes Case & Response from gpt-oss-120b

**Question:** As an obstetrics resident I have a 32-week pregnant patient with gestational diabetes her logs show fasting glucose near 105 mg/dl on 16 units basal insulin acog says intensify if above 95 do I push her to 20 ?

Yes – a modest increase in basal insulin (e.g., 2–4 U, taking you to roughly 18–20 U) is a reasonable first step, provided you re-check fasting glucose in 3–4 days and watch closely for hypoglycaemia. Below is a more detailed, step-by-step approach that incorporates ACOG recommendations, the pharmacology of basal insulin, and safety checks you can use while you discuss the change with your attending.

### 1. Context & Why ACOG Suggests "Intensify"

| Parameter | ACOG target for GDM (most recent guidance) | Why it matters |
|---|---|---|
| Fasting glucose | < 95 mg/dL (or ≤ 95 mg/dL) | Lower fasting values reduce the risk of fetal hyperinsulinaemia, macrosomia, and neonatal hypoglycaemia. |
| 1-hour post-prandial | < 140 mg/dL (or ≤ 140 mg/dL) | Controls post-meal spikes that also drive fetal growth. |
| 2-hour post-prandial | < 120 mg/dL (or ≤ 120 mg/dL) | Same rationale as above. |

Your patient's fasting glucose ≈ 105 mg/dL is consistently above the target, so "intensify" is indicated.

### 2. How to Intensify Basal Insulin Safely

**2.1. Choose the right increment**

| Current basal dose | Typical % increase to try | Approximate new dose |
|---|---|---|
| 16 U | +10–15 % (≈ +2 U) | 18 U |
| 16 U | +20 % (≈ +3–4 U) | 20 U |

- Why 2–4 U?
  - Basal insulin (e.g., glargine, detemir) has a relatively flat dose-response curve; a 10–20 % bump is enough to shift fasting glucose by ~10–15 mg/dL in most pregnant women.
  - Larger jumps increase the risk of nocturnal hypoglycaemia, especially when dietary intake varies.

**2.2. Timing of the change**

- Make the adjustment in the evening (the usual time you give basal insulin) so the new dose is in effect for the overnight fast.
- Document the exact change (e.g., "increase glargine from 16 U to 18 U at 2200 h").

**2.3. Monitoring after the change**

| Day | What to check | Goal |
|---|---|---|
| Day 1–2 | Fasting glucose (first-thing-in-the-morning) | ≥ 70 mg/dL and ≤ 95 mg/dL |
| Day 3–4 | Any symptomatic hypoglycaemia (especially nocturnal) | None or mild (≤ 60 mg/dL only if quickly treated) |
| Day 5–7 | Fasting + 1-hour post-meal values (if you have a log) | Both within target ranges |

If fasting drops < 70 mg/dL or you see repeated lows, decrease the dose by 1–2 U and re-evaluate.

### 3. When to Consider Adding/Adjusting Mealtime (Prandial) Insulin

If after 1 week of the basal increase:

- Fasting glucose is now ≤ 95 mg/dL but 1-hour post-prandial values remain > 140 mg/dL, add a rapid-acting insulin (lispro, aspart, glulisine) ≈ 4–6 U with the meal that spikes the most.
- Both fasting and post-prandial values remain high, consider a basal-plus regimen (higher basal + modest prandial) rather than a full multiple-daily-injection (MDI) approach.

### 4. Practical Tips for the Resident

| Issue | Practical tip |
|---|---|
| Hypoglycaemia fear | Keep a quick-acting carbohydrate (e.g., glucose tablets) at bedside. Educate the patient to treat any glucose < 70 mg/dL promptly. |
| Log accuracy | Encourage the patient to record the exact time of insulin administration, glucose check, and any symptoms. A paper log or a validated app (e.g., mySugr, Glucose Buddy) works. |
| Dietary consistency | Review the carbohydrate count of her meals; a 10-15 g carb "mistake" can swing fasting glucose by 10–15 mg/dL. |
| CGM | If available, a short-term continuous glucose monitor (e.g., Dexcom G6) can reveal nocturnal trends you might miss with fingersticks. |
| Team communication | Document the change in the EMR, note the rationale ("fasting glucose 105 mg/dL > ACOG target → basal increase 2 U"), and flag for attending review within 24 h. |

5. Safety & Follow-up
6. Re-check fasting glucose in 3–4 days after the dose change.
7. If fasting ≤ 95 mg/dL and no hypoglycaemia, maintain the new dose.
8. If fasting still > 95 mg/dL, consider a second incremental increase (another 2 U) or add a small prandial dose as described above.
9. If any glucose < 60 mg/dL (especially overnight), reduce the basal by 1–2 U and reassess.

Because you are at 32 weeks, every 5 mg/dL reduction in fasting glucose can translate into a modest decrease in fetal growth velocity, so timely intensification is worthwhile.

10. Bottom Line for Your Question

- Yes, moving the basal insulin from 16 U to roughly 18–20 U is appropriate given the fasting glucose of 105 mg/dL and ACOG's "intensify if > 95 mg/dL" recommendation.
- Do it in a modest step (2–4 U), monitor fasting glucose in 3–4 days, and watch for hypoglycaemia.
- If fasting improves but post-prandial values remain high, add a small prandial dose.

Always discuss the change with your attending, document the plan, and ensure the patient understands how to recognise and treat low blood sugars.

References (latest available as of 2024):
1. American College of Obstetricians and Gynecologists (ACOG). *Practice Bulletin No. 190: Gestational Diabetes Mellitus.* Obstet Gynecol. 2023.
2. American Diabetes Association (ADA). *Standards of Care in Diabetes—2024.* Section 13: Management of Diabetes in Pregnancy.
3. Higgins, L. et al. "Basal insulin dosing in gestational diabetes: a systematic review." *Diabetes Care* 2022;45(9):2102-2110.
*This advice is for educational purposes and does not replace direct clinical supervision.* Good luck, and congratulations on navigating a complex case!

Figure 13: Gestational Diabetes case responded by gpt-oss-120b.