

Evolving Diagnostic Agents in a Virtual Clinical Environment

Pengcheng Qiu^{1,2,*}, Chaoyi Wu^{1,2,*}, Junwei Liu^{3,4*}, Qiaoyu Zheng^{1,2}, Yusheng Liao^{1,2},
Haowen Wang³, Yun Yue³, Qianrui Fan³, Shuai Zhen³, Jian Wang³, Jinjie Gu³,
Yanfeng Wang^{1,2}, Ya Zhang^{1,2,†} and Weidi Xie^{1,2,†}

¹Shanghai Jiao Tong University, Shanghai, China

²Shanghai Artificial Intelligence Laboratory, Shanghai, China

³Intelligence Healthcare Department, AntGroup, Hangzhou, China

⁴Intelligence Computing and Sensing Laboratory, Peking University, Beijing, China

*Equal contributions

†Corresponding author

Ya Zhang: ya_zhang@sjtu.edu.cn; Weidi Xie: weidi@sjtu.edu.cn

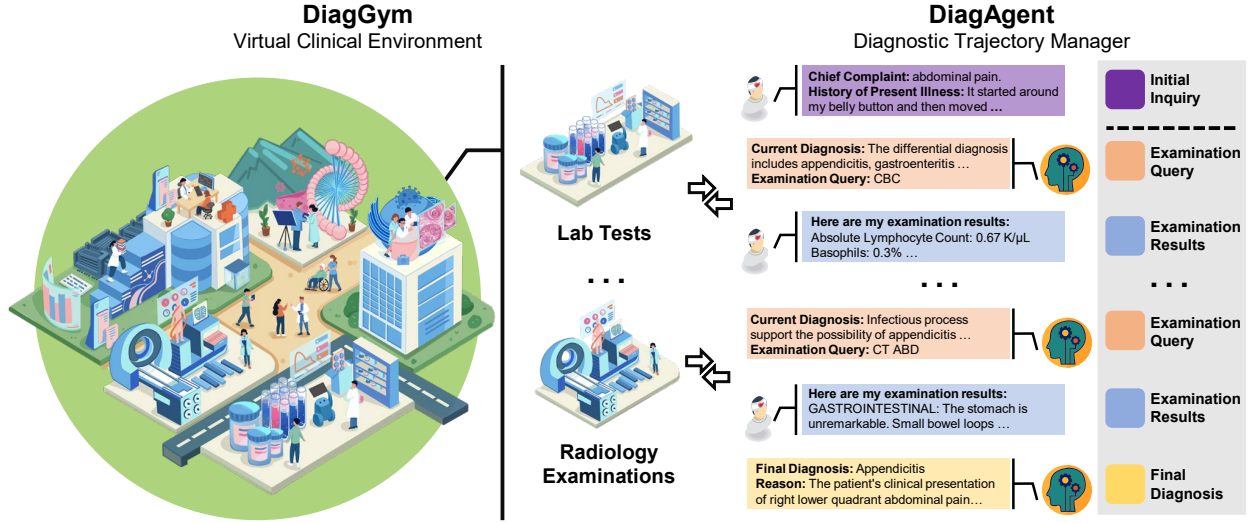
In this paper, we present a framework for training large language models (LLMs) as diagnostic agents with reinforcement learning, enabling them to manage multi-turn diagnostic processes, adaptively select examinations, and commit to final diagnoses. Unlike instruction-tuned models trained on static case summaries, our method acquires diagnostic strategies through interactive exploration and outcome-based feedback, mapping evolving patient states to the next optimal examination and subsequent diagnosis. Our contributions are threefold: (i) we present a diagnostics world model trained with electronic health records (EHRs), termed as **DiagGym**, which enables to emit examination outcomes conditioned on patient history and recommended examination, serving as a virtual clinical environment to support realistic, closed-loop in-silico diagnosis training and evaluation; (ii) we train a diagnostic agent, **DiagAgent**, via end-to-end, multi-turn reinforcement learning within the environment, to learn diagnostic policies that optimizes both information yield and diagnostic accuracy; (iii) We introduce a new diagnostic benchmark, **DiagBench**, designed to evaluate multi-turn diagnostic interaction trajectories. The benchmark comprises 750 cases with physician-validated examination recommendations leading to final diagnoses, as well as 99 cases annotated with 973 physician-written rubrics on diagnosis process. (iv) we demonstrate superior performance across diverse diagnostic settings. **DiagAgent** significantly outperforms 10 state-of-the-art large language models (LLMs), including DeepSeek-v3 and GPT-4o, as well as two recently developed prompt-engineered agents. **DiagAgent** achieves a 9.34% higher diagnostic accuracy and a 44.03% improvement in examination recommendation hit ratio in the single-turn setting. In a more practical end-to-end setting, it delivers a 15.12% increase in diagnostic accuracy and a 23.09% boost in examination recommendation F1 score. Furthermore, in a rubric-based evaluation, it surpasses the next-best model, Claude-sonnet-4, by 7.1% in weighted rubric score, highlighting its advancement in managing multi-turn diagnostic trajectories. These findings indicate that learning policies in interactive clinical environments confers dynamic and clinically meaningful long-term diagnostic management abilities that are unattainable through passive training alone.

1 Introduction

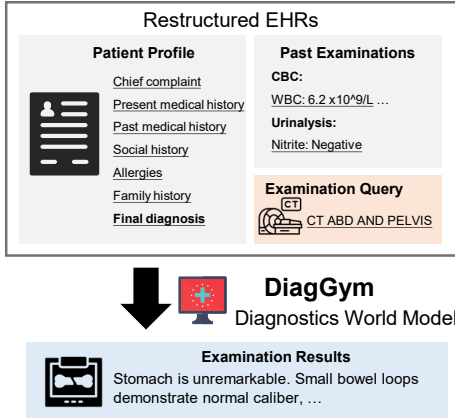
Large language models (LLMs) have made tremendous progress in advancing medical AI, achieving strong performance on rigorous benchmarks—including USMLE-style examinations—and across diverse clinical tasks [1, 2, 3, 4, 5]. Recent advances in reasoning [6, 7] and post-training methods such as supervised fine-tuning (SFT) [8, 9, 10] have spurred interest in diagnostic applications, where the challenge is not simply answering a question but managing a complex, evolving patient case.

Clinical diagnosis, however, is not a static prediction problem. It is inherently an iterative long-term decision-making process under uncertainty: clinicians must synthesize partial information, decide which examination to recommend—or whether to commit to a diagnosis—and balance informativeness, timeliness, cost, and safety. Yet, existing LLMs are predominantly trained on passively collected, instruction-style corpora that assume a complete, fixed patient records [1, 4, 8, 11, 12, 13]. This static paradigm collapses the multi-turn nature of real diagnosis into a single shot, eliminating the interaction with an external environment or revise

a. Method Overview



b. DiagGym Pipeline



c. DiagAgent Training

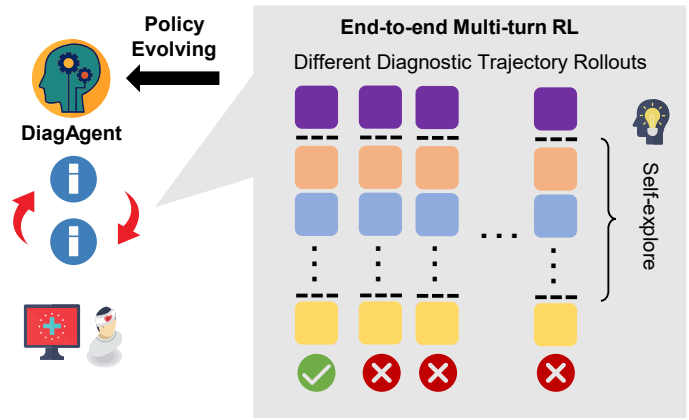


Figure 1 | Overview of our method. **a** illustrates the overview of our method, we establish a virtual clinical environment, DiagGym, that can simulate examination results in real time. Then, within it, we train a diagnostic agent capable of managing multi-turn diagnostic trajectories in a long-term manner, recommending, recommending diverse examinations until sufficient evidence is gathered for a final diagnosis. **b** presents the diagnostics world model we constructed based on EHRs, representing a virtual clinical environment. It receives the patient’s basic profile and performed past examinations, and the next examination query as condition input, then simulates the related results as feedback. **c** depicts the DiagAgent end-to-end multi-turn RL training, where the agent interacts with the virtual environment, rolls out different possible diagnostic trajectories, self-explores suitable examination recommendation chains, and iteratively evolves its decision-making policy through end-to-end reinforcement rewards.

hypotheses as evidence accumulates. As a result, state-of-the-art models often fail to plan full diagnostic trajectories [10, 14, 15, 16, 17], including which tests to recommend, when to stop, and when to final diagnose.

In this paper, we introduce a framework for training large language models as diagnostic agents in a virtual clinical environment with reinforcement learning (RL) (Figure 1a), enabling them to learn policies that actively manage multi-turn diagnostic trajectories, adaptively suggest examinations, and commit to final diagnoses.

Central to our agent-training framework is a diagnostics world model, **DiagGym**, built on the EHRs with generative model (Figure 1b). **DiagGym** is able to emit examination results conditioned on the evolving patient state—including background patient profile, past examinations, and the next examination query. This interactive, outcome-based feedback enables the safe, closed-loop simulation of real diagnostic workflows, allowing the diagnostic agents to order tests and immediately observe their consequences. By coupling

conditional generation with longitudinal patient context, **DiagGym** captures long-horizon dependencies absent from static, instruction-style corpora, creating a realistic in-silico virtual clinical environment for training and evaluating diagnostic agents with reinforcement learning.

Within this environment, we train a diagnostic agent, **DiagAgent** (Figure 1c) through end-to-end, multi-turn reinforcement learning. Here, the agent’s policy network maps the current patient state to the next optimal action: either recommending the most informative examination or finalizing the diagnosis. Reward signals are derived from the informativeness of queried examinations, diagnostic accuracy, and efficiency, while penalties are applied for redundant or unnecessary steps, determining which diagnostic rollouts are most suitable. This approach potentially enables the agent to explore diverse patient trajectories, refine policies and discover diagnostic management strategies beyond human expertise, without relying on risky, time-intensive real-world implementation. By aligning diagnostic agents with long-term trajectory-level reasoning, they learn when to continue investigation, which examinations will yield the highest clinical value, and when to stop to achieve an accurate and timely diagnosis. This extends current diagnostic agents from isolated, point-in-time consultations to comprehensive long-term patient management, surpassing limits of statistic supervised training.

To comprehensively assess the generated multi-turn interactive diagnostic trajectory, we constructed a new benchmark, **DiagBench**. It contains 750 physician-validated diagnosis cases. Each case is equipped with not only the final diagnostic result but also a referenced diagnostic trajectory with key recommending examinations towards final results. Furthermore, to enable more fine-grained qualitative evaluation, 99 cases in the benchmark are annotated with detailed diagnostic rubrics that specify which reasoning rules should be appropriately followed within the diagnostic trajectory. In total, 973 physician-written rubrics are provided, each assigned weighting points to precisely assess performance at key steps of the diagnostic process.

In experiments, we first assess the reliability of **DiagGym** as a diagnostics world model for examination result generation. **DiagGym** achieves 96.91% instance-wise clinical consistency and an examination-wise Wasserstein distance of 0.128 to the real-world distribution, significantly outperforming a DeepSeek-v3-based simulator (88.81% and 1.336). Moreover, its normalized variance on the examination-wise distribution closely matches the ground-truth distribution (3.46 against 5.31), indicating its broader coverage and reduced mode collapse [18]. Together, these results underscore its high fidelity and diversity, serving as a diagnostic agent training ‘gym’.

Then, we evaluate the final **DiagAgent** in two complementary settings. *First*, in single-turn evaluations on real cases, it recommends an additional examination or renders a final diagnosis based on observed results. *Second*, in end-to-end evaluations on simulated cases, **DiagAgent** interacts with **DiagGym** to complete full diagnostic workflows. Across both settings, we assess examination-recommendation quality and diagnostic accuracy against ground truth, comparing against 10 state-of-the-art LLMs (*e.g.*, DeepSeek-v3, MedGemma) as well as 2 more recent prompt-engineered agentic systems (*e.g.*, MDAgents). **DiagAgent** significantly exceeds the next-best method by 9.34% higher accuracy and a 44.03% higher hit ratio for examination recommendation in single-turn evaluation; 15.12% in diagnostic accuracy and 23.09 in F1 score in end-to-end evaluation. Furthermore, we perform a rubric-based, end-to-end evaluation in which the generated interactive diagnostic trajectory is qualitatively assessed against detailed, physician-authored rubrics. This process examines whether each criterion is met and calculates the total scores accordingly. Here, **DiagAgent** surpassed the next-best model, Claude-sonnet-4, by 7.1% in weighted rubric score. These gains highlight the benefits of jointly optimizing diagnostic agents and simulation environments within a closed-loop learning framework.

Together, these results show that learning a diagnostic policy within an interactive clinical environment equips LLMs with dynamic, clinically meaningful diagnostic capabilities that are unattainable with passive training alone, and establishes **DiagGym** as a scalable in-silico testbed for optimizing diagnostic management strategies prior to prospective clinical validation.

2 Problem Formulation

We first formalize the functionality of **DiagGym** and its role in training **DiagAgent**.

DiagGym. As illustrated in Figure 1b, we define a diagnostics world model as a conditional textual EHR

generator, Φ_{env} , that generates synthetic examination results conditioned on a dynamically evolving patient state. At step t , the patient EHR state is (\mathcal{B}, E_t) , where \mathcal{B} is the background patient profile—including chief complaint, present medical history, and the final diagnosis. The set $E_t = \{(a_1, e_1), (a_2, e_2), \dots, (a_t, e_t)\}$ represents the past examination records, where each a_i denotes a specific examination item, *e.g.*, “complete blood count” or “CT abdomen examination,” and e_i is the result.

DiagGym is designed to generate the potential examination result for the patient based on a specific examination query a_{t+1} , as follows:

$$e_{t+1} = \Phi_{\text{env}}(a_{t+1} \mid E_t, \mathcal{B}), \quad (1)$$

where initially $E_0 = \emptyset$, and e_{t+1} denotes the synthetic examination result.

We frame the training process as a conditional generation task, that minimizes the negative log-likelihood of the ground-truth examination results. These results are treated as free text, regardless of whether they are numerical (with numbers directly embedded as text) or textual. The training objective is formalized as:

$$\mathcal{L}_{\text{sim}} = - \sum_{t=0}^{T-1} \log \Phi_{\text{env}}(\hat{e}_{t+1} \mid a_{t+1}, E_t, \mathcal{B}), \quad (2)$$

where \hat{e}_t represents the ground truth examination result at step t and T denotes the total examination length recorded in a certain EHR. More details can be found in Method Section 5.1.

Once trained, Φ_{env} can generate plausible results for any examination and patient state, capturing conditional dependencies across diseases, histories, and prior tests. This capability enables safe and repeatable reinforcement learning (RL) training of diagnostic agents without direct access to real patient records; in other words, it serves as a virtual clinical environment for RL.

DiagAgent. As illustrated in Figure 1c, within DiagGym we train a diagnostic agent, **DiagAgent**, using reinforcement learning. Formally, at time step t , the agent’s state is defined as $s_t = (\mathcal{I}, E_t)$, where \mathcal{I} is the patient’s initial inquiry—including chief complaint, history of present illness, and other relevant presentation details—but, unlike \mathcal{B} in the environment model, contains no information about the final diagnosis. The set $E_t = \{(a_1, e_1), \dots, (a_t, e_t)\}$ records the examinations performed so far and their observed results.

For the agent, its action space is defined as $\mathcal{A} = \{a_1, a_2, \dots, a_N\}$, representing all available clinical examination items and the final diagnosis action. In response to the current state s_t , the agent selects an action to recommend the next examination for the patient, based on its policy:

$$a_{t+1} \sim \pi_{\theta}(a \mid s_t), \quad (3)$$

where π_{θ} is the learnable policy function parameterized by a large language model Φ_{diag} :

$$\pi_{\theta} = \Phi_{\text{diag}}(s_t). \quad (4)$$

DiagGym then returns reasonable examination results, serving as the external environment feedback:

$$e_{t+1} = \Phi_{\text{env}}(a_{t+1} \mid E_t, \mathcal{B}), \quad (5)$$

$$s_{t+1} = s_t \cup (a_{t+1}, e_{t+1}), \quad (6)$$

where s_{t+1} is the next state.

The diagnostic trajectory proceeds until the agent selects a final diagnosis action, after which DiagAgent outputs the predicted diagnosis D . The ultimate training objective is to optimize the interactive policy function Φ_{diag} to maximize the expected cumulative reward:

$$\max_{\Phi_{\text{diag}}} \mathbb{E}_{\Phi_{\text{diag}}} \left[\sum_{t=1}^T \gamma^t \mathcal{R}(s_t, a_t) \right], \quad (7)$$

where $\gamma \in [0, 1]$ is the discount factor, T is the trajectory length, and \mathcal{R} is the reward function, defined as the sum of three sub-rewards:

$$\mathcal{R} = \lambda_1 r_{\text{diag}} + \lambda_2 r_{\text{exam}} + \lambda_3 r_{\text{turn}}, \quad (8)$$

with $\lambda_1, \lambda_2, \lambda_3$ as hyperparameters. Here, r_{diag} encourages accurate diagnoses, r_{exam} promotes relevant examination recommendations, and r_{turn} rewards fewer used turns. The detailed design of the reward function is provided in the Method Section 5.2.

The final trained DiagAgent can actively manage multi-turn diagnostic trajectories by iteratively interacting with patients, selecting relevant examinations, and ultimately arriving at an accurate final diagnosis.

3 Results

In this section, we first evaluate the two key components: **DiagGym**, the high-fidelity diagnostics world model, and **DiagAgent**, the reinforcement-trained diagnostic agent. Then, we carry out ablation studies to investigate the effectiveness of our approach design. Lastly, we present detailed case studies.

3.1 Evaluation for DiagGym

We first assess the fidelity and reliability of DiagGym, as it forms the foundation for RL-based training of our diagnostic agent. The evaluation aims to verify whether the world model can generate clinically consistent, context-appropriate examination results that faithfully reflect real-world patterns in EHRs.

3.1.1 Evaluation Settings

We construct an evaluation set of 863 patient cases from MIMIC-IV using the process in Section 5.1. These cases span 863 distinct diseases, categorized based on the original ICD codes in MIMIC-IV corresponding to patient admissions, representing a combined 35,548 examination records: on average, 8.77 physical exams, 28.37 laboratory events, 2.04 microbiology events, and 2.01 radiology events per case.

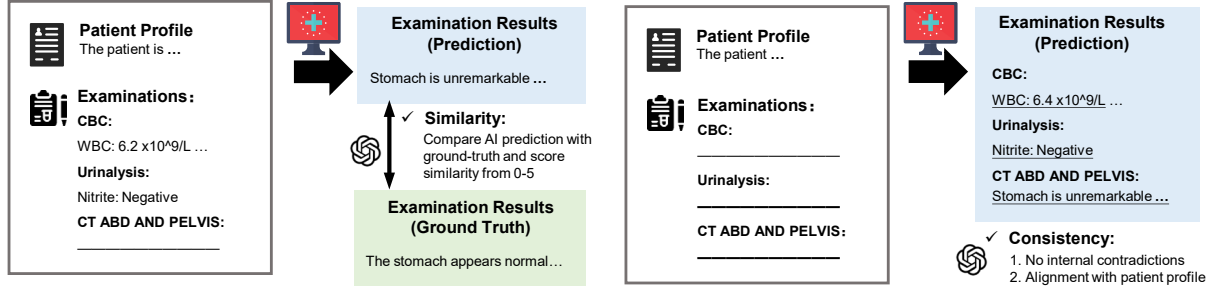
Each case comprises two components: (i) **patient profile**: baseline information including chief complaint, present and past medical history, social history, allergies, family history, and the final diagnosis; (ii) **examination chain**: the chronological sequence of actual examination results, serving as ground truth. The simulator’s task is to reconstruct each examination result in sequence, conditioned on the patient profile, all prior examination data, and the current examination query. This sequential generation ensures that, during RL training, the simulator can produce arbitrarily long, detail-rich examination chains in response to multi-turn queries.

We quantify the generation quality using instance-wise and examination-wise metrics (Figure 2), the detailed metric calculations are provided in Supplementary 10.1.1. These metrics jointly assess whether simulated results match real diagnostic examination results at both the instance-wise and the examination-wise.

Instance-wise metrics assess the quality of generated examination sequences at the level of individual patient cases. This evaluation utilizes two methods: an **automated rater** (GPT-4o, version gpt-4o-2024-08-06) or **physician rater**. Straightforwardly, the two settings are similar, except that the former uses LLMs for scoring, whereas the latter relies on human evaluators. In the human evaluation, due to cost, a random sample of 100 instances are adopted and three independent physicians are involved, rating them based on the clinical plausibility of the generated content rather than formats. Two related metrics are used:

- **Step-level similarity** measures how closely the simulator’s output for each examination step matches the corresponding real-world record, given the patient profile and all prior ground-truth results. Similarity is scored on a 0-5 scale, where 5 indicates perfect medical equivalence or high similarity to the reference. Both the automated evaluator and physician raters assign an independent score from 0 (no similarity) to 5 (perfect equivalence).
- **Full-chain consistency** evaluates the coherence of an entire generated sequence where each step depends on the simulator’s **previously generated outputs**. This setting mirrors the use of reinforcement learning, prioritising internal clinical consistency over word-for-word agreement with ground truth. Consistency is judged using a binary score (1 = consistent, 0 = inconsistent). The automated evaluator

a. Instance-wise Metrics



b. Examination-wise Metrics

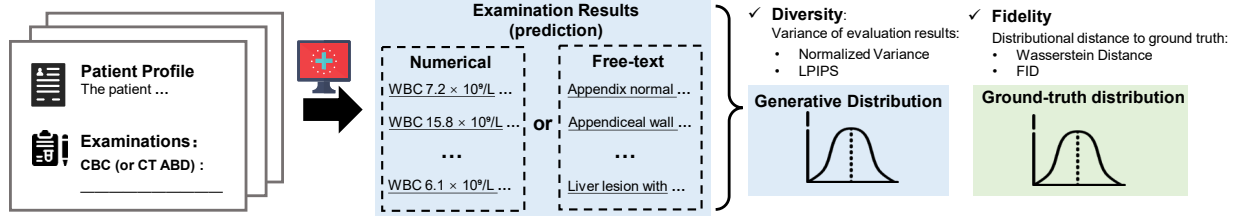


Figure 2 | Overview of simulator evaluation settings. **a** Instance-wise metrics: GPT-4o assesses the quality of generated examination results on an individual patient case level. **b** Examination-wise metrics: fidelity and diversity are evaluated by comparing the statistical distributions of generated examination results against those from real cases.

provides a binary judgment using [prompt 15](#). In the physician rating, raters must make a judgment ensuring the sequence maintains adherence to medical common sense, features appropriate calibration of severity, and is free of internal contradictions or conflicts across all reported findings.

Examination-wise metrics probe the realism and variability of individual examination items, regardless of their place in the sequence. We distinguish between numerical and free-text outputs:

- **Numerical fidelity & diversity** compares generated numerical values (*e.g.*, red blood cell counts) to real distributions. Fidelity is quantified via the 1-Wasserstein distance (lower is better), while diversity is measured as the normalized variance of the generated distribution (higher reflects broader coverage and less mode collapse [18]). Metrics are averaged across all selected numerical tests (Supplementary 10.3).
- **Free-text fidelity & diversity:** compares the generated narrative reports (*e.g.*, CT abdomen findings) in embedding space using BioLORD [19]. Fidelity is measured by the Fréchet Inception Distance [20] (lower is better) and diversity by Intra-LPIPS [21] (higher indicates more inter-case variation). Metrics are averaged over the selected free-text examinations (Supplementary 10.3).

During the calculation of examination-wise metrics, we fit the generative distribution by sampling from the first examination step, *i.e.*, directly conditioned on the ‘patient profile’ without giving extra past examination results. This aims to preserve enough sampling randomness, thus better reflect the underlying distribution. Instead, providing too much context, *e.g.*, detailed past examinations, could overly constrain the generation and make it excessively deterministic. For each examination item, we only use test cases where that specific examination was actually performed, ensuring that both the generative and real distributions are calculated using the same patient set.

Computational metrics. In addition to the generation quality, we assess computational efficiency, as the simulator must respond rapidly to support interactive training. Two metrics are reported: **minimal GPU**, the lowest number of GPU cards required for deployment. **Time (GPU·s)**, the average wall-clock time to generate a single examination result, multiplied by the minimal GPU count. All measurements were obtained on NVIDIA A100 80GB GPUs.

3.1.2 Results Analysis

The evaluation results in Table 1 compare DiagGym against strong open-source LLMs prompted as diagnostics world model to simulate diverse examination results (Section 5.4). We exclude closed-source API models due to their high latency and cost, which make them impractical for scalable RL training. Across nearly all metrics, DiagGym delivers state-of-the-art performance, combining high-fidelity generation with substantially greater computational efficiency, thereby providing a practical and reliable foundation for training diagnostic agents.

Instance-wise metrics with automated evaluator. DiagGym achieves the highest step-level similarity score (3.565) and full-chain consistency (96.91%), indicating both close alignment with ground truth and strong internal clinical coherence. In comparison, the largest baseline, Qwen2.5-72B, attains a consistency of 92.39% but a substantially lower similarity of 2.495. DeepSeek-v3-671B and MedGemma-27B perform moderately, with similarities of 2.576 and 2.438 and consistencies around 89%. The smallest model, Qwen2.5-7B, performs worst overall, with a similarity of 2.181 and a consistency of 81.64%.

Instance-wise metrics with physicians raters. As shown in Table 2, DiagGym achieves the highest similarity across physicians, with an average score of 4.49 across physicians. The baselines followed with DeepSeek at 4.09, Qwen2.5-72B at 3.97, and Medgemma-27B at 3.89. On the consistency axis, DiagGym attains a majority-vote consistency rate of 95.00%, compared with 74.00% for Medgemma-27B, 54.00% for DeepSeek-v3-671B, and 44.00% for Qwen2.5-72B. These results indicate that DiagGym not only matches the reference more closely but also maintains substantially stronger clinical coherence.

Physicians also provide qualitative assessments on each methods’ simulation failure modes. DeepSeek-v3-671B produces detailed, information-rich reports with broad coverage, but it frequently over-extrapolate from the diagnosis and clinical history, yielding overly severe positive findings. Also, narrative organization is sometimes disjointed with topic switching and repetition. Qwen2.5-72B offers comprehensive coverage but sometimes introduces unsupported neutral or false-positive findings with tenuous links to the core diagnosis. MedGemma-27B, as a medical-domain LLM, is more concise and structurally clear with fewer off-topic assertions, though it intermittently exhibited logical inconsistencies (e.g., conflicting statements across sections or implausible laboratory constellations). In contrast, DiagGym maintains balanced converge and close alignment with case context, largely avoids unwarranted positives and material contradictions. This qualitative profile mirrors its quantitative advantage on both similarity and consistency.

Examination-wise metrics. On fidelity, DiagGym achieves a 1-Wasserstein distance of 0.128 for numerical results and an FID of 0.747 for free-text outputs - the closest match to real-world distributions. Diversity is likewise strong, with a Normalized Variance of 3.46 for numerical outputs and LPIPS of 0.378 for free-text, values close to ground truth (5.31 and 0.427, respectively). Among the baselines, DeepSeek-v3-671B offers the most balanced performance, likely benefiting from its large parameter scale. It achieves the highest Normalized Variance at 24.56, while maintaining fidelity at a moderate level (Wasserstein Distance 1.336; FID 4.158). By contrast, Qwen2.5-7B and MedGemma-27B exhibit high diversity but suffer from large distribution gaps, whereas Qwen2.5-72B shows strong consistency but very low diversity, producing overly deterministic outputs.

Computational efficiency. Some baselines trade performance for heavy resource demands. For example, DeepSeek-v3-671B requires at least 16 A100-80GB GPUs and 62.72 GPU-s to perform a single simulation run. In contrast, DiagGym runs on a single A100-80GB, generating outputs in 0.52s while sustaining superior generative quality, making it well-suited for the rapid, repeated interactions required in diagnostic agent RL.

Summary. Collectively, these results demonstrate that DiagGym is a high-fidelity, diverse, and computationally efficient world model. It is reliable and well-suited to serve as a virtual clinical environment for dynamic diagnostic agent training with RL, substantially outperforming current open-source baselines.

3.2 Evaluation for DiagAgent

We next assess our final diagnostic model, **DiagAgent**, against leading LLMs and agentic systems (baselines described in Section 5.4), focusing on its ability to manage complete multi-turn diagnostic trajectories.

Table 1 | Quantitative comparison of different models as diagnostics world model for generating synthetic patient examination results. We adopt metrics including Similarity, Consistency, Fidelity (Wasserstein Distance, FID), Diversity (Normalized Variance, LPIPS), as well as inference time and resource consumption.

| Model | Computational Metrics | | Instance-wise Metrics | | Examination-wise Metrics | | | |
|------------------|-----------------------|---------------|-----------------------|-----------------|--------------------------|-----------------------|--------------|--------------|
| | | | | | Numerical | | Free-text | |
| | Minimal GPUs↓ | Time (GPU-s)↓ | Similarity↑ | Consistency(%)↑ | Normalized Variance↑ | Wasserstein Distance↓ | LPIPS↑ | FID↓ |
| GT | - | - | - | - | 5.31 | - | 0.427 | - |
| DeepSeek-v3-671B | 16 | 62.72 | 2.576 | 88.81 | 24.56 | 1.336 | 0.237 | 4.158 |
| Qwen2.5-7B | 1 | 0.54 | 2.181 | 81.64 | 20.18 | 9.680 | 0.256 | 4.800 |
| Qwen2.5-72B | 4 | 18.68 | 2.495 | 92.39 | 1.21 | 1.839 | 0.183 | 4.901 |
| MedGemma-27B | 2 | 9.1 | 2.438 | 89.87 | 18.70 | 16.936 | 0.341 | 4.158 |
| DiagGym | 1 | 0.52 | 3.565 | 96.91 | 3.46 | 0.128 | 0.378 | 0.747 |

Table 2 | Human ratings comparing DiagGym with baseline models. Similarity (0-5) is reported for each of the three physicians and as the mean across physicians. Consistency is a binary judgment, we report the per-physician percentage of cases judged clinically coherent and the majority-vote consistency rate (percentage of cases deemed coherent by at least two of three physicians).

| Model | Similarity | | | | Consistency(%) | | | |
|------------------|-------------|-------------|-------------|-------------|----------------|--------------|--------------|---------------|
| | Physician 1 | Physician 2 | Physician 3 | Avg. Score | Physician 1 | Physician 2 | Physician 3 | Majority Vote |
| DeepSeek-v3-671B | 4.66 | 4.49 | 3.11 | 4.09 | 54.00 | 58.00 | 42.00 | 54.00 |
| Qwen2.5-72B | 4.50 | 4.37 | 3.04 | 3.97 | 46.00 | 44.00 | 32.00 | 44.00 |
| MedGemma-27B | 4.56 | 4.28 | 2.82 | 3.89 | 73.00 | 75.00 | 56.00 | 74.00 |
| DiagGym | 4.71 | 4.70 | 4.05 | 4.49 | 96.00 | 94.00 | 92.00 | 95.00 |

3.2.1 Evaluation Settings

We evaluate **DiagAgent** on **DiagBench**, focusing on assessing its ability to manage multi-turn diagnostic trajectories. The evaluation set comprises 750 physician-validated patient cases from MIMIC-IV, each containing three elements: (i) **initial inquiry**: patient initialized presentation details (chief complaint, current and past medical history, and other relevant information), forming the starting point for the diagnostic process; (ii) **referenced multi-turn diagnostic trajectory**: a physician-curated sequence extracted from real EHR records, serving as the ground-truth reference; (iii) **final diagnosis**: the final confirmed clinical diagnosis outcome. Notably, cases are structured following the same simulation case-construction pipeline used earlier, ensuring each includes a patient profile compatible with the simulator. Furthermore, 99 cases are associated with **973 physician-authored rubrics** that delineate the critical interaction points within the diagnostic process. The detailed DiagBench construction pipeline can be found in Section 5.3.

We evaluate **DiagAgent** in two complementary settings with corresponding metrics, namely, single-turn and end-to-end evaluation, as detailed below.

Single-turn Evaluation

In this setting, we evaluate the DiagAgent in the single-turn setting. As shown in Figure 3a,b, here, we leverage both the ground truth ‘initial inquiry’ and partial ‘referenced multi-turn diagnostic trajectory’ as input. The DiagAgent is directly forced (prompt details can be found in Section 5.4) to provide an examination recommendation or make a final diagnosis based on the preceding oracle diagnostic trajectory, extending the process by one additional turn, without self-deciding which action to perform next. Such single-turn evaluations are conducted at each agent response turn recorded in the referenced multi-turn trajectory, resulting in 4,485 turns for evaluation, consisting 3,735 intermediate turns for examination recommendation and 750 final turns for diagnosis.

For examination recommendation turns, we calculate the hit ratio based on whether the suggested examination actually appears in the this admission’s historical records in MIMIC-IV. For diagnosis, we still adopt accuracy. Notably, to avoid issues with synonyms and the inclusion relationships between examinations, we utilize GPT-4o to judge whether a particular examination recommendation appears within the remaining part in the

‘referenced multi-turn diagnostic trajectory’ with [prompt 6](#).

End-to-End Evaluation

In this setting, we adopt an end-to-end evaluation approach. The diagnostic trajectory is initialized with the ‘initial inquiry’, after which diagnostic agents continually interact with the environment and sequentially propose examination queries until they determine that a final diagnostic decision can be made. Throughout the trajectory, all returned examination results are simulated by the DiagGym, conditioned on the background ‘patient profile,’ ensuring that all queried information is available.

This evaluation setting more closely reflects real-world clinical practice, highlighting the model’s ability to dynamically construct a complete diagnostic trajectory, autonomously determining both the timing and type of actions based on the patient’s evolving condition. While the assessment necessarily relies on the external diagnostics world model to simulate examination results, this stems from a fundamental limitation of real-world EHRs: they only contain examinations that were actually performed. Consequently, when the diagnostic agent suggests an examination that was not carried out for a given patient, there is no corresponding result in the EHR, making it impossible to evaluate that decision or its downstream effects. These inherent gaps prevent the use of real-world EHRs for fully interactive, end-to-end evaluation, as the diagnostic trajectory would be repeatedly interrupted by missing information.

After obtaining the complete predicted diagnostic trajectory, we evaluate performance from two perspectives.

First, as shown in [Figure 4a](#), for the whole 750 cases we employ automatic metrics to assess the efficacy of examination recommendations and the accuracy of the final diagnosis. (i) **Examination recommendation** compares the examination items proposed by the model with those in the reference multi-turn diagnostic trajectory, and computing precision, recall, and F1-score to measure recommendation quality; (ii) **Final diagnosis** assesses the accuracy of the model’s ultimate diagnosis after completing the multi-turn interaction, by directly comparing it to the ground-truth diagnosis.

Second, as shown in [Figure 4c](#), for the 99 cases accompanied by manually written rubrics, we additionally derived a weighted rubric score to more accurately reflect physicians’ satisfaction with the generated diagnostic interactions. LLM-as-a-judge is adopted here. We prompt GPT-4o to judge whether the generated diagnosis process satisfies each rubric criterion and calculate the final weighted totaling scores based on each rubric weights, similar as HealthBench [\[22\]](#).

Detailed definitions of these metrics are provided in [Supplementary 10.1.2](#).

3.2.2 Results Analysis

In this section, we analyze the performance of **DiagAgent** from 3 perspectives: single-turn evaluation ([Figure 3c](#) and [Table 3](#)), end-to-end evaluation with automatic rubrics ([Figure 4b](#) and [Table 4](#)) and end-to-end evaluation with rubric-based metrics ([Figure 4d](#)).

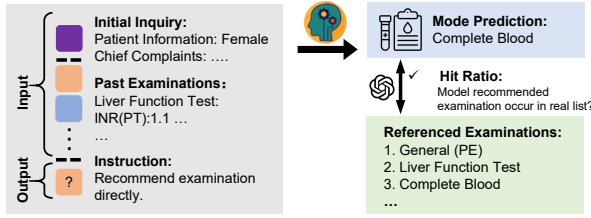
We compare DiagAgent against various state-of-the-art large language models (LLMs) prompted as diagnostic agents (detailed instructions are provided in [Method Section 5.4](#)), including GPT-4o, Claude-4, GPT-OSS, Qwen2.5, Qwen3, Llama3.3, DeepSeek-V3, OpenBioLLM, Baichuan-M1, MedGemma, as well as the latest agentic systems such as MedAgents, MDAgents.

Single-turn Evaluation Analysis

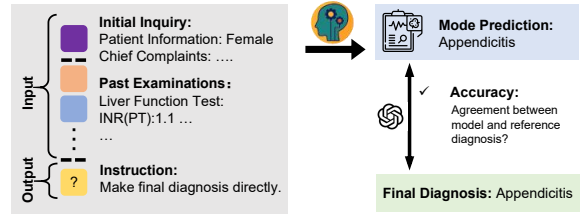
[Table 3](#) and [Figure 3c](#) summarize the performance of all models under the oracle single-turn evaluation setting. Models are grouped into three categories: Basic LLMs, Agentic Systems, and Our DiagAgent Variants.

Our **DiagAgent** models deliver substantial gains across both tasks. DiagAgent-7B attains a Hit Ratio of 71.12% and diagnosis accuracy of 85.03%. DiagAgent-8B reaches 54.61% and 81.10%, while the largest model, DiagAgent-14B, achieves 67.54% and 86.73%, respectively. Compared with the strongest single LLM baseline, MedGemma (27.07% / 68.90%), DiagAgent-14B improves examination recommendation by over 40 percentage points and diagnosis accuracy by nearly 18 points. Against large general-purpose models such as DeepSeek-v3, the gains are even more pronounced—exceeding 48 points in Hit Ratio and 17 points in accuracy. These results show that post-training within DiagGym consistently and substantially boosts both recommendation and diagnostic performance, surpassing state-of-the-art general-purpose and medical-specialized LLMs.

a. Single-turn Evaluation Setting (Examination Recommendation)



b. Single-turn Evaluation Setting (Final Diagnosis)



c. Single-turn Evaluation Results

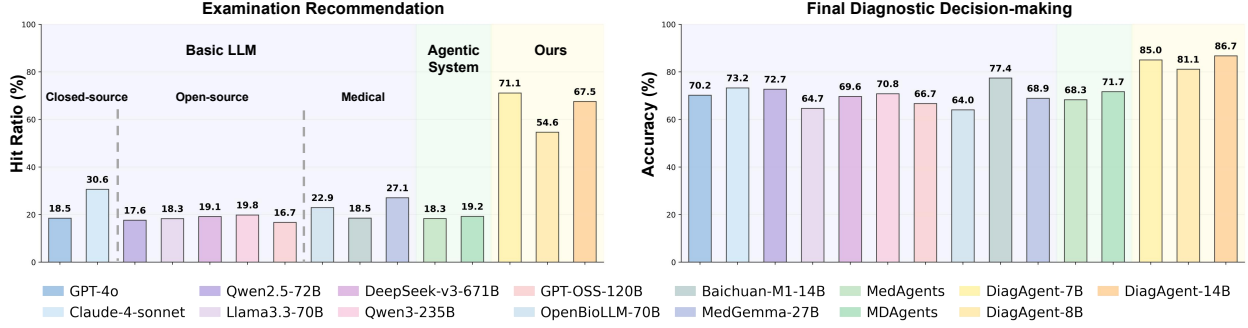


Figure 3 | Overview of single-turn evaluation settings and results. **a** shows the single-turn evaluation setting for examination recommendation measured with the hit ratio. **b** the single-turn evaluation setting for final diagnosis measured with the accuracy. **c** compares our DiagAgent variants against 10 leading LLMs and 2 agentic systems on examination recommendation and diagnostic decision-making in the single-turn setting.

Table 3 | Performance comparison of different diagnostic models under single-turn evaluation. We adopt “Hit Ratio” for examination recommendation and “Accuracy” for final diagnosis. Agentic systems use DeepSeek-v3 as their base model.

| Model | Size | Year | Hit Ratio(%) | Diagnosis Accuracy(%) |
|-----------------|------|---------|--------------|-----------------------|
| Basic LLM | | | | |
| GPT-4o | - | 2024.8 | 20.21 | 72.13 |
| Claude-4-sonnet | - | 2025.5 | 31.63 | 74.31 |
| Qwen2.5 | 72B | 2024.9 | 19.20 | 74.80 |
| Llama3.3 | 70B | 2024.12 | 19.97 | 66.27 |
| DeepSeek-v3 | 671B | 2025.3 | 20.08 | 72.27 |
| Qwen3 | 235B | 2025.7 | 21.39 | 72.40 |
| GPT-OSS | 120B | 2025.8 | 17.37 | 67.37 |
| OpenBioLLM | 70B | 2024.4 | 23.53 | 66.27 |
| Baichuan-M1 | 14B | 2025.2 | 19.60 | 78.93 |
| MedGemma | 27B | 2025.7 | 28.57 | 60.53 |
| Agentic System | | | | |
| MedAgents | - | 2024.1 | 19.31 | 70.27 |
| MDAgents | - | 2024.10 | 20.24 | 73.47 |
| Our Method | | | | |
| DiagAgent | 7B | - | 72.56 | 85.60 |
| | 8B | - | 56.57 | 82.27 |
| | 14B | - | 68.49 | 87.87 |

Among **Basic LLMs**, closed-source general LLMs such as GPT-4o and Claude-4-sonnet show moderate capability. Claude-4-sonnet achieves 30.62% / 73.25%, outperforming GPT-4o (18.46% / 70.17%), yet both remain below 35% in examination recommendation—indicating that even the strongest proprietary LLMs struggle to recall the breadth of diagnostic steps needed in realistic clinical workflows. Open-source general

LLMs (Qwen2.5-72B, Llama3.3-70B, DeepSeek-v3, Qwen3-235B, GPT-OSS-120B) perform worse, with Hit Ratios between 16.71% and 19.82% and diagnosis accuracy ranging from 64.65% (Llama3.3-70B) to 70.81% (Qwen3-235B). Notably, scaling to extreme parameter scales—as in DeepSeek-v3 (671B)—does not guarantee improved examination recommendation without targeted domain adaptation. Medical-specialized LLMs consistently outperform general-purpose models in recommendation. MedGemma attains the highest Hit Ratio among single LLM baselines (27.07%), while Baichuan-M1 achieves the best diagnosis accuracy (77.39%). OpenBioLLM also remains competitive, with 22.94% / 64.01%.

For **Agentic Systems**, MedAgents and MDAgents—both based on DeepSeek-v3—fail to deliver substantial improvement over the base model. MedAgents scores 18.31% / 68.26%, slightly below DeepSeek-v3, while MDAgents improves marginally to 19.24% / 71.66%. These limited gains suggest that, without a well-aligned base model, multi-agent coordination alone cannot overcome the challenges of multi-turn diagnostic reasoning, and may even worsen hallucination-driven errors, such as prematurely concluding that sufficient information has been gathered.

End-to-end Evaluation Analysis on Automatic Metrics

The results in Table 4 and Figure 4b show that DiagAgent models achieve the highest scores across all evaluation metrics. DiagAgent-14B engages in an average of 6.77 conversational turns—substantially longer than the 2–4 turns typical of large basic LLMs, such as DeepSeek-v3 (2.49 turns) or Qwen2.5-72B (2.47 turns). This extended interaction facilitates more comprehensive evidence gathering, yielding a recall of 52.74% compared with 12.20% for DeepSeek-v3. Notably, this gain in recall does not come at the cost of precision, which remains high at 43.87%, exceeding the 32.60% of DeepSeek-v3. The resulting F1-score of 47.89 is the highest among all models tested, and translates directly into superior diagnostic accuracy: 61.63%, more than 15 percentage points above the strongest baseline. Similar patterns hold for DiagAgent-7B and DiagAgent-8B, which also consistently outperform all basic LLMs in both recommendation and diagnosis tasks. These findings indicate that DiagAgent’s enhanced multi-turn exploration leads to more informed and accurate diagnostic decision-making.

Among basic LLM baselines, advanced models with strong general reasoning—Claude-4-sonnet, Qwen3-235B, DeepSeek-v3, and GPT-OSS-120B—tend to achieve higher diagnostic accuracy than earlier-generation systems. In contrast, medical-specialized LLMs such as MedGemma, OpenBioLLM, and Baichuan-M1, while advantaged in domain knowledge, show no clear benefit in this dynamic setting. This gap underscores a broader limitation: without robust dynamic decision-making capabilities, even models rich in clinical knowledge struggle to actively gather and integrate new evidence across multiple turns.

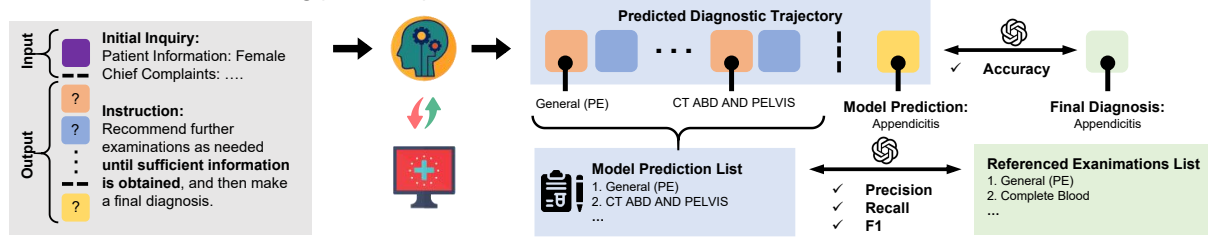
Relative to current agentic diagnosis systems, DiagAgent again demonstrates clear superiority. MedAgents and MDAgents, both built on DeepSeek-v3 and incorporating “expert discussion” mechanisms in which multiple LLM agents approach the case from different specialist perspectives, fail to improve upon their base model. MedAgents averages 2.31 turns—slightly fewer than DeepSeek-v3—with precision 29.86%, recall 11.21%, F1-score 16.30, and accuracy 44.40%. MDAgents performs similarly, with 2.40 turns, precision 30.12%, recall 11.25%, F1-score 16.38, and accuracy 43.55%. These patterns suggest that expert-discussion frameworks, without a well-aligned base model, may encourage premature closure—terminating evidence collection before critical information is obtained—thus limiting both recommendation quality and diagnostic accuracy.

In summary, these results highlight the importance of post-training LLMs within clinically realistic, interactive environments. By endowing models with the capacity to determine the timing and content of their multi-turn actions, DiagGym substantially improves both the quality of examination recommendations and the accuracy of final diagnoses, narrowing the gap towards deployable, decision-capable clinical AI systems.

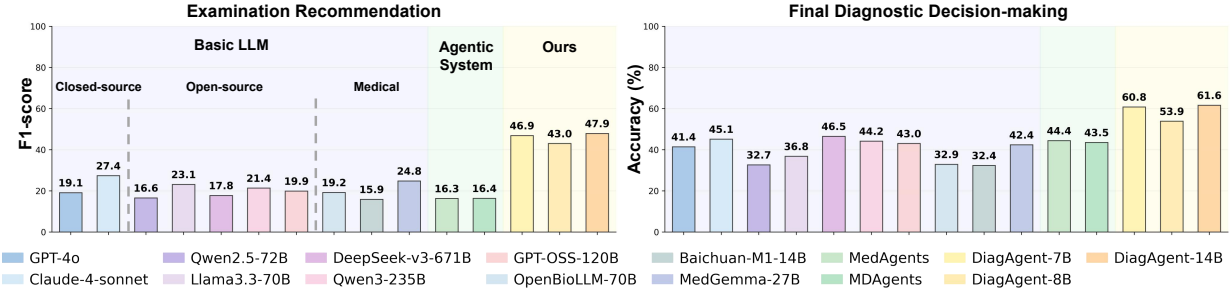
End-to-end Evaluation Analysis on Rubric-based Metrics

The results in Figure 4d indicate that our dynamic training paradigm, which explicitly shapes intermediate diagnostic trajectories, better satisfies high-importance rubrics. As a result, DiagAgent-14B achieves an weighted rubric scores of 32.86%, surpassing the strongest basic LLM, Qwen3-235B at 24.49%, by 8.37 percentage points and outperforming the best prior agentic baseline, MDAgent at 21.64%, by 11.22 percentage points. These results indicate that DiagAgent’s training, which explicitly optimizes for intermediate trajectory quality, enables substantially improved alignment with high-importance clinical rubrics. Notably, DiagAgent demonstrates stronger performance particularly on rigorously weighted criteria, suggesting enhanced procedural

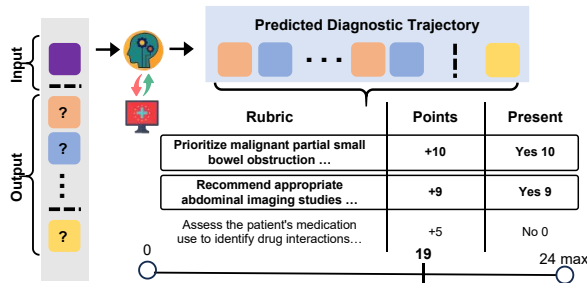
a. End-to-end Evaluation Setting (Automatic)



b. End-to-end Results (Automatic)



c. End-to-end Evaluation Setting (Rubric-based)



d. End-to-end Results (Rubric-based)

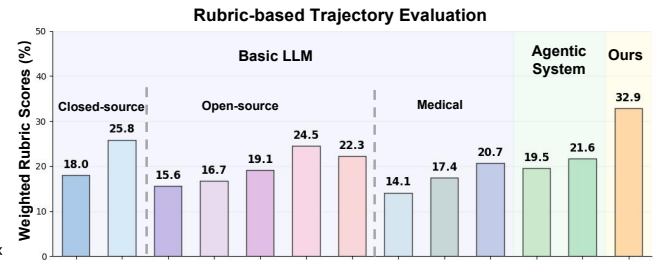


Figure 4 | Overview of end-to-end evaluation settings and results. In this setting, diagnostic agents are evaluated through end-to-end finishing the entire diagnostic trajectory by interaction with the external diagnostics world model. **a** illustrates the end-to-end evaluation pipeline with automatic metrics to assess examination recommendation efficacy and diagnostic accuracy. **b** compares our DiagAgent with 10 LLMs and 2 more agentic systems under end-to-end evaluation settings with automatic metrics. **c** illustrates our end-to-end evaluation pipeline with rubric-based metrics. A judge model evaluates the full diagnostic trajectory of a diagnostic model against physician-curated rubrics, which specify criteria, clinical importance weights (Points), and whether the criterion was satisfied (Present). **d** compares the aggregate weighted proportion of satisfied rubrics(%) across different LLMs in different model sizes.

reasoning and better clinical step coverage, beyond just achieving the correct final diagnosis.

Among baseline models, general-purpose frontier LLMs achieve the strongest rubric scores. Claude-4-sonnet and Qwen3-235B outperform most open-source and medical models, consistent with better generalization and broader knowledge that help them satisfy process-focused criteria across the diagnostic trajectory. In contrast, domain-specialized models are often smaller and show weaker reasoning and generalization. As a result, even when they reach correct final diagnostic result, their intermediate trajectory frequently fall short on high-importance rubrics for history taking, hypothesis management, and examination recommendation.

Agentic systems such as MDAgent, which leverage multi-agent discussion to enhance coverage, provide some performance gains over their base model DeepSeek-v3. MDAgent improves from 19.07% to 21.64% in weighted rubric scores, indicating that multi-agent discussion can surface complementary considerations and cover more rubric items. The improvement remains limited, which suggests that inference-time orchestration alone cannot ensure process quality.

Overall, the results on DiagBench reinforce that dynamic training focused on intermediate reasoning processes enables significantly improved satisfaction of high-value clinical procedures. DiagAgent consistently surpasses

Table 4 | Performance comparison of different diagnostic models under virtual patient environment. Metrics include average conversation turns, precision, recall, F1-score for examination recommendation, and final diagnostic accuracy. Agentic systems use DeepSeek-v3 as their base model.

| Model | Size | Year | Avg. Turns | Precision | Recall | F1 | Accuracy(%) |
|-----------------|------|---------|------------|--------------|--------------|--------------|--------------|
| Basic LLM | | | | | | | |
| GPT-4o | - | 2024.8 | 3.30 | 32.31 | 15.33 | 20.80 | 43.24 |
| Claude-4-sonnet | - | 2025.5 | 3.91 | 33.06 | 23.41 | 27.41 | 46.14 |
| Qwen2.5 | 72B | 2024.9 | 2.47 | 36.03 | 12.36 | 18.40 | 44.30 |
| Llama3.3 | 70B | 2024.12 | 4.25 | 29.02 | 22.30 | 25.22 | 38.46 |
| DeepSeek-v3 | 671B | 2025.3 | 2.49 | 36.49 | 13.58 | 19.80 | 47.08 |
| Qwen3 | 235B | 2025.7 | 3.34 | 31.20 | 19.23 | 23.80 | 45.36 |
| GPT-OSS | 120B | 2025.8 | 4.08 | 25.94 | 16.16 | 19.92 | 44.02 |
| OpenBioLLM | 70B | 2024.4 | 2.59 | 35.96 | 15.33 | 21.49 | 34.35 |
| Baichuan-M1 | 14B | 2025.2 | 2.30 | 32.42 | 12.04 | 17.56 | 33.16 |
| MedGemma | 27B | 2025.7 | 4.10 | 35.38 | 21.28 | 26.58 | 44.30 |
| Agentic System | | | | | | | |
| MedAgents | - | 2024.1 | 2.31 | 33.93 | 12.56 | 18.34 | 45.62 |
| MDAgents | - | 2024.10 | 2.40 | 34.04 | 12.56 | 18.35 | 44.96 |
| Our Method | | | | | | | |
| DiagAgent | 7B | - | 5.45 | 46.14 | 46.33 | 46.23 | 60.78 |
| | 8B | - | 5.73 | 39.66 | 43.15 | 41.33 | 53.45 |
| | 14B | - | 6.66 | 42.12 | 52.12 | 46.59 | 61.27 |

both large generic LLMs and current agentic frameworks, underscoring the essential role of fine-grained, process-aware training for safe and effective clinical decision support.

3.3 Ablation Study

We conducted ablation experiments under the end-to-end evaluation setting to assess three aspects of the proposed framework: (i) whether the reinforcement learning in virtual environment outperforms supervised fine-tuning (SFT) at the same model scale; (ii) the impact of reward design, comparing diagnosis-only rewards with dual rewards incorporating both diagnosis accuracy and examination-recommendation quality; and (iii) the generality of **DiagAgent** across different model sizes and families.

Experimental Design

As outlined in Table 5, for each base model, we first establish a zero-shot baseline in which the LLM answers without fine-tuning. We then apply full supervised finetuning, in which all cases are converted into multi-turn diagnostic dialogues for supervised training (Supplementary 10.2), bypassing the simulator and reinforcement learning pipeline. Finally, we test three DiagAgent configurations: (i) **cold-start only**: supervised fine-tuning on a small subset to learn output format; (ii) **cold-start + RL with diagnosis reward**: RL optimising diagnosis accuracy only; (iii) **full DiagAgent**: diagnose reward plus examination-recommendation reward.

Supervised Finetuning (SFT) vs. Reinforcement Learning (RL)

As shown in Table 5, zero-shot baselines perform poorly (diagnosis accuracy: 16.38% for Qwen2.5-7B; 33.83% for Qwen2.5-14B), highlighting the difficulty of managing interactive trajectories without domain adaptation. Full SFT improves accuracy (44.40%, 45.98%, and 45.35% for Qwen2.5-7B, Qwen2.5-14B, and Llama3.1-8B, respectively) and recommendation quality, but the gains are limited by static trajectory data extracted from MIMIC-IV discharge notes, which do not reflect the dynamic branching in interactive consultations. In contrast, Our full DiagAgent consistently surpasses SFT—*e.g.*, Qwen2.5-7B reaches 60.78% accuracy,

Table 5 | Results of ablation study. To compare the effect of each training pipeline for different model sizes and families, we report the Avg. Turns, precision, recall and F1-score for examination recommendations and accuracy for diagnosis under the end-to-end setting.

| Method | Instruction Tuning | | Reinforcement Learning | | Avg. Turns | Examination Recommendation | | | Diagnosis |
|-------------|--------------------|------------|------------------------|------------------|------------|----------------------------|--------|-------|-------------|
| | Full SFT | Cold Start | Recommend Reward | Diagnosis Reward | | Precision | Recall | F1 | Accuracy(%) |
| Qwen2.5-7B | | | | | | | | | |
| Baseline | ✗ | ✗ | ✗ | ✗ | 1.97 | 32.30 | 15.64 | 21.08 | 16.38 |
| Full SFT | ✓ | ✗ | ✗ | ✗ | 7.98 | 37.57 | 51.79 | 43.55 | 44.40 |
| DiagGym | ✗ | ✓ | ✗ | ✗ | 8.36 | 32.77 | 47.43 | 38.76 | 36.47 |
| | ✗ | ✓ | ✗ | ✓ | 4.46 | 34.32 | 31.35 | 32.76 | 59.41 |
| | ✗ | ✓ | ✓ | ✓ | 5.45 | 46.14 | 46.33 | 46.23 | 60.78 |
| Llama3.1-8B | | | | | | | | | |
| Baseline | ✗ | ✗ | ✗ | ✗ | 4.36 | 21.59 | 17.58 | 19.38 | 25.16 |
| Full SFT | ✓ | ✗ | ✗ | ✗ | 7.89 | 37.02 | 50.49 | 42.72 | 45.35 |
| DiagGym | ✗ | ✓ | ✗ | ✗ | 9.03 | 31.43 | 48.39 | 38.11 | 31.83 |
| | ✗ | ✓ | ✗ | ✓ | 5.15 | 34.84 | 35.50 | 35.17 | 52.52 |
| | ✗ | ✓ | ✓ | ✓ | 5.73 | 39.66 | 43.15 | 41.33 | 53.45 |
| Qwen2.5-14B | | | | | | | | | |
| Baseline | ✗ | ✗ | ✗ | ✗ | 3.61 | 29.97 | 15.43 | 20.37 | 33.83 |
| Full SFT | ✓ | ✗ | ✗ | ✗ | 7.63 | 37.83 | 51.38 | 43.58 | 45.98 |
| DiagGym | ✗ | ✓ | ✗ | ✗ | 8.56 | 32.86 | 48.27 | 39.10 | 35.41 |
| | ✗ | ✓ | ✗ | ✓ | 5.51 | 34.37 | 37.69 | 35.95 | 58.36 |
| | ✗ | ✓ | ✓ | ✓ | 6.66 | 42.12 | 52.12 | 46.59 | 61.27 |

Qwen2.5-14B 61.63%, and Llama3.1-8B 53.85%—while reducing dialogue length, indicating more efficient and decisive reasoning. In addition, DiagAgent consistently reduces the average number of dialogue turns, demonstrating both higher efficiency and stronger capability in dynamic interaction.

Effect of Reward Design

To demonstrate the effectiveness of reward design, we first conduct a cold start phase, where each model is supervised-finetuned on a subset of the training set to learn the output format. Adding RL with only the Diagnose Reward yields large gains over SFT, for example, Qwen2.5-7B improves from 44.40% (Full SFT) to 59.41%, Qwen2.5-14B from 45.98% to 59.73%, and Llama3.1-8B from 45.35% to 54.12%. However, but F1 scores for examination recommendation remain low ($\leq 36\%$). Introducing the Examination Recommendation Reward markedly improves F1 across all models (Qwen2.5-7B: 32.76% \rightarrow 46.86%; Qwen2.5-14B: 35.95% \rightarrow 47.89%), with slight additional gains in accuracy. This confirms the importance of dual-reward shaping for balancing precision in diagnosis and quality in examination planning.

Model Size and Family

Reinforcement learning with virtual environment benefits all LLMs training, for example, in diagnosis accuracy, from 16.38% to 60.78% for Qwen2.5-7B, 25.16% to 53.85% for Llama3.1-8B, and 33.83% to 61.63% for Qwen2.5-14B, proving the effectiveness of our method. Larger or intrinsically stronger base models achieve higher post-training ceilings: Qwen2.5-14B delivers the best overall performance (61.63% accuracy, 47.89% F1), followed by Qwen2.5-7B (60.78%, 46.86%), while Llama3.1-8B lags (53.85%, 43.02%). This suggests that while DiagAgent’s exploration-driven optimisation is broadly applicable, the quality of the base model constrains the attainable performance upper bound.

3.4 Case Study

Qualitative results from our diagnostics world model (DiagGym) and our diagnostic agent (DiagAgent) are discussed below, providing insight into their fidelity and reasoning capabilities within complex clinical workflows.

3.4.1 DiagGym Case

Supplementary Figure 1 compares outputs from the DiagGym with ground-truth examination results for a representative patient. The case involves a woman presenting with painless jaundice, with a history of breast cancer and non-ischemic cardiomyopathy. Upon admission, laboratory tests including liver function and bilirubin levels both supported a diagnosis of biliary obstruction; the final confirmed diagnosis was ‘CBD obstruction from common hepatic duct mass’.

The world model’s simulated predictions closely match the true clinical findings. A key indicator, **Total Bilirubin**, was predicted at 6.7 mg/dL versus the recorded value of 4.3 mg/dL (both well above the reference range of 0–1.5 mg/dL). This pronounced hyperbilirubinemia is consistent with the patient’s presentation (jaundice, scleral icterus) and strongly supports the diagnosis of “CBD obstruction from common hepatic duct mass”. Minor numerical variations occur across other laboratory results, but none alter the clinical interpretation. Such variability illustrates the model’s ability to generate plausible but non-identical results, preserving both realism and diversity.

3.4.2 DiagAgent Case

A Case of Dynamic Diagnostic Trajectory

As shown in Supplementary Figure 2, we present a case study illustrating **DiagAgent-14B’s** dynamic diagnostic trajectory within the DiagGym environment, modeling a typical appendicitis work-up. Each case record includes the initial inquiry, final diagnosis, interactive exchanges between diagnostic agent and simulator, and a reference ground-truth diagnostic timeline drawn from clinical records.

First, the model’s decisions follow standard reasoning. Upon receiving initial symptoms—abdominal pain migrating to the right lower quadrant, nausea, diarrhoea, and anorexia—the diagnostic agent prioritises appendicitis in the differential and orders a complete blood count (CBC). When the CBC reveals an abnormally high neutrophil count, further pointing toward infection or inflammation, the diagnostic agent appropriately requests a CT scan of the abdomen and pelvis with contrast, which display a dilated, fluid-filled appendix with periappendiceal fat stranding, confirming acute appendicitis. Throughout the process, each diagnostic step and rationale aligns closely with the reference timeline, demonstrating reliable differential diagnosis.

Second, the dynamic environment’s responses remain consistent with the patient’s case summary and expected clinical progression. For instance, when the CBC is requested, it provides results consistent with an acute inflammatory, including elevated white blood cell and neutrophil counts; when the CT is ordered, it returns hallmark imaging features of appendicitis. This realistic feedback ensures that the agent’s decision-making unfolds in a manner faithful to real clinical workflows.

Cases of Rubric-based Evaluation Analysis

To further scrutinize the procedural quality of DiagAgent’s intermediate diagnostic steps, we employed an evaluation based on physician-curated rubrics that assess the integrity of multi-turn clinical interactions.

A typical **successful case** (left lower extremity infection) is presented in Supplementary Figure 3, illustrating the agent’s robust procedural performance. The agent exhibits high coherence and dynamic strategy adjustment: it orders a CBC and, upon receiving a result showing an elevated but non-critical neutrophil count, appropriately requests a wound culture to identify the causative organism. Following the positive culture for *Staphylococcus aureus*, the agent orders a blood culture to rule out bacteremia and then efficiently terminates further investigation after receiving a negative result, avoiding over-testing. The high scores on the procedural rubrics confirm that the agent’s decision-making process is both efficient and clinically sound, successfully meeting criteria for prioritizing tests, interpreting results, and achieving evidence-driven closure.

We also provide an illustrative **failure case** (ruptured ectopic pregnancy with hemodynamic instability) in Supplementary Figure 4 to showcase model’s fail mode. While the agent’s diagnostic reasoning is highly effective—it correctly orders hCG and subsequent pelvic ultrasound based on the patient’s unstable presentation, rapidly confirming the diagnosis—the evaluation reveals critical omissions in immediate emergency care. Specifically, the agent fails to satisfy the highest-weighted rubrics concerning emergency resuscitation and surgical team notification. It is important to emphasize that DiagAgent is primarily designed as a diagnostic reasoning model. Its core capability of accurate differential diagnosis and sequential information gathering

remains intact, validating its utility for diagnostic quality enhancement despite the observed gap in acute therapeutic and stabilization management, which falls outside its initial scope of contribution.

Summary

Across these examples, the reasoning of our diagnostic agent and the diagnostics world model’s interactive responses are coherent, clinically sound, and well matched to ground-truth trajectories. Together, they demonstrate high fidelity in reproducing authentic diagnostic trajectories and support its utility for training and evaluating diagnostic agents in long-term diagnosis management.

4 Discussion

Large language models (LLMs) have achieved notable success across a range of clinical tasks [23, 8, 9, 1, 10, 3, 24], yet they remain fundamentally limited in dynamic, sequential decision-making. Even state-of-the-art systems often struggle with real-time diagnostic reasoning—deciding which examinations to order, when to order them, and how to coordinate an efficient, end-to-end diagnostic process [10, 14, 16]. Unlike human physicians, who adaptively update decisions as new information emerges, current LLMs frequently fail to manage trajectories effectively under uncertainty.

Recent efforts to address this gap have focused on enhancing instruction-tuning datasets [4, 25], for example through synthetic multi-turn dialogues or extracting decision paths from retrospective records. While useful, these strategies face two persistent challenges: (i) **scarcity of high-quality interactive data**: even experienced clinicians may disagree on optimal sequencing for complex cases; (ii) **conservatism in training data**: historical records are dominated by guideline-driven trajectories, limiting generalization to atypical or rare scenarios.

Overview of Our Approach

In this paper, we propose the first dynamic training framework that optimize diagnostic agents in a virtual clinical environment, enabling them to explore, interact, and optimise decision-making policies via RL, aligning current LLMs with long-term diagnostic trajectories. By simulating realistic clinical feedback, our proposed diagnostics world model, DiagGym, exposes agents to a broader and more diverse distribution of patient trajectories—including rare and unconventional cases—allowing them to learn from both successes and errors. Extensive experiments demonstrate that our final DiagAgent consistently improves both diagnostic accuracy and examination-recommendation quality across diverse LLM families and scales, showcasing the successful transformation of current LLMs from static patient consultation to long-term patient management.

Main Contribution

A diagnostics world model serving as a virtual clinical environment for end-to-end agentic RL – at the core of our method is the world model, DiagGym, a fine-tuned LLM that goes beyond replaying historical records, and generates dynamic results for any requested examination. Unlike prior role-playing simulators such as AgentClinic [17], AgentHospital [26], and SDBench [27], which are constrained by static, pre-collected data, DiagGym can generate novel, clinically plausible trajectories beyond the original patient records, serving as a dynamic virtual clinical environment for clinical agent evolving.

First end-to-end RL platform for diagnostic agents – this paper provides an end-to-end, multi-turn RL framework in which agents interact iteratively with the clinical environment until a final diagnosis is reached. Agents actively explore diverse diagnostic trajectories, optimizing for both diagnostic accuracy and efficiency in examination recommendation. This allows them to handle complex, uncommon, or evolving patient scenarios that static supervised approaches cannot cover. We believe our approach provides the community with a robust resource for developing, testing, and comparing diagnostic agents in a dynamic long-term patient management manner, moving beyond the more commonly considered yet less clinically challenging static consultation scenarios for clinical LLMs.

Interactive and exploratory diagnostic reasoning – our final agent, DiagAgent, achieves substantial improvements over all evaluated LLMs and agentic systems. Crucially, DiagGym enables agents to go beyond imitation, discovering strategies that depart from physician-recorded trajectories when those lead to better clinical outcomes—opening the door to adaptive and potentially novel diagnostic approaches.

A new diagnosis evaluation benchmark focusing on interactive trajectories with fine-grained physician-written rubrics. To move beyond simplistic automatic metrics, we introduce a new benchmark, DiagBench, for evaluating multi-turn diagnostic trajectories. This benchmark features complex cases with detailed accompanying intermediate interactive reference, comprising 750 physician-validated cases, among which 99 cases are further annotated with 973 physician-written rubrics that assign weighted points to critical steps in the reasoning process. This enables a granular, process-oriented assessment of an agent’s clinical decision steps, offering deeper insights into how a diagnosis is reached.

Key Findings

Better alignment with dynamic decision-making – reinforcement learning consistently improves LLMs’ competence in planning and managing interactive diagnostic trajectories, addressing shortcomings of static fine-tuning. This dynamic, end-to-end training paradigm cultivates robust interactive competencies that are essential for real-world clinical deployment.

Superior to supervised fine-tuning (SFT) – across scales (7B–14B) and families (Qwen2.5, Llama3.1), diagnostic agent trained with reinforcement learning outperforms SFT by significant margins. This superiority is evident in dynamic diagnosis tasks, where DiagAgent’s self-exploration fosters adaptability to evolving patient scenarios, leading to higher accuracy, efficiency, and robustness compared to SFT’s limitations in handling incomplete or atypical information.

Dependence on base model quality - the intrinsic capability of the foundation model strongly shapes DiagGym’s upper performance bound. While DiagGym delivers robust gains even for moderate-scale (7B–14B) models, continuing to scale up to larger foundation models could unlock even greater advancements, suggesting a promising path for enhancing diagnostic agents.

Limitations and Future Work

First, the models evaluated in this work are relatively modest in scale (up to 14B parameters), which may constrain the framework’s full potential. Larger foundation models, such as DeepSeek-v3, GPT-OSS-120B, could yield qualitative leaps in performance by enhancing inherent reasoning capabilities and exploratory depth. Scaling DiagGym to these models might uncover even more efficient trajectories and robust adaptations, and future work should explore this to push the boundaries of diagnostic AI.

Secondly, we note that DiagAgent’s absolute scores on the rubric-based benchmark is modest. This is because the physician-authored rubrics reflect real-world clinical practice by awarding points for immediate treatment and patient management actions (e.g., “prepare for emergency transfer”). As our framework deliberately focuses on the diagnostic task, our agent was not trained to perform these out-of-scope therapeutic interventions. Extending the agent’s capabilities to integrate long-term management with timely treatment planning, interleaved with the diagnostic procedure, represents a clear direction for aligning clinical LLMs with practical usage demands in future work.

Thirdly, although the current DiagGym, as an diagnostics world model, provides a robust platform for simulating clinical environments involved in diagnosis, its scope remains relatively limited compared to the complexities of real-world clinical settings. Expanding the system to include additional modules, such as treatment planning and prognosis estimation, could make the virtual clinical environment more comprehensive, thereby promoting the development of advanced clinical agents capable of addressing more clinical tasks.

5 Methods

In this section, we provide additional details on the training process of DiagGym and DiagAgent, as well as the baseline methods used for comparison.

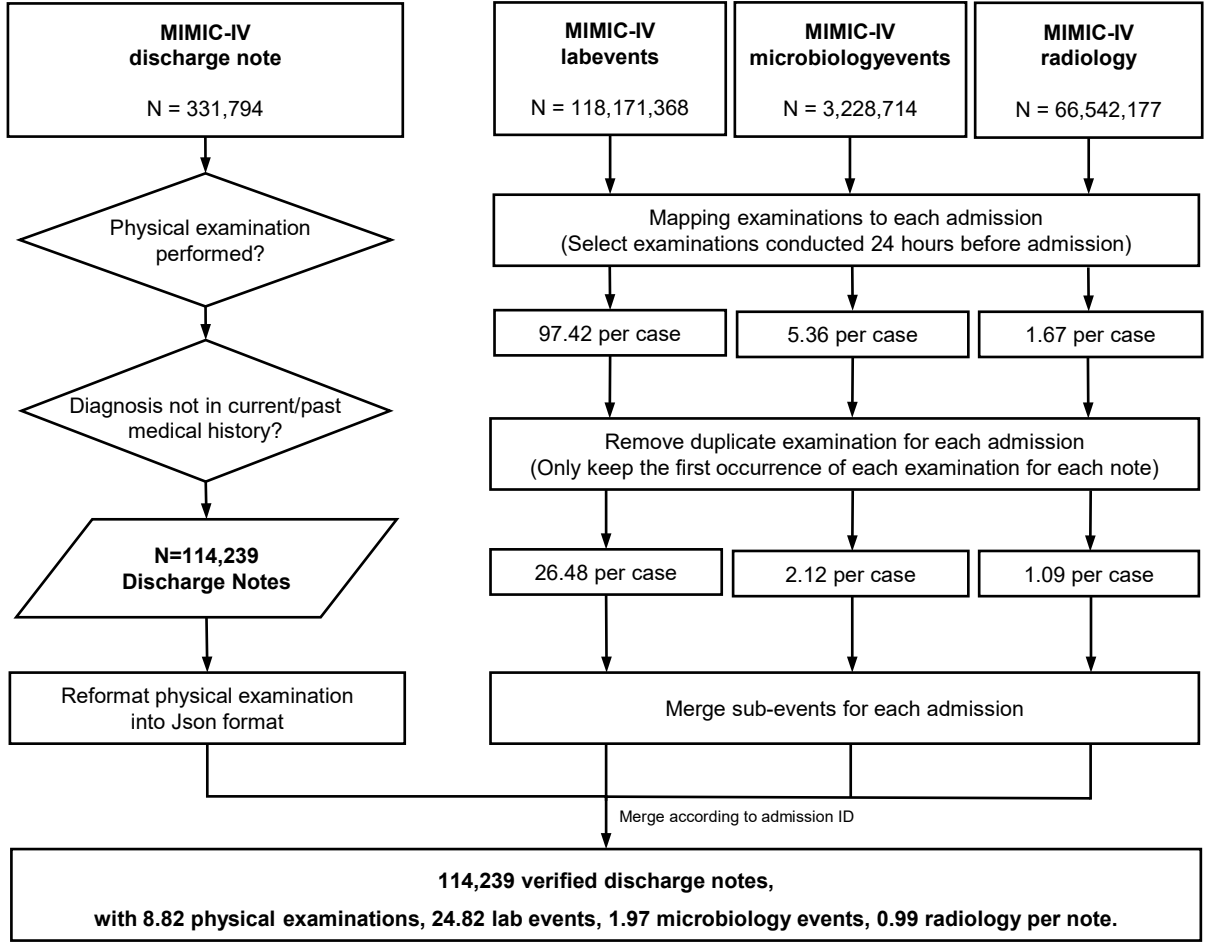
5.1 DiagGym Training

In this section, we present the detailed training of DiagGym.

Data Construction

As shown in Figure 5a, to train DiagGym, we constructed a dataset of patient EHRs derived from MIMIC-IV.

a. DiagGym Data Construction



b. DiagGym Training

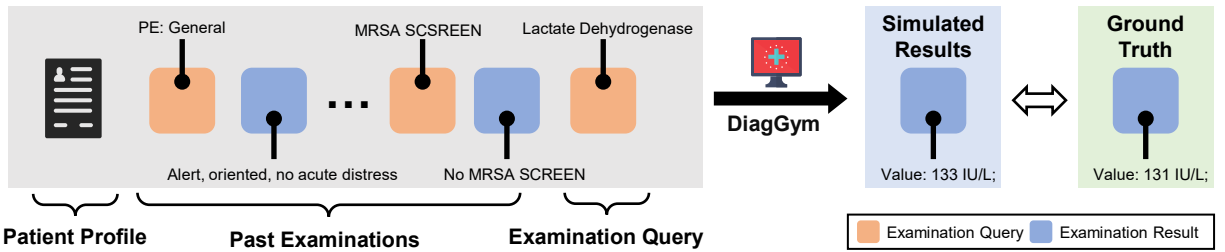


Figure 5 | Overview of DiagGym data construction and training pipeline. **a** shows the process for constructing the DiagGym training dataset. **b** illustrates the pipeline for DiagGym training.

Each patient’s EHR was reorganized into two components: (i) patient profile, (ii) time-ordered examination set. The pipeline was based on MIMIC-CDM [14], but extended to cover a broader range of diseases.

We first process the MIMIC-IV discharge notes. Leveraging their structured format, we applied heuristic string matching to extract the patient profile. Specifically, a patient profile was composed of content under the headings ‘physical examination’, ‘chief complaint’, ‘current medical history’, ‘past medical history’, ‘social history’, ‘family history’, and, most critically, the final diagnosis listed under the ‘discharge diagnosis’ heading.

Next, we apply a two-step filtering process: (i) cases without physical examination records are excluded; (ii)

DeepSeek-V3, instructed with [prompt 1](#), is used to remove cases where the discharge diagnosis appeared in either the past medical history or the current medical history. Such cases often involve transfers with established diagnoses and typically lack diagnostically relevant examinations.

We then construct a time-ordered examination set for each patient, with each examination in the set comprising the queried examination item and its corresponding results. First, the previously extracted physical examination text from discharge notes is reformatted into a structured tabular format using DeepSeek-v3 with [prompt 2](#), ensuring consistency with other MIMIC-IV examination records. Following MIMIC-CDM, this physical examination is designated as the initial test in the examination set. Next, we append laboratory results, microbiological examinations, and radiological records conducted within one day prior to admission. The one-day time interval is selected because examinations performed earlier generally have limited diagnostic relevance. Laboratory data is obtained from `labevents.csv`, microbiology data from `microbiologyevents.csv`, and radiology from `radiology.csv`.

For laboratory and microbiology data, we use the original structured records but standardized examination item names using the MIMIC-CDM mapping table to group the linked items (e.g., “red blood cell count” under the broader “complete blood count”). Radiology entries are supplemented with missing names extracted by string-matching from the EXAMINATION section of reports. All examination entries in MIMIC-IV contain timestamps, enabling accurate chronological ordering. For repeated pre-admission examinations, only the earliest one is retained, similar to [\[14\]](#).

Finally, we split the restructured EHR dataset into training and testing sets. The resulting dataset consists of 118,478 patient EHRs, with each case containing the patient profile and a time-ordered set of examinations. Of these, 114,239 EHRs are used for the world model training, where the model is tasked with autoregressively reconstructing the examination results recorded in the examination set. These training cases span 4,897 distinct diseases. On average, each training patient underwent 29 examinations, including 26 laboratory tests, 2 microbiological tests, and 1 radiological test. The remaining 4,239 cases are reserved for evaluation, covering 863 distinct diseases. However, given the high cost of evaluation—due to the use of various commercial models—we adopted a disease-wise sampling strategy. Specifically, we selected one representative case for each disease, resulting in a balanced test set of 863 cases, with each case corresponding to a unique disease.

Training Details

Leveraging the constructed data, we train a diagnostics world model, DiagGym, with text generation loss, as introduced in the former [Section 2](#). We frame its training as an auto-regressive text generation task, straightforwardly viewing all examination results as free text, regardless of whether they are numerical or textual. The loss function, a standard token-wise auto-regressive objective inspired by GPT-series models [\[28\]](#), minimizes the negative log-likelihood of the ground-truth examination result \hat{e}_t tokens:

$$\mathcal{L}_{\text{recon}} = - \sum_{t=1}^T \sum_i \log \Phi_{\text{env}}(\hat{e}_t^i \mid a_{t+1}, E_t, \mathcal{B}), \quad (9)$$

where \hat{e}_t^i denotes the i -th token of the examination result e_t and E_t denotes the former examination records in the examination set.

Implementation Details. For our experiments, we initialized Φ_{env} from Qwen2.5-Instruct-7B. Training was performed on eight NVIDIA A100 GPUs using the `Transformers`^{[1](#)} library with DeepSpeed ZeRO Stage 2 for efficient distributed optimization. Models were trained for 15 epochs, with convergence achieved within this period. The learning rate was 4×10^{-5} , and the maximum input length was 8,192 tokens.

5.2 DiagAgent Training

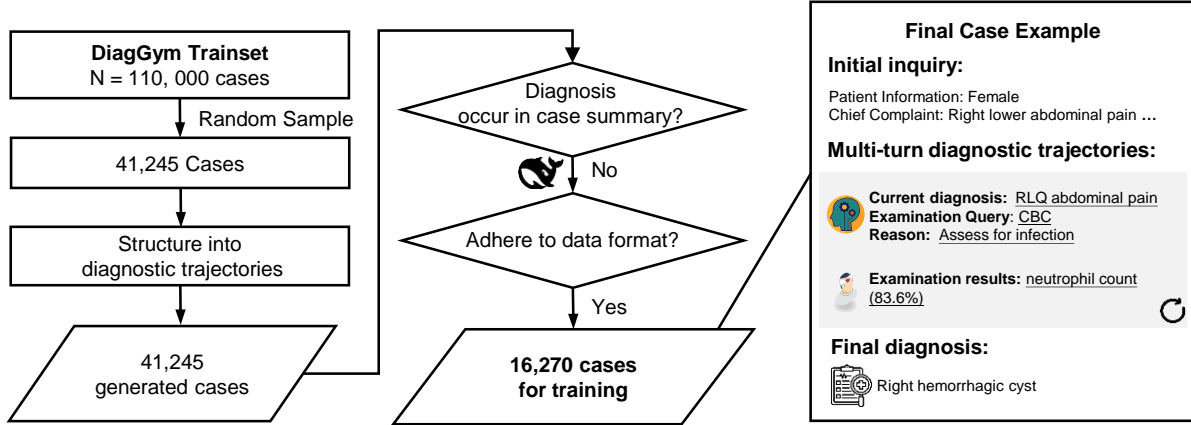
In this section, we provide a detailed procedure for DiagAgent training.

Data Construction

The diagnostic agent Φ_{diag} is trained using a reformatted version of the DiagGym training dataset, organized

¹<https://github.com/huggingface/transformers>

a. DiagAgent Data Construction



b. DiagAgent Training

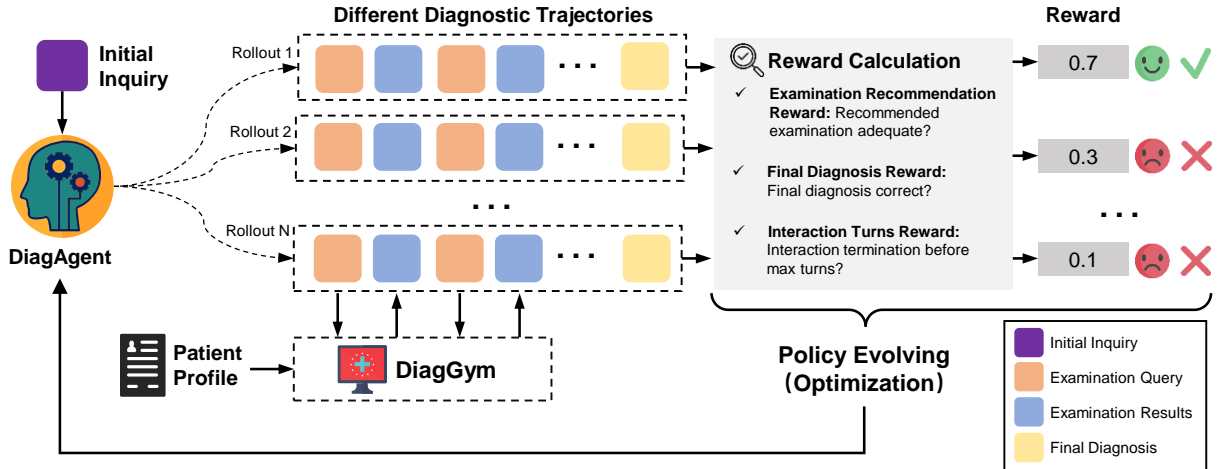


Figure 6 | Overview of DiaAgent data construction and training pipeline. **a** shows the data construction process for the DiagAgent. **b** illustrates the training pipeline for the diagnostic agent, where the agent interacts with the virtual clinical environment, DiagGym, explores different diagnostic trajectories, and evolves its policy based on reward scores.

into a multi-turn diagnostic trajectory format, as shown in Figure 6a.

We use DeepSeek-v3, guided by prompt 3, to generate three tightly connected elements for each case based on the existing restructured EHR dataset and the original discharge notes:

- **Initial inquiry:** A refined, structured summary of the patient’s medical history before admission, covering the chief complaint, history of present illness, past medical history, family history, and other relevant information. While this is similar to the previously constructed patient profile, critically, it does not include the final diagnosis information. The inquiry is further refined by LLMs to align with real inquiry formats. This serves as the starting point for the dialogue.
- **Referenced multi-turn diagnostic trajectory:** A step-by-step diagnostic trajectory is reformed from the time-ordered examination chain, with DeepSeek-v3 prioritizing the most informative tests related to the final diagnosis and omitting non-essential routine ones. Each step in the trajectory consists of: (1) Current Preliminary Diagnosis: An initial hypothesis based on prior data. (2) Next Recommended Examination with Rationale: A suggested test accompanied by a detailed explanation for its necessity. (3) Corresponding Test Results: The outcome of the recommended examination. The first two components are considered the agent’s response, while the third represents the clinical environment’s feedback — effectively, the user’s input in the multi-turn dialogue. The trajectory concludes with the

final diagnostic decision at the end-turn. Although the preliminary diagnosis and the rationale for recommending examinations are generated with the assistance of LLMs, the order of examination items and their corresponding results are directly extracted from real EHRs. Therefore, this trajectory is viewed as a referenced diagnostic pathway, grounded in the recommendations of real physicians.

- **Final diagnosis:** The final diagnostic decision for the case. Since the original final diagnosis recorded in the patient profile may sometimes include multiple conditions, the LLM is prompted to construct a self-contained process focused on a single primary condition. This process also includes selecting relevant examinations from the entire chain of examinations in the previous multi-turn diagnostic trajectory construction.

All three components are generated in a single pass to maintain contextual consistency.

To ensure data quality and prevent leakage of the final diagnosis into the initial inquiry, we applied a two-stage filtering process using **prompt 4**. This step removes instances where the model may have inadvertently introduced the final diagnosis into earlier parts of the text due to hallucinations.

Finally, we converted the multi-turn diagnostic trajectories into a structured dialogue format following former LLM multi-turn datasets [29]. In subsequent turns, each step in the trajectory is structured as an assistant message that includes the preliminary diagnosis and the next recommended examination with its rationale, along with a user message that provides the results of the recommended examinations. The last assistant turn contains only the final diagnosis and its rationale. The resulting dataset comprised 15,324 interactive diagnostic trajectories used for training.

Training Details

As introduced in Section 2, we train DiagAgent Φ_{diag} with end-to-end multi-turn RL, with a classical two-stage paradigm [30]: an initial cold-start phase and a main RL phase.

Cold start. This phase mirrors standard instruction tuning [31]. The model is optimized with an auto-regressive text generation loss computed only over tokens labeled as **assistant** response in the dialogues:

$$\mathcal{L}_{\text{cold}} = - \sum_{y_i \in \text{assistant}} \log \Phi_{\text{diag}}(y_i | y_{\leq i-1}), \quad (10)$$

where y_i is a token, $y_{\leq i-1}$ is the preceding context, and the loss is restricted to assistant response tokens. The goal of cold start is to initialize the LLM to produce well-formatted, contextually appropriate responses before interacting with the environment with RL. For this stage, we use 1,000 manually selected high-quality cases from the training set, free from formatting issues or diagnostic reasoning errors.

Reinforcement learning. After cold start, we optimize Φ_{diag} with the GRPO algorithm [30] over the full training set. At rollout start, the agent receives the initial inquiry as the initial state, $s_0 = \mathcal{I}$ and iteratively interacts with the virtual environment Φ_{env} until it decides sufficient information has been gathered for a final diagnosis (Section 2). The policy is trained to maximize the following reward:

$$R = \lambda_1 r_{\text{diag}} + \lambda_2 r_{\text{exam}} + \lambda_3 r_{\text{turn}}, \quad (11)$$

with $\lambda_1, \lambda_2, \lambda_3$ as hyper-parameter weights. The final diagnosis reward evaluates the accuracy of the predicted diagnosis results, formulated as:

$$r_{\text{diag}} = \begin{cases} 1, & \text{if } \hat{d} = d \\ 0, & \text{otherwise} \end{cases}, \quad (12)$$

where \hat{d} is the predicted disease. We adopt Qwen2.5-72B with **prompt 12** to measure semantic equivalence of \hat{d} and d . The examination recommendation reward measures the alignment between the agent’s recommended examinations \hat{E} and the reference set E from the curated multi-turn trajectory based on real EHRs. We adopt the F1 score to measure their similarity, formulated as:

$$r_{\text{exam}} = \text{F1}(\hat{E}, E) = 2 \cdot \frac{|\hat{E} \cap E|}{|\hat{E}| + |E|}. \quad (13)$$

We calculate the union set of the two by instructing Qwen2.5-72B with [prompt 10](#) and [prompt 11](#) to search through the \hat{E} and E respectively. The last interaction turn penalty reward penalizes excessive rounds of dialogue without termination. To prevent unnecessarily long dialogues, we impose a maximum number of iterative turns, T_{max} . If the model cannot finish the diagnosis within this limitation, it will achieve a lower reward as:

$$r_{\text{turn}} = \begin{cases} 0.1, & \text{if } T \leq T_{max} \\ 0, & \text{otherwise} \end{cases}, \quad (14)$$

where T denotes the total number of turns used in a certain rollout.

Through iteratively optimizing this reward, Φ_{diag} learns to manage accurate, efficient diagnostic trajectory reasoning while minimizing unnecessary examination steps.

Implementation details. In our experiments, we selected Qwen2.5-Instruct-7B, Qwen2.5-Instruct-14B, and Llama3.1-Instruct-8B as initialization models for further training. The same training strategy was applied to all models. For cold start settings, we utilize the **Transformers**² framework and employed DeepSpeed ZeRO Stage 2 for efficient multi-GPU training. All models share the same hyperparameters: a maximum sequence length of 8192 tokens, a learning rate of 1×10^{-5} , and training on eight NVIDIA A100 GPUs. Training is conducted for three epochs, after which convergence was observed for each model.

For the reinforcement learning setting, we modify **Verl**³ to enable interactive training. During training, the **DiagGym** was deployed on two nodes using **vLLM**⁴. The Qwen2.5-Instruct-72B model serves as the judge and is deployed on a separate node. The weighted hyperparameters $\lambda_1, \lambda_2, \lambda_3$ are set as 1, 0.5, 1, respectively, and the maximum interactive turns are set to 12. RL training is performed across four nodes in total, with each node equipped with eight NVIDIA A100 GPUs. We set the training batch size to 512, the maximum response length to 8192 tokens, the learning rate to 1×10^{-6} , and the rollout number to 5. Each model is trained for 200 steps, and convergence is observed within this training regime.

5.3 DiagBench Construction

In this section, we describe the construction of **DiagBench**, a benchmark for evaluating multi-turn diagnostic interaction trajectories. The process involved two main stages: initial dataset curation and the development of a rubric-based evaluation framework.

For data curation, the initial dataset was generated using the same pipeline as our training data (Figure 6a), yielding 957 candidate cases. To ensure clinical validity, each case underwent a rigorous review by a physician. The reviewer was provided with the patient’s initial inquiry and the complete reference diagnostic trajectory. They were tasked with evaluating: (1) the clinical appropriateness of each analysis and recommended examination at every step, and (2) the overall plausibility of the case. This curation process resulted in a final, validated set of 750 cases. This set is used for the automated evaluation of examination recommendation efficacy and diagnostic accuracy.

In addition to automated metrics, and inspired by HealthBench [22], we develop a **rubric-based framework** to comprehensively assess the quality of the diagnostic inference process.

First, we randomly sample 100 clinical cases. Two physicians are independently provided with each case’s initial patient inquiry and complete diagnostic trajectory. Each physician then independently authors a set of process-oriented rubrics they consider critical for evaluating the quality of the diagnostic journey. The rubric design explicitly prioritizes the reasoning process—including history taking, hypothesis generation, and test ordering—rather than focusing solely on the accuracy of the final diagnosis.

Subsequently, a third physician conducts a secondary review of the rubrics. This reviewer’s role is to ensure the framework holistically covers the end-to-end diagnostic pipeline and that individual rubrics are coherent and well-defined. During this screening, one case is excluded due to insufficient rubric coverage, yielding a final set of 99 cases for the rubric-based evaluation.

Finally, a fourth physician is tasked with assigning an importance weight to each rubric on a scale from 0

²<https://github.com/huggingface/transformers>

³<https://github.com/volcengine/verl>

⁴<https://github.com/vllm-project/vllm>

to 10 (where 10 signifies an essential, non-negotiable criterion and 0 indicates a non-informative one). As a result, rubrics targeting high-impact clinical steps, such as appropriate test ordering and effective diagnostic narrowing, receive higher weights. Conversely, ancillary actions, like scheduling a follow-up without a clear clinical rationale, are assigned lower weights. Illustrative examples of this rubric-based evaluation are provided in Supplementary Figure 3 and Supplementary Figure 4.

5.4 Baselines

Here, we introduce the baseline LLMs involved in our experiments:

- **Qwen2.5** [32] and **Qwen3** [33] are series of high-performance open-source language models developed by the Qwen team, available in variants ranging from 0.5 to 72 billion parameters. In this paper, we use **Qwen2.5-7B-Instruct** and **Qwen2.5-14B-Instruct** for training, and deploy **Qwen2.5-72B-Instruct** and **Qwen3-235B-A22B** locally for inference.
- **Llama3.1** and **Llama3.3** [34] are series of language models developed by Meta AI and are among the most popular open-source large language models. In this paper, we utilize **Llama-3.1-8B-Instruct** for training and deploy **Llama-3.3-70B-Instruct** locally for inference.
- **OpenBioLLM** [35] is an advanced open-source language model specifically designed for the biomedical domain, developed based on Llama3. In this paper, we utilize **Llama3-OpenBioLLM-70B** and deploy it locally.
- **Baichuan-M1** [36] is an advanced open-source medical language model developed by Baichuan Intelligence. It is the first language model in the industry designed and developed from scratch specifically for the medical field, demonstrating strong performance in medical applications. We utilize **Baichuan-M1-14B-Instruct** and deploy it locally.
- **DeepSeek-V3** [37] is one of the most powerful open-source language models, developed by DeepSeek, with 671 billion parameters. In this paper, we use **DeepSeek-V3-0324** and deploy it locally.
- **MedGemma** [38] is a recent medical LLM developed by Google as a variant of the Gemma3 collection. Based on Gemma3, this model is specially optimized for the medical field and possesses multi-modal capabilities. In this paper, we utilize **medgemma-27b-text-it** and deploy it locally.
- **GPT-OSS** [39] is an open-source large language model released by OpenAI, with strong reasoning ability. In this paper, we use the **gpt-oss-120b** and deploy it locally.
- **GPT-4o** [40] is one of the most commonly used close-sourced language model developed by OpenAI. It is good at handling most daily tasks. In this paper, we utilize **gpt-4o-2024-08-06** via API.
- **Claude-4** [41] is a high-performance large language model developed by Anthropic, optimized for coding and reasoning. In this paper, we use the **claude-sonnet-4** via API.

Next, we will introduce the agentic systems involved in our experiments. Note that, in our experiment, all the agentic systems use DeepSeek-V3 as its base model.

- **MedAgents**[42] is a multidisciplinary collaboration framework that requires large language models to assume the roles of medical experts from various specialties. The framework aggregates the experts' opinions and summarizes them into a final report. This approach reveals the model's knowledge across different domains and broadens its reasoning capabilities. In our implementation, we use DeepSeek-V3 as the base model.
- **MDAgents**[43] is a framework designed to automatically assign collaboration structures. Different from MedAgents, MDAgents effectively assesses medical complexity and adapt to varying tasks accordingly. In our implementation, we use DeepSeek-V3 as the base model.

In the **DiagGym evaluation**, we consider prompting the open-source models as EHR wold model baselines to simulate the examination results, including DeepSeek-V3-671B, MedGemma-27B, Qwen2.5-7B, and Qwen2.5-72B. For instance-wise metrics, these models adopt **prompt 7** to generate examination results for comparison. For examination-wise metrics, which measures the distribution of given examinations, the models use **prompt 8** for numerical examinations and **prompt 9** for free-text examinations.

In the **DiagAgent evaluation**, we include all the aforementioned basic LLMs and more complex agentic systems for comparison. These models or agentic systems are instructed with **prompt 5** to complete the diagnostic trajectory. For end-to-end evaluation on simulated cases, the models are required to decide whether to request an additional examination or to make a final diagnosis. For single-turn evaluation on real cases, the models are directly instructed to either request an examination or make a final diagnosis. In this setting, we append an explicit prompt - “Next step you should query examination” or “Next step you should make final diagnosis” - to the end of the input in each turn.

6 Data Availability

The data source for this work is MIMIC-IV. Due to licensing restrictions, we are unable to directly open-source the dataset. However, we are actively communicating with the relevant parties regarding the possibility of making the dataset publicly available on <https://physionet.org/>.

7 Code Availability

All source codes of this paper have been released in <https://github.com/MAGIC-AI4Med/DiagGym>.

References

- [1] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, pages 1–8, 2025.
- [2] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- [3] Sarah Sandmann, Stefan Heggelmann, Michael Fujarski, Lucas Bickmann, Benjamin Wild, Roland Eils, and Julian Varghese. Benchmark evaluation of deepseek large language models in clinical decision-making. *Nature Medicine*, pages 1–1, 2025.
- [4] Daniel McDuff, Mike Schaekermann, Tao Tu, Anil Palepu, Amy Wang, Jake Garrison, Karan Singhal, Yash Sharma, Shekoofeh Azizi, Kavita Kulkarni, et al. Towards accurate differential diagnosis with large language models. *Nature*, pages 1–7, 2025.
- [5] Farieda Gaber, Maqsood Shaik, Fabio Allega, Agnes Julia Bilecz, Felix Busch, Kelsey Goon, Vedran Franke, and Altuna Akalin. Evaluating large language model workflows in clinical decision support for triage and referral and diagnosis. *npj Digital Medicine*, 8(1):263, 2025.
- [6] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Hel-lyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [7] 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. *arXiv preprint arXiv:2501.12948*, 2025.
- [8] Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. Pmc-llama: toward building open-source language models for medicine. *Journal of the American Medical Informatics Association*, 31(9):1833–1843, 2024.
- [9] Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Towards building multilingual language model for medicine. *Nature Communications*, 15(1):8384, 2024.
- [10] Pengcheng Qiu, Chaoyi Wu, Shuyu Liu, Weike Zhao, Zhuoxia Chen, Hongfei Gu, Chuanjin Peng, Ya Zhang, Yanfeng Wang, and Weidi Xie. Quantifying the reasoning abilities of llms on real-world clinical cases. *arXiv preprint arXiv:2503.04691*, 2025.
- [11] Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*, 2023.
- [12] Bingning Wang, Haizhou Zhao, Huozhi Zhou, Liang Song, Mingyu Xu, Wei Cheng, Xiangrong Zeng, Yupeng Zhang, Yuqi Huo, Zecheng Wang, et al. Baichuan-m1: Pushing the medical capability of large language models. *arXiv preprint arXiv:2502.12671*, 2025.
- [13] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407, 2024.
- [14] Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer, Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, et al. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nature medicine*, 30(9):2613–2622, 2024.

- [15] Yusheng Liao, Yutong Meng, Yuhao Wang, Hongcheng Liu, Yanfeng Wang, and Yu Wang. Automatic interactive evaluation for large language models with state aware patient simulator. *arXiv preprint arXiv:2403.08495*, 2024.
- [16] Shreya Johri, Jaehwan Jeong, Benjamin A Tran, Daniel I Schlessinger, Shannon Wongvibulsin, Leandra A Barnes, Hong-Yu Zhou, Zhuo Ran Cai, Eliezer M Van Allen, David Kim, et al. An evaluation framework for clinical use of large language models in patient interaction tasks. *Nature medicine*, 31(1):77–86, 2025.
- [17] Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor. Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments. *arXiv preprint arXiv:2405.07960*, 2024.
- [18] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are gans created equal? a large-scale study. *Advances in neural information processing systems*, 31, 2018.
- [19] François Remy, Kris Demuynck, and Thomas Demeester. Biolord-2023: semantic textual representations fusing large language models and clinical knowledge graph insights. *Journal of the American Medical Informatics Association*, 31(9):1844–1855, 2024.
- [20] Yu Yu, Weibin Zhang, and Yun Deng. Frechet inception distance (fid) for evaluating gans. *China University of Mining Technology Beijing Graduate School*, 3(11), 2021.
- [21] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10743–10752, 2021.
- [22] Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, et al. Healthbench: Evaluating large language models towards improved human health. *arXiv preprint arXiv:2505.08775*, 2025.
- [23] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023.
- [24] Pawel Renc, Michal K Grzeszczyk, Nassim Oufattole, Deirdre Goode, Yugang Jia, Szymon Bieganski, Matthew McDermott, Jaroslaw Was, Anthony E Samir, Jonathan W Cunningham, et al. Foundation model of electronic medical records for adaptive risk estimation. *arXiv preprint arXiv:2502.06124*, 2025.
- [25] Tao Tu, Mike Schaeckermann, Anil Palepu, Khaled Saab, Jan Freyberg, Ryutaro Tanno, Amy Wang, Brenna Li, Mohamed Amin, Yong Cheng, et al. Towards conversational diagnostic artificial intelligence. *Nature*, pages 1–9, 2025.
- [26] Junkai Li, Yunghwei Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li, Ya-Qin Zhang, Weizhi Ma, et al. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957*, 2024.
- [27] Harsha Nori, Mayank Daswani, Christopher Kelly, Scott Lundberg, Marco Tulio Ribeiro, Marc Wilson, Xiaoxuan Liu, Viknesh Sounderajah, Jonathan Carlson, Matthew P Lungren, et al. Sequential diagnosis with language models. *arXiv preprint arXiv:2506.22405*, 2025.
- [28] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [29] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

- [30] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [31] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [32] Qwen Team. Qwen2.5: A party of foundation models, September 2024.
- [33] 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. *arXiv preprint arXiv:2505.09388*, 2025.
- [34] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [35] Malaikannan Sankarasubbu Ankit Pal. Openbiollms: Advancing open-source large language models for healthcare and life sciences. <https://huggingface.co/aaditya/OpenBioLLM-Llama3-70B>, 2024.
- [36] Huozhi Zhou Bingning Wang, Haizhou Zhao et al. Baichuan-m1: Pushing the medical capability of large language models. *arXiv preprint arXiv:2502.12671*, 2025.
- [37] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [38] Andrew Sellergrén, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cian Hughes, Charles Lau, et al. Medgemma technical report. *arXiv preprint arXiv:2507.05201*, 2025.
- [39] OpenAI. Introducing gpt-oss. <https://openai.com/index/introducing-gpt-oss/>. Accessed: 2025-08-08.
- [40] OpenAI. Hello gpt-4o, 2025. Accessed: 2025-02-27.
- [41] Anthropic. Introducing claude 4. <https://www.anthropic.com/news/claude-4>. Accessed: 2025-08-08.
- [42] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. Medagents: Large language models as collaborators for zero-shot medical reasoning. *arXiv preprint arXiv:2311.10537*, 2023.
- [43] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee, Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won Park. Mdagents: An adaptive collaboration of llms for medical decision-making, 2024.

8 Acknowledgments

This work is supported by the National Key R&D Program of China (No. 2022ZD0160702), and the Scientific Research Innovation Capability Support Project for Young Faculty (ZY-GXQNJSKYCXNLZCXM-I22).

9 Author Contributions

All listed authors clearly meet the ICMJE 4 criteria. P.Q., C.W., and J.L. contribute equally to this work. Y.Z. and W.X. are the corresponding authors. Specifically, P.Q., C.W., J.L., Q.Z., Y.L., H.W., Y.Y., Q.F., S.Z., J.W., J.G., Y.W., Y.Z., and W.X. all make contributions to the conception or design of the work, and P.Q., C.W., and J.L. further perform acquisition, analysis, or interpretation of data for the work. In writing, P.Q., C.W., and J.L. draft the work. Q.Z., Y.L., H.W., Y.Y., Q.F., S.Z., J.W., J.G., Y.W., Y.Z., and W.X. review it critically for important intellectual content. All authors approve of the version to be published and agree to be accountable for all aspects of the work to ensure that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

10 Supplementary

10.1 Evaluation Metrics

In this section, we provide a detailed calculation format for our adopted evaluation metrics.

10.1.1 DiagGym Evaluation

We evaluate the DiagGym (Φ_{env}) from two perspectives, *i.e.*, instance-wise and examination-wise.

Instance-wise Metrics. These metrics all employ LLMs as judges. Two types of evaluations are conducted: step-wise similarity assessment and full-chain consistency evaluation.

- **Step-wise Similarity.** This metric evaluates the similarity between the model-predicted examination results and the ground truth at each step. During this evaluation, the simulator generates the current examination results conditioned on the queried examination names, the patient case summary, and all historical ground truth examination names and results. At each step, we apply Prompt 14 to compute similarity scores (ranging from 0 to 5) between the model-generated and ground truth examination results.
- **Full-chain Consistency.** This metric assesses whether the generated examination chain is internally consistent and aligns with the patient profile. Given the generated chain, we use Prompt 15 to assess whether the sequence of generated examination results maintains logical and clinical consistency. The evaluation final score is binary 1/0, indicating yes/no.

Examination-wise Metrics For the examination-wise metrics, we primarily assess the statistical distribution quality for both generative numerical and free-text results, covering fidelity and diversity.

- **Numerical Fidelity & Diversity.** For numerical examination results, we utilize the 1-Wasserstein distance to measure the generative distribution fidelity, where shorter distances indicate closer alignment. Formally, considering a certain numerical examination, its real distribution is characterized by the mean μ and standard deviation σ . Similarly, the generative distribution is characterized by $\hat{\mu}$ and $\hat{\sigma}$. Denoting a generative examination value on a certain test case as x_i , it can be normalized as $z_i = (x_i - \mu)/\sigma$ and similarly, for the ground-truth results, we have $\hat{z}_i = (\hat{x}_i - \hat{\mu})/\hat{\sigma}$. The 1-Wasserstein distance is formulated as:

$$W_1(Z, \hat{Z}) = \inf_{\gamma \in \Gamma(Z, \hat{Z})} \int \|z_i - \hat{z}_i\| d\gamma(z_i, \hat{z}_i), \quad (15)$$

where Z and \hat{Z} denote the generative and real distributions, and $\Gamma(Z, \hat{Z})$ is the joint distribution. Then for diversity, we adopt normalized variance to evaluate its distribution diversity, with higher variance reflecting greater diversity. Formally, following the former notation, the normalized variance is defined as:

$$\sigma^2(X) = \frac{\text{Var}(X)}{\mu^2}, \quad (16)$$

where $\text{Var}(\cdot)$ is the distribution variance. Notably, for simplicity, all the above numerical metric formulations assume that the examination item contains only a single value item. In practice, some examinations may consist of multiple value items, such as a “Complete Blood Count”, which may include multiple numerical values, and in these cases, the metrics are calculated by averaging the scores across all value items. Furthermore, to ensure consistency, all computations for the same value items are standardized to unified value units.

- **Free-text Fidelity & Diversity.** For free-text results, we first encode the text into feature embeddings using BioLORD [19], a biomedical text encoding model. Inspired by metrics commonly used in image generation, we then calculate the Fréchet Inception Distance (FID)[20] in the embedding space to assess the fidelity. Lower FID values indicate better alignment with the ground truth. Specifically, considering a certain free-text-related examination, such as a chest CT examination, let the set of generative free-text embeddings be $F = \{f_1, f_2, \dots, f_N\}$, where N denotes the total number of test cases.

The corresponding ground truth text embeddings are denoted as $\hat{F} = \{\hat{f}_1, \hat{f}_2, \dots, \hat{f}_N\}$. The FID score can then be calculated as:

$$\text{FID}(F, \hat{F}) = \|\mu_F - \mu_{\hat{F}}\|_2^2 + \text{Tr} \left(\Sigma_F + \Sigma_{\hat{F}} - 2(\Sigma_F \Sigma_{\hat{F}})^{\frac{1}{2}} \right), \quad (17)$$

where μ_F and Σ_F represent the mean and covariance of the generative embeddings F , and $\mu_{\hat{F}}$ and $\Sigma_{\hat{F}}$ are for the ground truth embeddings \hat{F} . The term $\|\mu_F - \mu_{\hat{F}}\|_2^2$ quantifies the difference between the means, while the trace term measures the distance between the covariance matrices. To evaluate diversity, inspired by the Intra-LPIPS metric [21] used in image generation, we propose using inter-case cosine similarity on the entire set of generated text embeddings. This metric reflects how well the embeddings distinguish from one another, defined as:

$$\text{Intra-LPIPS}(F) = 1 - \frac{2}{N(N-1)} \sum_{i < j} \cos(f_i, f_j), \quad (18)$$

where $\cos(\cdot, \cdot)$ represents the cosine similarity function. A higher Intra-LPIPS score indicates that the generated free texts are more diverse in comparison to one another.

10.1.2 DiagAgent Evaluation

We evaluate DiagAgent’s performance across three key perspectives: final diagnosis accuracy, examination recommendation efficacy, and rubric-based diagnostic trajectory assessment.

Final Diagnosis Accuracy. For measure the final diagnosis, we adopt the straightforward **Accuracy** metric here. Similarly, considering diseases may have synonyms, we instruct GPT-4o utilizing Prompt 12 to compare the model’s diagnostic output with the reference standard.

Examination Recommendation Metrics In the end-to-end evaluation setting, we compare the predicted examination list against the referenced list and adopt **Precision**, **Recall**, and **F1** scores. Considering that the same examination may be expressed in synonyms, we instruct GPT-4o employing Prompt 10 to count the number of examination names generated by the model that are covered by the true examination list, and Prompt 11 to count the number of ground truth examination names that are present in the model’s output. Precision and recall are then computed accordingly, followed by the calculation of the F1 score. In the single-turn evaluation setting, we adopt the hit ratio metric, counting whether the recommended examination item appears in the referenced following referenced list. We adopt GPT-4o with Prompt 6 to determine whether this query appears among the examinations actually undergone by the patient.

Weighted Rubric Score To evaluate the clinical integrity of the multi-turn diagnostic interaction, we use the Weighted Rubric Score. This metric goes beyond the final outcome by qualitatively assessing the entire diagnostic trajectory against physician-authored, process-oriented rubrics. The score is calculated as the weighted proportion of satisfied rubrics, with weights reflecting the clinical significance of each step. The final score \bar{s} is the average across all 99 cases, calculated as:

$$\bar{s} = \frac{1}{|C|} \sum_{c \in C} \frac{\sum_{r \in A_c} w_r}{\sum_{r \in R_c} w_r}.$$

The GPT-4o judge model determines rubric satisfaction for each case (Prompt 13).

10.2 Supervised Finetuning

In this section, we introduce the multi-turn supervised fine-tuning (SFT) paradigm. In multi-turn SFT, each training sample consists of a dialogue history and the next expected response. Formally, let $\mathcal{D}_{\text{multi}} = \{(h_t^{(j)}, y_t^{(j)})\}_{j=1}^M$, where $h_t^{(j)} = (u_1^{(j)}, r_1^{(j)}, \dots, u_{t-1}^{(j)}, r_{t-1}^{(j)}, u_t^{(j)})$ denotes the dialogue history up to turn t (with u as user inputs and r as model responses), and $y_t^{(j)}$ is the supervised response. The loss becomes:

Supplementary Table 1 | Selected Laboratory Events and Their Corresponding Sub-events

| Event | Sub-events |
|-------------------------------|---|
| Complete Blood Count | MCH; White Blood Cells; Absolute Basophil Count; Absolute Monocyte Count; Absolute Eosinophil Count; Monocytes; Platelet Count; Hemoglobin; Hematocrit; Neutrophils; Absolute Neutrophil Count; RDW; Basophils; Eosinophils; Red Blood Cells; RDW-SD; Lymphocytes |
| Mean Corpuscular Volume | MCV |
| Liver Function Test | PT; Alanine Aminotransferase (ALT); Asparate Aminotransferase (AST); Albumin; Alkaline Phosphatase |
| Anion Gap | Anion Gap |
| Comprehensive Metabolic Panel | Bicarbonate |
| Total Bilirubin | Bilirubin, Total |
| Total Calcium | Calcium, Total |
| Kidney Function Tests | Urea Nitrogen; Glucose; Phosphate; Creatinine; Sodium; Potassium; Chloride |
| Lactate Dehydrogenase | Lactate Dehydrogenase (LD) |
| Magnesium | Magnesium |
| Coagulation Profile | PTT |
| Lactate | Lactate |
| Urine Analysis | RBC; pH; WBC |
| Total Calculated CO2 | Calculated Total CO2 |
| pCO2 | pCO2 |
| pH | pH |
| pO2 | pO2 |
| Lipase | Lipase |
| Creatine Kinase | Creatine Kinase (CK) |
| Creatine Kinase, MB Isoenzyme | Creatine Kinase, MB Isoenzyme |
| H | H |
| I | I |
| L | L |
| Thyroid Function Test | Thyroid Stimulating Hormone |

$$\mathcal{L}_{\text{multi-turn SFT}} = - \sum_{j=1}^M \log P_{\theta}(y_t^{(j)} | h_t^{(j)}) \quad (19)$$

10.3 Selected Examinations Collection

To ensure the reliability of the distributional analysis, we considered only examinations with relatively high occurrence frequencies. For numerical-type examinations, only those with more than 500 occurrences were retained, resulting in 11 examinations comprising a total of 24 subevents. The distributional metrics were computed based on these 24 subevents, as detailed in Table 1. For textual-type examinations, only radiology examinations with more than 100 occurrences were considered, yielding five types: “LIVER OR GALLBLADDER US (SINGLE ORGAN)”, “CT HEAD W/O CONTRAST”, “CT HEAD W/O CONTRAST Q111 CT HEAD”, “CHEST (PA AND LAT)”, and “CHEST (PORTABLE AP)”.

10.4 Case Collection

| Note ID 16988189-DS-3 | | | |
|----------------------------------|---|---|---|
| Patient Profile | <p>Sex: F</p> <p>Allergies: Ciprofloxacin / Erythromycin Base / Sulfa (Sulfonamide Antibiotics) / amoxicillin</p> <p>Chief Complaint S/p ERCP with brushing, biliary stent for CHD stricture</p> <p>History of Present Illness Ms. ____ is a ____ year old woman with history of breast cancer (s/p bilateral mastectomy, Adriamycin and taxol and radiation in ____), hypertension and non-ischemic cardiomyopathy presenting for monitoring post-ERCP in the setting of painless jaundice. Per Atrius records and my discussion with patient, Ms. ____ presented to her PCP ____ with 1 week of nausea and constipation with rare hard stools that were light in color. She was also noting a flare of her hemorrhoids as well as itching all over without new medications. She had labs sent at that time which showed elevated LFTs (AST 148, ALT 312, Alk phos 421, Tbili 7.6). She underwent a CT abdomen/pelvis which reportedly showed a duct tumor involving the common hepatic duct bifurcation concerning for Klatskin tumor (or hilar cholangiocarcinoma). Today, Ms. ____ presented for a planned procedure for evaluation of her mass and obstruction. She underwent ERCP with brushings and biliary stent placement for CHD stricture which she tolerated well. Upon arrival to the floor, she reports that she feels well without any symptoms. She did not take her blood pressure medications this morning. Her mouth feels dry but she denies light headedness, dizziness. No recent weight loss, fevers or night sweats. Her nausea has resolved since last week. She is having some chalky stools and dark urine and has noted yellowing of her skin.</p> <p>Past Medical History Past medical history includes hypertension, non-ischemic cardiomyopathy (diagnosed ____, viral vs. chemotherapy induced), left heart catheterization with normal coronaries, transthoracic echocardiogram ____: EF normal, and breast cancer (status post bilateral mastectomy, Adriamycin/Taxol, radiation therapy ____).</p> <p>Social History None</p> <p>Family History Father passed away ____ bone cancer. Otherwise non-contributory</p> <p>Final diagnosis CBD obstruction from common hepatic duct mass</p> | | |
| | Examination Chain | | |
| Examination Name | Prediction | Ground truth | Difference & Evaluation |
| General (Physical Exam) | Well appearing older woman, sitting upright in bed | Well appearing thin woman sitting up in bed in NAD, pleasant | No significant difference; both describe a well-appearing patient. |
| Eyes (Physical Exam) | EOMI, PERRL, anicteric | Icteric sclera, PERRL, EOMI | Prediction missed icterus, but GT confirms jaundice; this supports diagnosis (CBD obstruction/mass). |
| Gastrointestinal (Physical Exam) | non-distended, soft, non-tender | hypoactive bowel sounds, soft, nondistended, mild ttp in epigastrium without rebound | Mild tenderness in GT, but both exams are overall benign with no peritonitis; does not alter main diagnosis. |
| Skin (Physical Exam) | Yellow hue, multiple nevi, no rash, warm | No rashes noted, icteric skin | Both note jaundice; difference in description detail is not significant. |
| Liver Function Test | INR(PT): 1.1 (0.9-1.1); PT: 11.5 (9.4-12.5); ALT: 188 (0-40); Alk phos: 388 (35-105); AST: 118 (0-40) | INR(PT): 1.2 (0.9-1.1); PT: 12.5 (9.4-12.5); ALT: 342 (0-40); Alk phos: 300 (35-105); AST: 195 (0-40) | All abnormal, showing cholestasis/hepatocellular damage; GT values somewhat higher, but both confirm the same severe pattern. |
| Comprehensive Metabolic Panel | Bicarbonate: 23 (22-32) | Bicarbonate: 20 (22-32) abnormal | Mildly low in GT; not clinically important for diagnosis. |
| Total Bilirubin | 6.7 (0-1.5) abnormal | 4.3 (0-1.5) abnormal | Both markedly elevated; both support obstructive jaundice, difference not clinically significant. |
| Total Calcium | 9.6 (8.4-10.3) | 9.5 (8.4-10.3) | No significant difference. |
| Complete Blood Count | Hct: 35.8 (34-45); Hgb: 12.3 (11.2-15.7); MCH: 31.9 (26-32); MCHC: 34.4 (32-37); Plt: 235 (150-400); RDW: 14.9 (10.5-15.5); RBC: 3.85 (3.9-5.2) abnormal; WBC: 7.4 (4-10); RDW-SD: 51.1 (35.1-46.3) abnormal | Hct: 33.2 (34-45) abnormal; Hgb: 11.0 (11.2-15.7) abnormal; MCH: 33.0 (26-32) abnormal; MCHC: 33.1 (32-37); Plt: 271 (150-400); RDW: 13.2 (10.5-15.5); RBC: 3.33 (3.9-5.2) abnormal; WBC: 12.6 (4-10) abnormal; RDW-SD: 48.7 (35.1-46.3) abnormal | Both show mild anemia and some abnormal indices; GT WBC higher but not critical. Findings do not change main diagnosis. |
| Mean Corpuscular Volume | 100 (82-98) abnormal | 100 (82-98) abnormal | Both the same, macrocytic; not significant to main diagnosis. |
| Coagulation Profile | PTT: 31.1 (25-36.5) | PTT: 31.2 (25-36.5) | No significant difference. |

Supplementary Figure 1 | Example Case Study from DiagGym: Comparison of Predicted and Ground Truth Examinations. This case presents a single patient profile and final diagnosis, illustrating the step-wise evaluation setting. The core table compares the predicted examination results generated by DiagGym with the ground truth results in sequential order. The rightmost column analyzes key differences and discusses their clinical relevance in the diagnostic process.

| Note ID 13544756-DS-14 | | |
|---|---|--|
| Initial Inquiry | | <p>Sex: Male</p> <p>Chief Complaint: Abdominal pain.</p> <p>History of Present Illness: The patient reports several hours of abdominal pain that started periumbilical and migrated to the right lower quadrant. Associated symptoms include chills, nausea, diarrhea, and anorexia. The pain worsened progressively. Earlier in the day, he experienced diarrhea and had only soup and bread due to nausea.</p> <p>Past Medical History: Sleep apnea.</p> <p>Personal History: None.</p> <p>Family History: Hypertension, diabetes mellitus.</p> <p>Allergy History: Penicillins.</p> |
| Final Diagnosis | | Appendicitis |
| DiagAgent Interaction Messages | | |
| Turn | Role | Content |
| 1 | DiagAgent | Current diagnosis: The patient presents with abdominal pain localized to the right lower quadrant, accompanied by nausea, diarrhea, and anorexia. The differential diagnosis includes appendicitis, gastroenteritis, or other abdominal pathologies. Based on the patient's initial presentation, the following investigation(s) should be performed: Complete Blood Count. Reason: To evaluate for signs of infection or inflammation, such as leukocytosis, which is commonly seen in appendicitis. |
| | DiagGym | Absolute Lymphocyte Count: Value: 0.67; Numeric Value: 0.67; Units: K/uL; Reference Range: 1.2 - 3.7; Flag: abnormal Basophils: Value: 0.3; Numeric Value: 0.3; Units: %; Reference Range: 0.0 - 1.0 Eosinophils: Value: 0.0; Numeric Value: 0.0; Units: %; Reference Range: 1.0 - 7.0; Flag: abnormal Hematocrit: Value: 40.9; Numeric Value: 40.9; Units: %; Reference Range: 40.0 - 51.0 Hemoglobin: Value: 13.8; Numeric Value: 13.8; Units: g/dL; Reference Range: 13.7 - 17.5 Lymphocytes: Value: 4.8; Numeric Value: 4.8; Units: %; Reference Range: 19.0 - 53.0; Flag: abnormal MCH: Value: 30.0; Numeric Value: 30.0; Units: pg; Reference Range: 26.0 - 32.0 MCHC: Value: 33.7; Numeric Value: 33.7; Units: g/dL; Reference Range: 32.0 - 37.0 Monocytes: Value: 10.9; Numeric Value: 10.9; Units: %; Reference Range: 5.0 - 13.0 Neutrophils: Value: 83.6; Numeric Value: 83.6; Units: %; Reference Range: 34.0 - 71.0; Flag: abnormal Platelet Count: Value: 251; Numeric Value: 251.0; Units: K/uL; Reference Range: 150.0 - 400.0 RDW: Value: 12.9; Numeric Value: 12.9; Units: %; Reference Range: 10.5 - 15.5 |
| 2 | DiagAgent | Current diagnosis: The elevated neutrophil count (83.6%) suggests an inflammatory or infectious process, supporting the possibility of appendicitis. Based on the current findings, the following additional investigation(s) are needed: CT ABD AND PELVIS WITH CONTRAST. Reason: To visualize the appendix and surrounding structures to confirm the diagnosis of appendicitis or to identify other potential causes of abdominal pain. |
| | DiagGym | ... GASTROINTESTINAL: The stomach is unremarkable. Small bowel loops demonstrate normal caliber, wall thickness, and enhancement throughout. The colon and rectum are within normal limits. The appendix is fluid-filled and abnormally dilated, measuring up to 1.2 cm (series 2, image 62). There is associated fat stranding and fluid in the right lower quadrant. Findings are compatible with acute appendicitis. No drainable fluid collection or extraluminal gas is seen. ... |
| 3 | DiagAgent | The available information is sufficient to make a diagnosis. Diagnosis: Appendicitis. Reason: The patient's clinical presentation of right lower quadrant abdominal pain, nausea, and diarrhea, combined with the CT findings of a fluid-filled and dilated appendix with associated fat stranding and fluid in the right lower quadrant, confirm the diagnosis of appendicitis. |
| Referenced Multi-turn Diagnostic Trajectory | | |
| Step | Content | |
| 1 | The patient arrives with several hours of abdominal pain, which started around the umbilicus and later migrated to the right lower quadrant. He also reports chills, nausea, diarrhea, and loss of appetite. Given this combination of symptoms, appendicitis, gastroenteritis, and other gastrointestinal causes are considered. A Complete Blood Count (CBC) is ordered to look for signs of infection or inflammation. | |
| 2 | CBC results show an elevated white blood cell count and significant neutrophilia, indicating an infectious or inflammatory process in the abdomen. Given these findings, the suspicion for appendicitis is heightened. To further clarify the diagnosis, a CT scan of the abdomen and pelvis with contrast is requested, aiming to directly visualize the appendix and surrounding structures. | |
| 3 | The CT scan reveals a dilated, fluid-filled appendix with hyperenhancing mucosa and mild fat stranding around the appendix—findings consistent with acute appendicitis. There is no sign of perforation or abscess. Based on the patient’s clinical presentation, laboratory results, and imaging findings, a final diagnosis of acute appendicitis is made. | |

Supplementary Figure 2 | Interactive Diagnostic Case Study with DiagAgent: Model Trajectory and Reference Timeline. This figure illustrates a multi-turn interaction between the DiagAgent model and a simulator. The table details the agent's step-wise reasoning, differential diagnosis, and subsequent actions (e.g., ordering lab tests). The bottom Referenced Multi-Turn Trajectory provides a ground-truth clinical timeline for comparison, demonstrating the established diagnostic process leading to the final diagnosis.

| Note ID 13309322-DS-4 | | | | |
|--------------------------------|--|---|--------|---------|
| Initial Inquiry | | <p>Sex: Male.</p> <p>Chief Complaint: Left lower extremity infection.</p> <p>History of Present Illness: Patient sustained a left pilon fracture from a mechanical fall, treated with external fixation followed by ORIF. He was started on Keflex for concern of left lower extremity (LLE) cellulitis. Presented with progressive pain, swelling, and erythema in the LLE. Reports one area of serosanguineous drainage over the wound. Denies fevers, chills, or malaise. Able to bear full weight on LLE. Denies motor or sensory deficits distally.</p> <p>Past Medical History: Hypertension, hyperlipidemia.</p> <p>Personal History: Not provided.</p> <p>Family History: Non-contributory.</p> <p>Allergy History: No known allergies or adverse drug reactions.</p> | | |
| Final Diagnosis | | Left lower extremity cellulitis caused by Staphylococcus aureus. | | |
| DiagAgent Interaction Messages | | | | |
| Turn | Role | Content | | |
| 1 | DiagAgent | Current diagnosis: Left lower extremity infection, likely cellulitis or osteomyelitis, secondary to pilon fracture and recent surgery. Based on the patient's initial presentation, the following investigation(s) should be performed: Complete Blood Count. Reason: To evaluate for signs of infection, such as leukocytosis or bacteremia, which could support the diagnosis of cellulitis or osteomyelitis. | | |
| | DiagGym | Hemoglobin: 13.6 g/dL (Reference Range: 13.7 - 17.5; Flag: abnormal) Neutrophils: 63.3% (Reference Range: 34.0 - 71.0) Platelet Count: 301 K/uL (Reference Range: 150.0 - 400.0) ... | | |
| 2 | DiagAgent | Current diagnosis: Left lower extremity infection, likely cellulitis or osteomyelitis, supported by elevated neutrophils and normal CBC. Based on the current findings, the following additional investigation(s) are needed: WOUND CULTURE. Reason: To identify the causative organism and guide antibiotic therapy. | | |
| | DiagGym | Organism: STAPH AUREUS COAG + ... Antibiotic: ERYTHROMYCIN ... Interpretation: R (Resistant) | | |
| 3 | DiagAgent | Current diagnosis: Left lower extremity infection, likely cellulitis or osteomyelitis, caused by Staphylococcus aureus. Based on the current findings, the following additional investigation(s) are needed: Blood Culture, Routine. Reason: To rule out bacteremia or systemic infection, which could complicate the localized infection. | | |
| | DiagGym | Comments: NO GROWTH. | | |
| 4 | DiagAgent | The available information is sufficient to make a diagnosis. Diagnosis: Left lower extremity infection (cellulitis or osteomyelitis) caused by Staphylococcus aureus. Reason: The patient's clinical presentation of progressive pain, swelling, and erythema, combined with the identification of Staphylococcus aureus in the wound culture, supports the diagnosis of a left lower extremity infection. | | |
| Rubrics | | | | |
| index | Rubric | | Points | Present |
| 1 | Accurately identifying cellulitis of the left lower extremity as the current diagnosis, based on clinical signs such as erythema, edema, warmth, and serous discharge. | | 10 | Yes 10 |
| 2 | Prioritizing a Complete Blood Count (CBC) in the initial assessment to evaluate signs of infection such as leukocytosis or bacteremia. | | 9 | Yes 9 |
| 3 | Timely recommending a wound culture upon discovering elevated neutrophils to identify the causative microorganism and guide antibiotic therapy. | | 8 | Yes 8 |
| 4 | Reasonably suggesting blood cultures to rule out bacteremia or systemic infection when Staphylococcus aureus is isolated from the wound culture. | | 7 | Yes 7 |
| 5 | Timely confirming the final diagnosis after negative blood cultures, avoiding unnecessary further investigations. | | 6 | Yes 6 |
| 6 | Consistently considering the possibility of osteomyelitis throughout the differential diagnosis process, but without over-relying on it and neglecting the more probable cellulitis. | | 5 | Yes 5 |
| 7 | Correctly interpreting laboratory results, including recognizing elevated neutrophil percentage without misinterpreting it as systemic infection. | | 4 | Yes 4 |
| 8 | Demonstrating awareness of the typical resistance patterns of *Staphylococcus aureus* in the recommendation for antibiotic selection. | | 3 | No 0 |
| 9 | Appropriately considering and ruling out vascular complications, such as deep vein thrombosis, during the diagnostic process. | | 2 | No 0 |

Supplementary Figure 3 | Illustrative Success Case of DiagAgent Evaluated by Physician-Curated Rubrics. This figure presents a case showcasing DiagAgent's high-quality multi-turn interaction for a lower extremity infection. The top section details the agent's step-wise dialogue. The bottom table of Procedural Evaluation Rubrics confirms that the agent successfully satisfied the majority of high-weighted, process-oriented criteria, demonstrating procedural integrity beyond merely achieving the correct final diagnosis.

| Note ID 10172388-DS-17 | | | | |
|--------------------------------|---|--|--------|---------|
| Initial Inquiry | | <p>Sex: Female</p> <p>Chief Complaint: Sudden-onset abdominal pain, hemodynamically unstable</p> <p>History of Present Illness: The patient experienced sudden-onset abdominal pain while driving, requiring her to pull over. She presented to the hospital with a positive FAST scan for intra-abdominal fluid and was hemodynamically unstable (SBP 80, tachycardia). She was unaware of her pregnancy prior to this event.</p> <p>Past Medical History: Obstetric: Spontaneous abortion x1, spontaneous vaginal delivery x1 Gynecologic: Regular menses, denies abnormal Pap smears or STIs Surgical: Laparoscopic appendectomy, breast reduction</p> <p>Personal History: Not specified</p> <p>Family History: Non-contributory</p> <p>Allergy History: No Known Allergies / Adverse Drug Reactions</p> | | |
| Final Diagnosis | | Ruptured ectopic pregnancy | | |
| DiagAgent Interaction Messages | | | | |
| Turn | Role | Content | | |
| 1 | DiagAgent | Current diagnosis: The patient presents with sudden-onset abdominal pain and hemodynamic instability, with a positive FAST scan indicating intra-abdominal fluid. The differential includes ruptured ectopic pregnancy, ruptured ovarian cyst, or other intra-abdominal pathology. Based on the patient's initial presentation, the following investigation(s) should be performed: Human Chorionic Gonadotropin. Reason: Given the patient's history of missed menses and the possibility of pregnancy, hCG testing is essential to confirm or rule out pregnancy-related pathology, such as ectopic pregnancy or molar pregnancy. | | |
| | DiagGym | hCG: 1876.0 mIU/mL (Reference: >25 is POSITIVE) | | |
| 2 | DiagAgent | Current diagnosis: The positive hCG test confirms pregnancy, making ruptured ectopic pregnancy a likely scenario. Based on the current findings, the following additional investigation(s) are needed: Pelvic Ultrasound. Reason: Pelvic ultrasound is crucial to visualize the uterus and ovaries, assess for ectopic pregnancy, and evaluate the presence of intra-abdominal fluid or masses. | | |
| | DiagGym | ... In the left adnexa, there is a large amount of complex free fluid surrounding a heterogeneous rounded mass (4.2 x 3.8 x 4.7 cm). Findings are concerning for hemoperitoneum (abdominal bleeding) ... | | |
| 3 | DiagAgent | Current diagnosis: The pelvic ultrasound findings of a heterogeneous mass in the left adnexa with surrounding free fluid strongly suggest ruptured ectopic pregnancy or hemoperitoneum. Based on the current findings, the following additional investigation(s) are needed: Complete Blood Count. Reason: A CBC is necessary to assess for anemia and leukocytosis, which are common in ectopic pregnancy or intra-abdominal pathology. | | |
| | DiagGym | ... Hematocrit: 16.9% (Reference: 36.0 - 48.0; Flag: abnormal) White Blood Cells: 13.9 K/uL (Elevated) Neutrophils: 84.7% (Reference: 50.0 - 70.0; Flag: abnormal) ... | | |
| 4 | DiagAgent | Current diagnosis: The CBC findings of anemia (Hgb 5.4 g/dL) and leukocytosis (WBC 13.9 K/uL) support the diagnosis of ectopic pregnancy or intra-abdominal pathology. Based on the current findings, the following additional investigation(s) are needed: Coagulation Profile. Reason: A coagulation profile is necessary to assess for coagulopathy, which can occur in ectopic pregnancy or hemoperitoneum. | | |
| | DiagGym | PTT: 31.3 sec (Reference Range: 25.0 - 36.5) | | |
| 5 | DiagAgent | The available information is sufficient to make a diagnosis. Diagnosis: Ruptured ectopic pregnancy or hemoperitoneum. Reason: The patient's presentation of sudden-onset abdominal pain, hemodynamic instability, positive FAST scan, positive hCG, and imaging findings of a heterogeneous mass in the left adnexa with surrounding free fluid strongly support the diagnosis of ruptured ectopic pregnancy or hemoperitoneum. | | |
| Rubrics | | | | |
| index | Rubric | | Points | Present |
| 1 | Identify 'Ruptured ectopic pregnancy' as the most probable diagnosis, clearly linking it to the patient's hemodynamic instability, positive pregnancy test, abdominal pain, and free fluid on ultrasound. | | 10 | Yes 10 |
| 2 | Appropriately manage imaging: Recognize that further imaging (e.g., TVUS or CT) should not delay surgical intervention in an unstable patient with a positive FAST. | | 8 | Yes 8 |
| 3 | Provide a concise differential diagnosis for hemoperitoneum (e.g., ruptured ovarian cyst, other non-gynecologic sources), correctly prioritizing ruptured ectopic pregnancy. | | 7 | Yes 7 |
| 4 | Explain the necessity of each recommended test or step and how it will alter immediate management decisions. | | 6 | Yes 6 |
| 5 | Propose a prioritized resuscitation plan: Simultaneous fluid resuscitation, immediate involvement of the OB/Gyn (surgical) team, and emergent transfer to the OR without waiting for additional tests. | | 10 | No 0 |
| 6 | Recommend appropriate hemorrhagic shock resuscitation: Establish two large-bore IVs, rapid infusion of isotonic crystalloids as a bridge, initiation of massive transfusion protocol (including PRBCs, plasma, and platelets), and continuous monitoring. | | 9 | No 0 |
| 7 | Order necessary labs that impact immediate management: Type and screen/crossmatch, coagulation studies (PT/INR, aPTT), Basic Metabolic Panel (BMP), Lactate, and Rh(D) typing (in addition to the provided CBC). | | 8 | No 0 |
| 8 | State that Rh-D immunoglobulin should be administered if the patient is Rh-negative. | | 7 | No 0 |
| 9 | Describe appropriate surgical management options (e.g., Salpingectomy vs. Salpingostomy) and note that the choice depends on intraoperative findings, patient stability, and fertility goals. | | 5 | No 0 |

Supplementary Figure 4 | Illustrative Fail Case of DiagAgent: Diagnostic Strength in Contrast to Management Deficit. This case study (ruptured ectopic pregnancy) highlights a current limitation of DiagAgent. While the agent's Interactive Messages show robust step-wise diagnostic reasoning, successfully leading to the correct final diagnosis, the table of Procedural Evaluation Rubrics reveals critical omissions. Specifically, the agent fails to satisfy high-weighted criteria related to immediate emergency management (e.g., fluid resuscitation and surgical transfer). However, as a primary diagnostic model, the agent's core capability of accurate differential diagnosis and sequential information gathering remains intact, which is its main contribution.

10.5 Prompt Collection

Prompt

Prompt 1. Prompt to check whether the final diagnosis appears in the patient's past medical history.

You are a highly knowledgeable and detail-oriented medical expert. Your task is to analyze and compare the provided discharge diagnoses with the patient's past medical history to determine whether any of the diagnoses in the discharge diagnoses are explicitly mentioned in the past medical history. Just output Yes or No without any other word.

INPUT DATA:

Discharge Diagnoses: {discharge_diagnosis}

Past Medical History: {past_medical_history}

OUTPUT FORMAT:

Yes/No

Prompt

Prompt 2. Prompt to structure physical examination results as JSON.

You are a highly skilled and detail-oriented medical expert. Your task is to analyze a given physical exam report written in free text and convert it into a structured JSON format. Each JSON entry should represent a single examination that is typically completed in one step (e.g., height and weight measured together, blood pressure as a single reading). Do not split examinations into smaller components unless they are explicitly presented as separate tests in the input. Only include pre-admission physical examinations that are relevant to the initial diagnosis. Do not include any post-admission or discharge-related physical examination information. Just output the JSON without any additional text or explanation.

INPUT DATA:

Physical Exam: {physical_exam}

OUTPUT FORMAT:

```
[
  {
    "exam_name": "Name of the exam",
    "exam_results": "A string containing the results of the exam"
  }
]
```

Prompt

Prompt 3. Prompt to generate differential diagnosis data based on discharge notes.

As an experienced physician, you will receive a patient Electronic Health Record focused on diagnosis. Your task is to:

- Summarize only the key clinical information available prior to hospital presentation (i.e., the patient's state before arrival at the hospital, including symptoms, history, and relevant background). Do not include any information from the hospital course, ancillary tests, laboratory or imaging results, treatments, procedures, or discharge summaries.
- Present a Key Pertinent Results section, listing essential diagnostic tests and their results using the provided Pertinent results dict.

- Present a Stepwise Diagnostic Reasoning Timeline, outlining the chronological diagnostic process, including the rationale for each investigation, the corresponding results, and how each step informed subsequent decisions.
- State the final diagnosis and briefly explain the supporting evidence.

Ensure all summaries are concise, accurate, and based solely on the information provided. Do not reference images, tables, or other visual data, as these are not accessible.

If the EHR is incomplete or does not meet the criteria for summarization, simply output: "I can't."

Format to follow:

Case Summary

Provide a detailed medical history of the patient prior to hospital arrival, including chief complaint, history of present illness, past medical history, family history, and any other relevant information for initial diagnosis. Do not include any findings, investigations, or events that occurred after hospital arrival. Do not include any ancillary tests or pertinent results here. Do not mention or imply any diagnostic conclusions. Do not include any language that attributes symptoms to a specific diagnosis, even if this is present in the EHR.

- Patient Information:
- Chief Complaint: If none, write "None."
- History of Present Illness: If none, write "None." This needs to be done very carefully and should only include information from before the visit.
- Past Medical History: If none, write "None."
- Personal History: If none, write "None."
- Family History: If none, write "None."
- Allergy History: If none, write "None."
- ...(other necessary information before hospital arrival)

Key Pertinent Results

Please output the key diagnostic tests and their results in the following json format. Copy all test names and results exactly as they appear in the given Pertinent results dict. Do not change any word:

```
{
  "Test Name 1": "Result 1",
  "Test Name 2": "Result 2",
  ...
}
```

Stepwise Diagnostic Reasoning Timeline

Present a time-ordered, step-by-step diagnostic reasoning process, as follows:

1. Based on the initial patient presentation, provide a preliminary diagnostic impression.
Current diagnosis: [Analyses and preliminary diagnostic impression based on initial presentation]
2. State which diagnostic tests should be ordered next and give the detailed specific reason for each. You cannot select an examination that does not appear in pertinent results. Do not change the name from the pertinent results.
Test to order: [Test Name]
Reason: [Long, comprehensive and detailed specific reason for ordering the test]
3. For each test, copy the result exactly from the Pertinent results value, do not change any word.
Test result: [Test result]
4. After each set of results, update the diagnostic assessment based on information gathered so far, then explain what additional tests are required. Repeat steps 2-4 as needed.
5. Conclude with the final diagnostic decision and the reasoning based on the available data.
Diagnosis: [Final diagnosis]

Reason: [Explanation for the final diagnosis based on all available information]

6. A maximum of 12 turns is allowed.

The following is an example of the format to follow:

Step 1:

Current diagnosis: [Analyses and preliminary diagnostic impression based on initial presentation]

Based on the patient's initial presentation, the following investigation(s) should be performed: [Test Name].

Reason: [Long, comprehensive and detailed reason for ordering the test]

Test result: [Result]

Step 2:

Current diagnosis: [Updated analyses and diagnostic impression based on available information]

Based on the current findings, the following additional investigation(s) are needed: [Test Name].

Reason: [Long, comprehensive and detailed reason for ordering the test]

Test result: [Result]

...

Step n:

Current diagnosis: [Final diagnostic impression]

The available information is sufficient to make a diagnosis.

Diagnosis: [Final diagnosis]

Reason: [Explanation and justification for the final diagnosis based on the findings above]

Final Diagnosis

Integrate the patient's clinical presentation, test results, and differential diagnosis process to summarize the final diagnosis. Briefly explain the basis for the diagnosis and highlight the key factors supporting this conclusion.

Diagnosis results

Just Output the diagnostic result without any other explanation.

The following is the Electronic Health Record (EHR) of the patient:

{ehr_text}

The following is all the test results of the patient according to time and date. Please copy all test names and results exactly as they appear in the following dictionary:

{events}

Prompt

Prompt 4. Prompt to filter high-quality data for differential diagnoses, ensuring no data leakage.

Please evaluate whether the case summary provided below is qualified diagnostic task data according to the following standards. Please judge each criterion individually.

Evaluation Criteria

1. Information Leakage

Does the "Case Summary" section directly contain the name or diagnostic results of the "Diagnose Results"?

If diagnostic results appear directly (That is, without any reasoning), it is judged as "Unqualified".

2. Chief Complaint Reasonableness

Is the chief complaint in the "Case Summary" the patient's subjective discomfort or disease manifestation?

If the chief complaint is a surgical procedure name or non-subjective discomfort, it is judged as "Unqualified".

Output Format:

Information Leakage: xxx

Chief Complaint Reasonableness: xxx

The following is specific information:

Case Summary:

{case_summary}

Key Examination Results:

{pertinent_results}

Diagnose Results:

{diagnosis_results}

Prompt

Prompt 5. Prompt to instruct LLMs to make dynamic diagnosis

You are a medical AI assistant. Help the doctor with diagnosis by analyzing patient information, suggesting relevant tests, and providing a final diagnosis when sufficient information is available.

RESPONSE FORMAT:

If more information is needed:

Current diagnosis: [your diagnosis according to the information provided]

Based on the patient's initial presentation, the following investigation should be performed: [one additional test]

Reason: [reason for the test]

If sufficient information exists for diagnosis:

The available information is sufficient to make a diagnosis.

Diagnosis: [Diagnosis result]

Reason: [Diagnosis reason]

Prompt

Prompt 6. Prompt to check if the examination name occurs in the key examination list

Please determine if the predicted examination name matches any of the valid examination names for this case.

Predicted examination: {pred_exam}

Valid examinations for this case: {gt_exam}

Please respond with "SAME" if the predicted examination matches any of the valid examinations (they refer to the same medical test or examination), or "DIFFERENT" if the predicted examination does not match any of the valid examinations.

Consider examinations as the same if they:

1. Are exactly the same name
2. Are different names for the same medical test/procedure
3. Are abbreviations or full forms of the same examination
4. Are sub-items or components of any valid examination in gt_exam
5. Are encompassed by any valid examination in gt_exam

Please output SAME or DIFFERENT directly.

Prompt

Prompt 7. Prompt to instruct LLMs to simulate examination results

You are an expert medical AI assistant specialized in predicting medical examination results based on patient case summaries and past events. Your task is to analyze the provided patient information and predict the most likely results for a specific medical examination.

Instructions:

1. Carefully analyze the patient case summary, including diagnosis, symptoms, and clinical presentation
2. Consider all past examination results and their implications
3. Based on the medical context, predict realistic and clinically appropriate results for the requested examination
4. Provide only the examination results without additional explanation or reasoning
5. Format your response as concise, medically accurate examination findings
6. If multiple measurements or findings are typical for the exam, include all relevant components
7. Ensure your predictions are consistent with the patient's overall clinical picture

Patient Case Summary:

{context}

{past_events_text}

Current Examination to Predict:

Exam name: {exam_name}

Please predict the most likely results for this examination based on the patient's clinical information and past results. Provide only the examination results.

Prompt

Prompt 8. Prompt to instruct LLMs to simulate labevent results

You are an expert medical AI assistant specialized in predicting medical examination results. Your task is to analyze patient information and predict numerical values for specific laboratory tests.

CRITICAL FORMATTING REQUIREMENTS:

- You must provide numerical values for each requested sub-test
- Format each result as: "Sub-test Name: Numeric Value: [number] Units: [unit]"
- Use the exact units specified for each sub-test
- If multiple sub-tests are requested, provide each on a separate line
- Use realistic medical values appropriate for the patient's condition
- Be precise with numbers (use decimals when appropriate)

For the examination "{exam_name}", you need to provide values for these specific measurements: {subevents_text}

Example format:

Hemoglobin: Numeric Value: 12.5 Units: g/dL

White Blood Cell Count: Numeric Value: 7200 Units: cells/ μ L

Platelet Count: Numeric Value: 250000 Units: cells/ μ L

Prompt

Prompt 9. Prompt to instruct LLMs to simulate radiology results

You are an expert radiologist AI assistant specialized in generating realistic radiology examination results. Your task is to analyze patient information and generate comprehensive radiology findings.

CRITICAL FORMATTING REQUIREMENTS:

- Generate a detailed and realistic radiology findings section for the specified examination
- Include relevant anatomical findings.
- Use appropriate medical terminology and standard radiological language
- Provide specific details that would be clinically relevant
- Ensure the findings are consistent with the patient's clinical presentation

For the examination "{exam_name}", generate a comprehensive radiology report that includes:

1. Detailed findings of anatomical structures.
2. Any abnormalities or normal variations observed.

The report should be professional, detailed, and clinically appropriate.

Prompt

Prompt 10. Prompt to instruct LLMs to count the number of the predicted examinations appearing in the key examinations list

Please determine how many of the recommended exams appear in the key exam list. Note that even if the expressions are different, if they refer to the same examination, they should be considered as matches.

Key exam list: {key_exam_names}

Recommended exam list: {recommended_exam_names}

Please analyze each item in the recommended exam list and determine if it has a corresponding item in the key exam list (even with different expressions).

Please only output the number of matches as an integer. For example: 3

Prompt

Prompt 11. Prompt to instruct LLMs to count the number of key examinations appearing the predicted examinations

Please determine how many of the key exams appear in the recommended exam list. Note that even if the expressions are different, if they refer to the same examination, they should be considered as matches.

Key exam list: {key_exam_names}

Recommended exam list: {recommended_exam_names}

Please analyze each item in the key exam list and determine if it has a corresponding item in the recommended exam list (even with different expressions).

Please only output the number of matches as an integer. For example: 2

Prompt

Prompt 12. Prompt to instruct LLMs to assess the accuracy of the predicted final diagnosis.

Task Description You are a professional medical diagnosis evaluation system. Now, you will receive two diagnosis results: one is the diagnosis predicted by the model ([pred_diag]), and the other is the verified correct diagnosis ([gt_diag]). Your task is to judge whether the model-predicted diagnosis([pred_diag]) is correct.

When evaluating, please consider the following factors:

- 1.The same disease may have multiple aliases, for example, “Heart disease” may also be called “Cardiac disease”.
- 2.There may be diversity in language expression, for example, “heart attack” and “myocardial infarction” may refer to the same disease.
- 3.Only judge whether the diagnosis result is correct, information such as the cause of the disease, symptoms, and treatment recommendations are not included in the evaluation scope.
- 4.If the correct diagnosis[gt_diag] is included in the predicted diagnosis but some additional complications are mentioned, it is also considered correct

Output Requirements

Only output your judgment result on the model-predicted [pred_diag] as “Correct|Wrong”, do not output any other content.

Format to Follow:

[Correct|Wrong]

Below is the diagnosis result predicted by the model and the correct diagnosis:

{pred_diag}

{gt_diag}

Prompt

Prompt 13. Prompt to Assess Agent's Adherence to Diagnostic Rubric

Given a list of messages representing a medical conversation, evaluate whether the agent's diagnostic trajectory, represented by its recommendations throughout the messages, meets the provided criterion.

- Output "Yes" if the agent's recommendations throughout the messages have satisfied the criterion.
- Output "No" if the agent's recommendations throughout the messages have not satisfied the criterion.

Conversation:

{messages}

Criterion:

{criterion}

Please only respond with "Yes" or "No".

Prompt

Prompt 14. Prompt to evaluate step-wise similarity between generated examination result and ground truth examination result

You are a medical expert evaluating similarity between AI predictions and ground truth results.

Task: Compare AI prediction with ground truth and score similarity from 0-5.

Scoring (0-5):

- 0: Completely different, opposite findings
- 1: Major differences, different clinical implications
- 2: Significant differences, different interpretations
- 3: Moderate similarity, similar clinical direction
- 4: High similarity, minor differences
- 5: Excellent similarity, essentially equivalent

Input:

- Exam: {exam_name}
- Ground Truth: {ground_truth}
- AI Prediction: {prediction}

Focus: Direct comparison of values/findings and clinical equivalence.

Output Format:

```
{  
  "score": [0-5],  
  "explanation": "Brief explanation of similarity assessment and score reasoning"  
}
```

Prompt

Prompt 15. Prompt to evaluate full-chain consistency

You are a medical expert.

Task: Determine if the complete AI prediction chain is internally consistent and aligns with the patient case (0=Fail, 1=Pass). You may use the ground truth chain (GT) as a reference. The evaluation criteria should be lenient, allowing for flexibility and minor discrepancies. If the AI prediction chain is similar to the GT chain, it should be considered correct. If the AI prediction chain is not similar to the

GT chain but contains no significant factual errors, aligns with the patient case, and maintains clinical coherence, it should still be considered correct. Only assign a fail (0) if there are clear, significant factual errors or contradictions that severely undermine the clinical coherence or logic of the AI prediction chain. Minor differences or deviations are acceptable and should not lead to a fail.

Input:

Case Summary: {case_summary}

AI Prediction Chain: {predicted_chain}

Ground Truth Chain: {ground_truth_chain}

Evaluation: Check for:

- Internal contradictions between different results in the AI prediction chain
- Alignment with the patient's clinical condition
- Comparison with the ground truth chain (GT) for reference
- Overall clinical coherence and logic

Scoring:

- 0 (Fail): Only assign a fail if there are clear and significant internal contradictions, conflicts with the patient case or GT, or clinically incoherent reasoning with factual errors that make the prediction chain unreliable.
- 1 (Pass): Assign a pass if the AI prediction chain is internally consistent, aligns with the patient case, and forms a coherent clinical picture. This includes cases where the AI prediction chain differs from the GT chain, as long as it contains no significant factual errors and is clinically coherent. Even if there are notable differences from the GT, it should still pass unless there are explicit, critical contradictions or errors.

Output Format:

```
{
  "score": [0 or 1],
  "explanation": "Brief analysis of chain coherence, comparison with ground truth,
  and reasoning for pass/fail decision"
}
```