

Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents

JUNKAI LI^{†#}, SIYU WANG[†], MENG ZHANG[†], WEITAO LI^{†#}, YUNGHWEI LAI[†], XINHUI KANG^{†#}, WEIZHI MA[†], and YANG LIU^{#†}

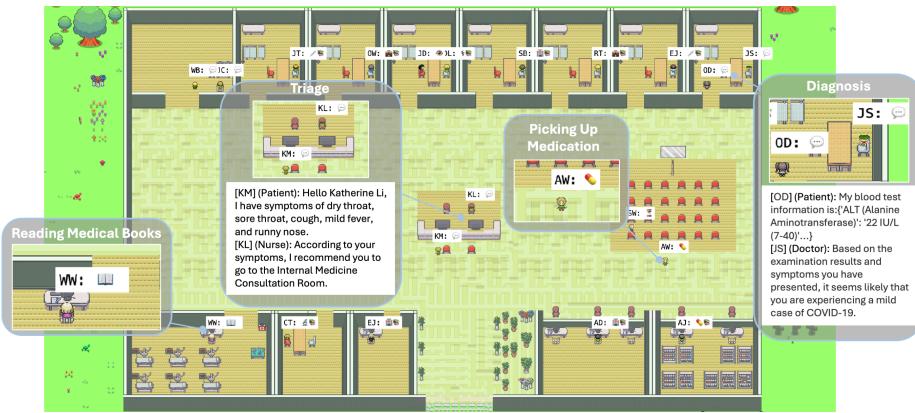


Fig. 1. An overview of Agent Hospital. It is a simulacrum of hospital in which patients, nurses, and doctors are autonomous agents powered by large language models. Agent Hospital simulates the whole closed cycle of treating a patient’s illness: disease onset, triage, registration, consultation, medical examination, diagnosis, medicine dispensary, convalescence, and post-hospital follow-up visit. An interesting finding is that the doctor agents can keep improving treatment performance over time without manually labeled data, both in simulation and real-world evaluations.

Abstract In this paper, we introduce a simulacrum of hospital called Agent Hospital that simulates the entire process of treating illness. All patients, nurses, and doctors are autonomous agents powered by large language models (LLMs). Our central goal is to enable a doctor agent to learn how to treat illness within the simulacrum. To do so, we propose a method called MedAgent-Zero. As the simulacrum can simulate disease onset and progression based on knowledge bases and LLMs, doctor agents can keep accumulating experience from both successful and unsuccessful cases. Simulation experiments show that the treatment performance of doctor agents consistently improves on various tasks. More interestingly, the knowledge the doctor agents have acquired in Agent Hospital is applicable to real-world medicare benchmarks. After treating around ten thousand patients (real-world doctors may take over two years), the evolved doctor agent achieves a state-of-the-art accuracy of 93.06% on a subset of the MedQA dataset that covers major respiratory diseases. This work paves the way for advancing the applications of LLM-powered agent techniques in medical scenarios.

Corresponding authors: Weizhi Ma (mawz@tsinghua.edu.cn) and Yang Liu (liuyang2011@tsinghua.edu.cn). [†]: Institute for AI Research (AIR), Tsinghua University, China. #: Department of Computer Science and Technology, Tsinghua University, China.

1 INTRODUCTION

Large Language Model (LLM) agents have demonstrated promising performance in various tasks, including code generation [18], information game [29], and question answering [20], etc. Motivated by the capabilities of LLM agents, some studies adopt them to simulate human interaction and behavior rather than dealing with single tasks, such as information spreading on social media [25] and the “Stanford Town” simulacrum project [17]. People from both academia and industry all believe that LLM agents will have significant impacts on various scenarios.

Despite the success achieved by existing studies, LLM agents are typically used either to solve specific tasks or for social simulation. This raises a question: can we combine the two capabilities? i.e., could the process of social simulation enhance the performance of LLM agents on specific tasks? Motivated by this assumption, we aim to verify it by designing a simulacrum of hospital for the evolution of medical LLM agents. There are mainly two reasons for why choosing the medical scenario: On one hand, AI for medicine is widely believed as an ideal scenario in which AI technology can make significant improvements, and effective simulation of a hospital could facilitate related studies [9]. On the other hand, the hospital setting involves multiple typical tasks, such as disease diagnosis and heterogeneous signal understanding, which serve as excellent benchmarks to verify if simulation environments can help LLM agents in evolution.

In this study, we develop a comprehensive simulation that covers nearly all medical processes in a hospital, which is named Agent Hospital. An overview of the simulated environment is shown in Figure 1. There are mainly two types of agents: residents (potential patients) and medical professionals. Each resident will choose to visit the hospital upon developing a disease. During their stay, they undergo a series of procedures as the real-world healthcare process in a hospital, including triage, registration, consultation, examination, diagnosis, and treatment. In Agent Hospital, after receiving their treatment plan, residents’ health status changes are predicted with the help of LLMs, and they will actively report back to the hospital once they recover as a follow-up. This simulation of the entire hospital interaction process provides an excellent platform for further study, e.g., treating infinite simulated patients by a single doctor agent for evolution.

Based on Agent Hospital, we aim to train proficient doctor agents to handle medical tasks such as diagnosis and treatment recommendation, which are critical in hospital settings. Traditional studies often integrate medical knowledge into LLMs/agents to construct powerful medical models through pretraining, supervised fine-tuning, or retrieval-augmented generation strategies. However, we propose a novel strategy that trains doctor agents by simulating doctor-patient interactions within the simulated environment. Due to no manually labeled data utilized, we name the proposed strategy as MedAgent-Zero. The doctor agent interacts with various patient agents in Agent Hospital, evolving into a more brilliant agent by accumulating records from successful cases and deriving experience from failed cases. Due to the low cost and high efficiency of doctor agent training, we can enable the agent to easily handle tens of thousands of cases within just a few days, achieving that would take a real-world doctor several years to manage.

We conduct two types of experiments to verify the effectiveness of the evolved doctor agent by MedAgent-Zero strategy in our hospital. On one hand, within the virtual hospital, we conduct generated patient interaction experiments ranging from 100 to 10,000 (a human doctor may treat around 100 patients in a week), covering 8 different respiratory diseases,

more than a dozen types of medical examinations, and three different treatment plans for each disease. The doctor agent trained via the MedAgent-Zero strategy continually self-evolved during the process of handling simulated patients, ultimately achieving accuracy rates of 88%, 95.6%, and 77.6% in the examination, diagnosis, and treatment tasks, respectively. On the other hand, we adopt the evolved doctor agent to attend an evaluation on a subset of the MedQA dataset [7]. Surprisingly, even without any manually labeled data, our agent achieved state-of-the-art performance after evolving in the Agent Hospital. These results demonstrate that the simulation environment can effectively assist the evolution of LLM agents in dealing with specific tasks.

The main contributions of our work are summarized as follows:

- To the best of our knowledge, this is the first simulacrum of hospital, which comprehensively reflects the entire medical process with excellent scalability, making it a valuable platform for the study of medical LLMs/agents.
- Based on this virtual environment, we propose the MedAgent-Zero strategy that is designed for the self-evolution of medical agents without manually labeled data. The record accumulation and experience summarization modules allow doctor agents to continuously evolve their capabilities by processing medical cases and engaging in self-feedback within Agent Hospital, thus enhancing their ability to handle various medical tasks.
- In experiments with simulated cases, MedAgent-Zero can handle tens of thousands of cases within several days (human doctors may take over two years) and demonstrates powerful performance. Furthermore, the experience that automatically accumulated in Agent Hospital even enables our doctor agent to achieve state-of-the-art performance on the real-world evaluation dataset - a subset of MedQA, even without any manually labeled data.

2 RELATED WORK

2.1 LLM-based Real World Simulation

Recent research initiatives have leveraged LLMs to replicate real-world dynamics. In multiple fields such as epidemiology, sociology, and economics, researchers are utilizing agents based on LLMs to simulate human decision-making, leading to many exciting emergence phenomena in various domains.

Agents based on LLMs demonstrate their capability to interact in a competitive, dynamic real-world simulation [34] by providing services and receiving feedback, aligning their behaviors with existing sociological and market theories to ensure fair competition among various agents. Recommendation systems powered by LLM-based agents [31, 32] engage in collaborative learning by understanding user queries and communicating needs, while also simulating user preferences and behaviors within recommendation frameworks, thereby offering deeply personalized recommendation services.

LLMs are utilized to simulate macroeconomic activities [11, 12] by developing agents that are capable of understanding human decision-making processes, thereby adapting to complex economic environments and improving the prediction of economic behaviors and trends. The social network simulation system [4, 13] uses LLMs to equip agents with the capability to mimic human behaviors such as individual emotions, attitudes, and interactive behaviors, using real social network data to build simulation environments.

The public administration crisis simulation system based on generative agents [28] enables non-technical social science researchers to conduct complex simulations easily and analyze individual and group behaviors in public management crises. Moreover, a multi-agent warfare simulation system based on LLMs [5] has been designed to simulate and analyze large-scale historical conflicts, such as World Wars. This system reconstructs historical military events, political negotiations, and strategic planning, exploring the processes leading to peace and future conflict prevention.

In the field of epidemiology, researchers have integrated LLMs to develop models that simulate individual behaviors. Each agent driven by its own personality, health status, and perception of epidemic spread, mimics real-world behavioral patterns, providing an understanding of the dynamic relationship between social behaviors and disease transmission [27].

These applications demonstrate the capability of LLMs to simulate real-world dynamics, providing us with many exciting insights. Utilizing agents based on LLMs for healthcare simulations is also a promising direction. However, existing studies [8, 24] mostly focus on simulating the treatment tasks rather than fully simulating the whole closed cycle of treating a patient's illness. They also fail to allow doctors to evolve throughout this process, let alone simulate societal healthcare events such as seasonal flu outbreaks. Our work aims to fill the gaps in existing research, hoping to better leverage LLMs for diagnosis and treatment, providing more accurate and effective support for medical decision-making.

2.2 Evolution of Agents

Recent advancements in LLMs have considerably inspired their integration and application in various fields. Current LLMs have achieved multiple breakthroughs through methodologies such as pre-training [2], fine-tuning [19], and other forms of human-supervised training [16]. However, current LLMs may encounter limitations in performance as task complexity and diversity escalate. The existing training paradigms, which require the use of extensive data corpora or heavy human supervision, are deemed costly. Therefore, the development of self-evolutionary approaches has gained momentum. These approaches enable LLM-powered agents to autonomously acquire, refine, and learn through self-evolving strategies.

One approach involves the direct and explicit integration of external knowledge to enhance task-solving capabilities. For instance, Self-Align [23] uses a topic-guided method to collect topics across 20 scientific domains, including scientific and legal expertise. SOLID [1] generates structured knowledge about entities to initiate conversations. UltraChat [3] compiles unstructured knowledge from 20 different types of textual materials, organized around 30 meta-concepts, to design conversational tasks. The accumulation of such experience empowers agents to address a broader array of tasks.

Agents can also engage in a reflection process during solution generation to facilitate self-evolution. STaR [30] produces rationales when addressing tasks. Should an error occur, it revises both the rationale and the response, which then informs the fine-tuning of the model to enhance model performance. LSX [22] introduces two interconnected modules working in tandem to evolve: a learner module that executes a foundational task and a critic module that evaluates the quality of explanations provided by the learner. Furthermore, SelfEvolve and LDB [6, 35] enhance an agent's capability in code generation by enabling it to reflect on and learn from feedback generated during operation. Through such reflective processes, agents can self-evolve, refine their methodologies, and thus achieve improved performance.

The development of agent self-evolution strategies appears promising. However, current studies on agent evolution predominantly concentrate on isolated and standalone tasks, with

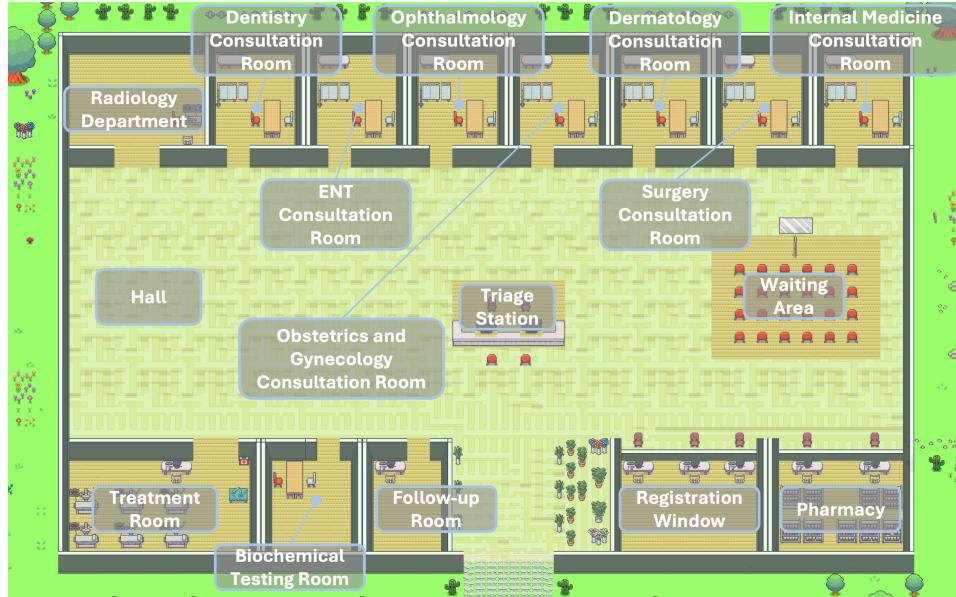


Fig. 2. The distribution of various areas within Agent Hospital.

insufficient integration into world simulations, particularly in contexts such as healthcare simulations where there is a dynamic progression in the patient's condition over time. There is hence considerable potential in employing self-evolving LLM agents without real-world environments.

3 HOSPITAL SIMULACRUM

3.1 Environment Settings

To visualize the entire consultation process, we first design a hospital sandbox simulation environment. Inspired by previous studies [17], the hospital sandbox is implemented by the Tiled¹ and Phaser², where Tiled is a highly flexible map designing tool and Phaser is a framework to manage the movements and interactions of agents on the sandbox. Finally, as shown in Figure 2, there are 16 areas with distinct functions in the Agent Hospital, including triage station, consultation rooms, examination rooms, etc.

3.2 Agent Roles

We designed two types of roles for interaction in the hospital, including medical professional agents and resident agents (who may become patients at any time). The information about these roles is generated using LLM (GPT-3.5) and can be easily expanded.

3.2.1 Medical Professional Agents. As shown in Figure 2, there are various consultation and examination rooms, so we need a series of medical professional agents to work at the Agent Hospital, including 14 doctors and 4 nurses. Our doctor (physician) agents are engineered to diagnose diseases and create detailed treatment plans, whereas our nursing agents focus on

¹<https://www.mapeditor.org>

²<https://phaser.io>

<p>Patient</p> <p>Name: Kenneth Morgan Age: 35 Gender: Male Disease: Acute Nasopharyngitis Medical History: Hypertension Symptoms: Diarrhea, persistent vomiting, enlarged cervical lymph nodes, recurrent fever, abdominal pain, headache</p> 	<p>Internal Medicine Doctor</p> <p>Name: Élise Martin Age: 32 Gender: Female Skill: Excellent communication and empathetic patient care abilities Duty: Diagnose, treat, and provide preventive care for adult patients with a broad range of acute and chronic illnesses</p> 
<p>Radiologist</p> <p>Name: Zhao Lei Age: 58 Gender: Male Skill: Strong analytical skills and detailed observational abilities Duty: Interprets medical images such as X-rays, MRIs, CT scans, and ultrasounds to diagnose patient conditions</p> 	<p>Receptionist</p> <p>Name: Fatoumata Diawara Age: 48 Gender: Female Skill: Excellent communication skills and proficiency with office software Duty: Manages appointment scheduling, patient check-in, and communication coordination</p> 

Fig. 3. Example agents in Agent Hospital.

triage, supporting the day-to-day therapeutic interventions. More information about several simulated medical professional agents is summarized in Figure 3, e.g., Elise Martin is a female internal medicine doctor, who works in the internal medicine consultation room, and Zhao Lei is a male radiologist that good at interpreting medical images.

3.2.2 🩺 Resident Agents. Another type of role in the hospital is the patient. Our simulation starts when patient agents are healthy, so we prefer to name them as resident agents (may transform into patient agents once get ill). To simplify the interaction, we assume that medical professional agents will not develop diseases as resident agents. Each resident agent has distinct demographic information, and may get a disease randomly. As shown in Figure 3, Kenneth Morgan is a male resident with the disease. Upon contracting an illness, these agents automatically initiate a process to seek medical assistance, reflecting typical patient behavior in a clinical setting.

3.3 Planning

In order to enhance the realism of our Agent Hospital, the actions and interactions of both resident agents and medical professional agents are strategically planned and dynamically rescheduled when the agent gets instructions (e.g., going to a specific consultation room). This approach goes beyond merely simulating the standalone treatment procedures. Instead, it comprehensively models disease progression and recover over time.

3.3.1 Planning of Residents. Residents (Patients) play the most active roles in hospitals, so we introduce their plans first. There are mainly two types of plans: 1) Daily Planning. Resident agents are randomly manifesting illnesses, and if failed into diseases, they will schedule hospital visits. 2) Dynamic Planning. Upon arrival at the hospital, the resident, also a patient now, will go to the triage station. The actions and movements of the patient agents are dynamically adapted according to the sequence of triage, registration, consultation, examination, diagnosis, and treatment processes. These processes are dynamically generated based on the patient profiles and the responses of medical professional agents to the evolving clinical situation. Patients navigate this adaptive sequence customized according to their specific profiles and the dynamic progression of their disease. This dynamic configuration

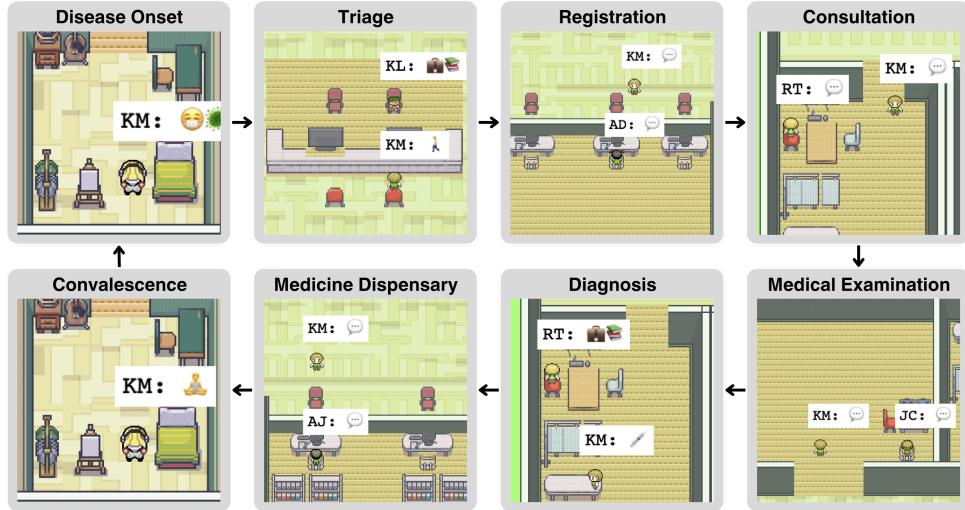


Fig. 4. Resident agent Kenneth Morgan falls ill and visits the Agent Hospital. Triage nurse Katherine Li conducts an initial evaluation of Mr. Morgan’s symptoms and refers him to the dermatology department. Mr. Morgan then registers at the hospital’s counter and is subsequently arranged for a consultation with dermatologist Robert Thompson. After undergoing the prescribed medical examination, Mr. Morgan receives a diagnosis and medication. He returns home to rest and monitor the improvement of his condition.

enables a comprehensive evaluation of treatment efficacy and disease management strategies, providing an accurate simulation of a controlled yet realistic virtual environment.

Moreover, the agent’s health is continuously monitored. Depending on the effectiveness of the treatment and medication, their condition may improve or deteriorate. If the condition worsens, the agent will schedule another hospital visit for the next day. Conversely, if there is an improvement, the agent remains at home, recovering progressively each day until fully restored to health. For agents who recover and maintain a healthy status, the system randomly generates new symptoms and diseases daily, reflecting the unpredictability of real-world medical conditions. This initiates another cycle of hospital visits and consultations.

3.3.2 Planning of Medical Professionals. Medical professional agents are assigned to specific stations within the hospital where they fulfill their responsibilities based on their designated roles. They have fewer action types than residents, but they should sharpen their expertise to achieve better treatment levels in the Agent Hospital. They are trained from two types of actions: 1) **Practice**. Doctor agents are positioned in their respective departments, where they manage clinical sessions and provide care to patients assigned to them during their shifts. The follow-up information from patients will help them polish their medical records experience. 2) **Learning**. Outside of working hours, they engage in studying past medical records to gain clinical experience, and reading medical textbooks to expand their knowledge.

3.4 Patient Events/Interactions

Patient agents typically experience eight main types of events or interactions, which are the most active roles in the Agent Hospital. For clarity, we will begin with an example. Figure 4 illustrates a case study featuring resident agent Kenneth Morgan, who woke up with a skin condition, and decides to seek medical attention at a hospital. Initially, he undergoes a preliminary evaluation at the triage station, where his symptoms are assessed. Based on this evaluation, he is referred to the dermatology department for a consultation with a specialist. Upon arrival, Morgan registers at the hospital reception, which organizes his consultation schedule. He then waits in the designated area until he is called to the dermatologist's office. During the consultation, the dermatologist agent determines the need for a medical examination which Morgan undergoes. The dermatologist agent provides a diagnosis, outlines a treatment strategy, and prescribes medication after reviewing the results. Finally, Morgan collects his medication from the hospital pharmacy and returns home to commence his recovery.

Next, we introduce the patient events and interactions in detail:

3.4.1 Disease Onset. Resident agents contract diseases from a predefined dataset at random upon awakening. Each disease is categorized into one of three severity levels: mild, moderate, or severe. The simulation assigns LLM-generated specific disease symptoms, diagnostic results, potential complications, all types of examination results, differential diagnoses, confirmed diagnoses, treatment protocols, and preventative measures to each resident. These details are encapsulated in the complete medical record, which is depicted in Figure 16 in the Appendix. Note that all disease information is unseen by medical professional agents, they can only ask patients/conduct medical examinations to get information.

Medical knowledge is involved in disease simulation, and generates patient case reports based on this information in conjunction with the severity of the patient's illness. The cases we generate strictly adhere to medical principles, and the data construction process involves manual verification of information to ensure that the content generated conforms to medical logic. For example, the applied medical knowledge in simulating COVID-19 patients is illustrated in Figure 15 in the Appendix.

3.4.2 Triage. Once arrived the Agent Hospital, the patient's journey begins at the triage station. Patients arrive and describe their symptoms to the nursing agents. The instructions guide the nursing staff in their decision-making, enabling them to direct patients to the appropriate specialist departments where medical professional agents are available to conduct further diagnostics.

3.4.3 Registration. After the initial assessment, patients follow the advice from the triage station and proceed to register at the registration counter. They then wait in the designated waiting area for their consultation turn with the specialists from the respective departments.

3.4.4 Consultation. When it is their turn for consultation, patients engage in a preliminary dialogue with the physician agents to describe their symptoms and the duration since onset. The physician then determines which medical examination is needed to investigate the cause and assist with diagnosis and treatment. In the current version, only one type of medical examination will be conducted for each patient based on the decisions made by doctor agents.

3.4.5 Medical Examination. After receiving the prescribed list of medical examinations, patients proceed to the relevant department to undergo the tests. The resulting medical data

which are pre-generated by LLM are subsequently presented to the patient and the doctor. This process designed to mimic real-time diagnostic feedback, aligns with the presentation of symptoms as illustrated in the patient profiles found in Figure 16.

3.4.6 Diagnosis. Subsequent to the medical examination, patients are guided to the respective department where physician agents undertake the diagnostic process. Patients disclose their symptoms and share the results of the medical examination with the physician agents, who then undergo diagnostic processes based on a predefined disease set. The diagnostic result is promptly communicated back to the patient, showcasing the model’s capacity to integrate complex medical data and its advanced diagnostic ability.

3.4.7 Treatment Recommendation. The medical agent is presented with the patient’s symptoms, results from medical examinations and the diagnosis of the disease they made. In addition, three distinct treatment plans tailored to mild, moderate, and severe conditions are also provided. The doctor is then tasked with selecting the appropriate plan from the mild, moderate, or severe options, according to the patient’s specific needs. If any medicine is prescribed, patients proceed to the dispensary to collect it.

3.4.8 Convalescence (Follow Up). At the end of the diagnostic and treatment process, the patient provides feedback or updates on their health condition for follow-up actions. To mimic the dynamic progression of diseases accurately, the LLM-enhanced simulation involves a few key steps: doctors devise treatment plans based on the patient’s detailed health information and test results, and then these details – specifically the patient’s symptoms, the prescribed treatment plan, and the diagnosed disease are incorporated into a template for simulation.

3.5 Medical Professional Events

Besides interacting with patient agents, medical professional agents, particularly doctor agents, mainly engage in the following two types of actions. Both of these actions are aimed at enabling the self-evolution of medical agents within Agent Hospital.

3.5.1 Practice. Doctor agents continuously learn and accumulate experience during the treatment process in Agent Hospital, thereby enhancing their medical capabilities similar to human doctors. We assume that doctor agents are constantly repeating this process during all working hours. A newly designed evolution strategy is introduced in detail in Section 4.3.

3.5.2 Learning. Apart from improving their skills through clinical practice, doctor agents also proactively accumulate knowledge by reading medical documents outside of work hours. This process primarily involves strategies to avoid parametric knowledge learning for agents, which we will also introduce in Section 4.

4 METHODOLOGY

4.1 Definition of Medical Tasks

LLM evaluation tasks generally adopt a multiple-choice format, the performance of which is easier to measure than generation tasks. Therefore, we also format some representative medical tasks in this way to assess the capabilities of medical agents, including examination judgment, diagnosis, and treatment plan. It is noteworthy that our methodology focuses on how to enhance the doctor (physician) agents.

We define three medical tasks here:

- (1) Examination Decision: The patient agent tells about her/his symptoms, and the doctor agent should select one medical examination from available options. As there are several examination options for each disease, if the examination selected by the doctor agent is in the ground-truth list, it can be seen as the right answer.
- (2) Diagnosis: Based on the patient’s symptoms and the medical examination results, the doctor agent should give a diagnosis to the patient. All candidate diseases are included in the prompt, and the answer of the doctor agent is correct only if the generated disease name is the same as the patient’s actual disease.
- (3) Treatment Plan: Based on the patient’s symptoms and the diagnosis result, the doctor agent should decide on an appropriate treatment plan for the patient. To avoid diverse outputs that are hard to evaluate, all diseases are categorized into three treatment levels, namely mild, moderate, and severe. The answer is correct if the output matches the ground truth.

4.2 Datasets

4.2.1 Simulated Medical Dataset. In Section 3.4, we introduced how to generate simulated electronic health records for patients based on LLMs. Generating such records also requires foundational domain knowledge. So we have primarily collected data on eight representative respiratory diseases from the 8th Edition of the “Infectious Diseases” [10], including 8 diseases: Acute Nasopharyngitis, Acute Rhinitis, Bronchial Asthma, Chronic bronchitis, COVID-19, Influenza A, Influenza B, and Mycoplasma infection. For each disease, the symptoms, laboratory test/examination results, and treatment plans are collected. Figure 15 shows the medical knowledge about COVID-19.

After that, the disease knowledge is added to the medical records generation prompt. We generate around ten thousand records by LLM, where 10,000 records are adopted for training, and 500 records are used for testing. Each record involves three medical tests (examination, diagnosis, and treatment) to help doctor agent evolve in practice. This dataset is named Simulated Medical Dataset.

4.2.2 Medical Document Dataset. Doctor agents also enhance themselves by learning, so that we collect some medical news/textbook datasets about respiratory diseases. Over 5M tokens from the medical news website³ and about 4M tokens from the Merck Manuals website⁴ are collected. To help the agents learn from them, these documents are utilized to generate multiple-choice questions as the simulated medical dataset. While different from the simulated medical datasets where choices are from a predefined set (e.g., the diagnosis choices are the eight mentioned diseases), the multiple-choice questions here are all LLM-generated. This dataset is named Medical Document Dataset.

4.3 Evolution

To facilitate the evolution of LLM-powered medical agents, we propose MedAgent-Zero strategy, which is shown in Figure 5. MedAgent-Zero is a parameter-free strategy, and no manually labeled data is applied as AlphaGo-Zero [21]. There are two important modules in this strategy, namely the Medical Record Library and the Experience Base. Successful cases, which are to be used as references for future medical interventions, are compiled and stored in the medical record library. For cases where treatment fails, doctors are tasked to reflect and

³<https://www.news-medical.net/medical>

⁴<https://www.merckmanuals.com/professional>

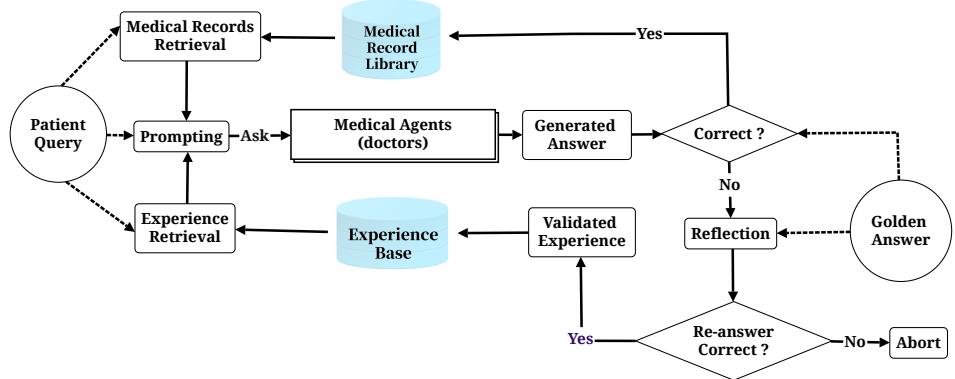


Fig. 5. The overview of the MedAgent-Zero method. This diagram illustrates the method by which doctors achieve self-evolution: 1) Accumulating examples and summarizing experience; 2) Adding correct responses directly to the example library; 3) Summarizing experience from incorrect responses and retesting them; 4) Incorporating successful experience into the experience library after further abstraction; 5) Utilizing both libraries to retrieve the most similar content for reasoning during the inference process.

analyze the reasons for diagnostic inaccuracies and distill a guiding principle to be used as a cautionary reminder for subsequent treatment processes. The construction details will be introduced in Sections 4.3.1 and 4.3.2.

In the course of patient treatment, we employ dense retrievers to retrieve related historical medical records and guiding principles, assisting doctors in delivering better patient care. As experience and records are accrued, they are actively applied, with both the medical record library and the experience base being perpetually updated.

4.3.1 Medical Record Library Building. In the process of administering treatment, it is highly beneficial for doctors to consult and reference previously validated medical records. These medical records contain abundant knowledge and demonstrate the rationale behind accurate and adequate responses to diverse medical conditions. Therefore, we propose to build a medical record library for doctor agents to sharpen their medical abilities, including historical medical records from hospital practices and exemplar cases from medical documents. The library is structured in the format of question-answer pairs, where the question details the medical condition requiring decision-making, and the answer contains the validated response and chain of thoughts.

As shown in the upper part of Figure 5, for each generated answer from doctor agents, the question-answer pair will be added to the medical record library if the answer is correct. When a new query is coming, we will search for related records from the library based on dense retrieval techniques use [15]. As there are various medical tasks, we prefer that each task should have its private medical record library to avoid irrelevant records utilization.

For the simulated hospital tasks, we aim to enable medical professional agents to gain clinical experience from historical medical records. We segment each hospital visit into three parts: *Examination*, *Diagnosis*, and *Treatment* (details in 4.1). To construct the query part for

retrieving in distinct tasks, we record symptoms for the *Examination* stage; symptoms and examination results for the *Diagnosis* stage; symptoms, examination results, and possible treatment plan for the *Treatment* stage. In the answer part, we record the chain of thoughts corresponding to a specific query.

The generated problems from medical documents are utilized to construct an extra record library for doctor agents to learn general medical knowledge. The correct question-answer pairs are also added to the learning record library, and will be adopted for general medical-related tasks. Note that we end up with two growing record libraries, one practice-based and one learning-based, that doctors can draw on at any time to diagnose patients or answer other medical questions.

4.3.2 Experience Base Expanding. Learning from diagnostic errors is also crucial for the growth of doctors. We believe that LLM-powered medical professional agents can engage in self-reflection from these errors, distilling relevant principles (experience) to ensure correct diagnoses when encountering similar issues in future cases.

We draw inspiration from a previous study [33] to allow doctor agents to learn from failures. As shown in the below part of Figure 5, if the answer is wrong, the agent will reflect the initial problem, generated answer, and golden answer to summarize reusable principles. All principles generated are subject to a validation process. Upon generation, the principle is integrated into the original question which was initially answered incorrectly, allowing medical professional agents to re-diagnose. Only if the diagnosis is correct will the principle be added to the experience base.

To eliminate the influence of noise and maximize the utilization of the experience base, we incorporate additional judgment when utilizing experience. This judgment involves evaluating whether the top- K experience retrieved based on semantic similarity are helpful for the treating process. Helpful experience will be incorporated into the prompt, while unhelpful experience will be excluded.

Our framework for utilizing and accumulating experience is dynamic. Specifically, once medical professional agents are initialized, they begin learning from errors and continuously accumulating experience. Once there are at least the top- k experience in the experience base, they start utilizing experience, mirroring real-world doctors' practices. The experience bases of distinct tasks are accumulated separately as the medical record library.

4.4 Inference

Based on the Medical Record Library and Experience Base introduced above, we enhance the prompt for the medical agents by using successful medical records and validated experience retrieved from them. For medical records, we get the most similar ones by comparing existing queries in the Medical Record Library with the current query. Then, the chosen records combined by query and answer are arranged for few-shot examples in the prompt.

For experience, in Real-World Evaluation, we identify the most relevant ones by calculating the similarity between the experience itself and the current query. In Simulation Evaluation, for the Examination Decision task, we calculate the similarity between the symptoms of the current patient and those of previous patients in Experience Base. For the Diagnosis and Treatment Plan tasks, we calculate the similarity between the current query and the queries of experience to retrieve experience. Before adding these retrieved experience to the prompt, we will judge whether they are helpful in answering the question with LLMs. Some valued experience is selected and others are dropped.

Both records and experience are retrieved using cosine similarity and text is embedded into vector space by "text-embedding-ada-002" model provided by OpenAI⁵.

5 SIMULATION EVALUATION

5.1 Experimental Settings

5.1.1 Dataset. As introduced in Section 4.2.1, leveraging insights from GPT Turbo-3.5 and our comprehensive medical database, we can dynamically generate detailed patient profiles and complete medical histories. We use the Simulated Medical Dataset for evaluation, which includes personal information such as names, ages, and genders, along with medical details like current diseases, their severity, physiological symptoms, and necessary diagnostic tests. To improve the reproducibility of our experiments, we have constructed a balanced training dataset consisting of 10,000 instances and a separate test set of 500 instances. Table 5 shows the detailed distributions of these records in the Appendix.

5.1.2 Evaluation Metrics. We propose an evaluation strategy to evaluate doctor agents in Agent Hospital with three primary capabilities: medical assessment, diagnosis, and treatment recommendation. Firstly, each agent is tested on the medical examination task, which aims to select a medical examination from sixteen candidates based on patient symptoms. The selection is deemed correct if it matches any item in the predefined ground truth of suitable medical examinations for the patient. Secondly, the agent's disease diagnostic skills are evaluated by prompting it to choose the correct one from eight diseases, given both the patient's symptoms and medical examination results. Lastly, the agent is asked to recommend the most suitable treatment plan based on the symptoms and examination results of the patient, and select from the three treatment plans tailored to mild, moderate, and severe conditions.

Accuracy is utilized as the metric to evaluate the agent's performance on each task. Note these metrics are designed to be integrally linked to replicate the sequential decision-making process observed in the real world, as the outcome at the previous stage impacts the next. Specially, as the treatment plan is highly related to the diagnosis result, once the diagnosis is wrong, the treatment result is seen as incorrect.

5.1.3 Implementation Details. For each query, the number of utilized medical records and principles after retrieving is set to 3, i.e., only the top 3 relevant experience and records are adopted in the prompt. The medical record library and experience base are training from empty, and will be updated dynamically during training to support further decisions. So the training of a doctor agent is similar to a new doctor improving her/his medical skills by practicing. All of our simulation experiments are based on gpt-3.5-turbo-1106 API.

5.2 Experimental Results

Based on the described settings, we conduct experiments to verify the effectiveness of the proposed MedAgent-Zero. The accuracy changes during training across three tasks are shown in Figure 6, and the accuracy changes on the test set during the training process are shown in Figure 7 (evaluated every 100 training samples).

From the experimental results, we have the following conclusions. Firstly, the proposed MedAgent-Zero strategy effectively enhances doctor Agents on the three tasks, with the cumulative accuracy on 10,000 training samples showing a continuous increase. The best performances of examination, diagnosis, and treatment are 88%, 95.6%, and 77.6%, respectively.

⁵<https://platform.openai.com/docs/guides/embeddings/embedding-models>

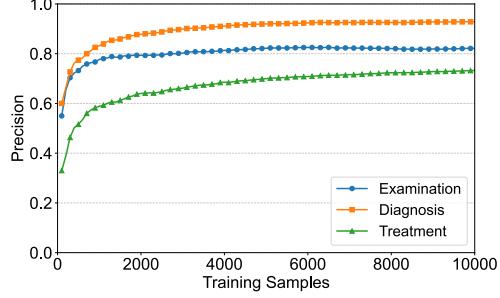


Fig. 6. MedAgent-Zero training performance on the three tasks with the increase of training samples.

It shows our agent evolves during the training phase, just as human doctors become experienced after treating thousands of patients. Furthermore, agent evolution is more efficient than human, as human doctors may take over two years to treat ten thousand patients.

Secondly, The original GPT-3.5 performs poorly on the three medical tasks (accuracy without training samples), with the precisions on the test set all below 0.4. However, after training, the test set performance of doctor agents improved rapidly. Although there were fluctuations, the accuracy of the diagnosis and treatment tasks continued to increase. The performance on the examination task showed greater variability, possibly due to the complexity of the task (each question may have multiple correct answers).

Thirdly, although accuracy also continuously improves when training is conducted solely using the medical record library or experience base, the performance on the test set is not as good as that achieved using both.

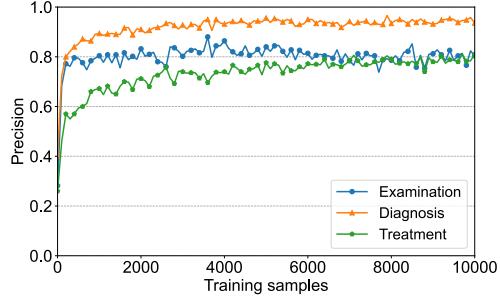


Fig. 7. MedAgent-Zero test performance on the three tasks with the increase of training samples.

5.3 Further Analyses

5.3.1 *Performance on distinct diseases.* To further verify the performance of MedAgent-Zero on distinct diseases, we draw Figures 8, 9, and 10 to demonstrate the examination, diagnosis, and treatment accuracy, respectively.

According to the results, a common trend is that more training samples contribute to better performance on distinct diseases of different tasks. Besides, influenza B is a disease that hard to deal with, as our doctor agent achieves the worst performance on this disease across all three tasks. While bronchial asthma is easy to handle, the examination and treatment precision of which is the highest.

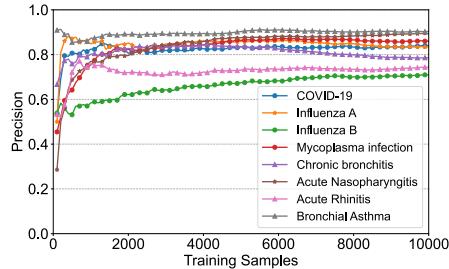


Fig. 8. Examination accuracy of distinct diseases with the increase of training samples.

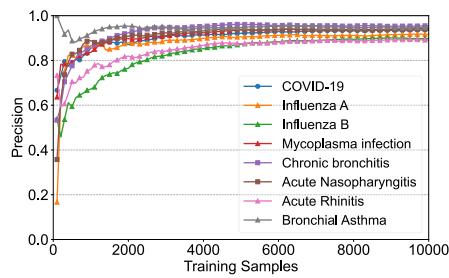


Fig. 9. Diagnosis accuracy of distinct diseases with the increase of training samples.

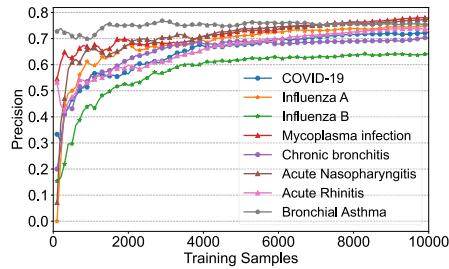


Fig. 10. Treatment accuracy of distinct diseases with the increase of training samples.

The experiments show that with the accumulation of experience and examples, the performance of examination, diagnosis, and treatment of different diseases all improve rapidly. Our system can continuously iterate and evolve through practice. The reason for the initially higher accuracy of bronchial asthma may be caused by the output bias. Among the first 100 training samples, there are 11 instances of bronchial asthma, with 10, 11, and 8 correct answers in the examination, diagnosis, and treatment, respectively. As the number of training samples increases, the accuracy of bronchial asthma tends to normalize.

5.3.2 Experience Accumulation. Figures 11, 12, and 13 show the accumulation of validated experience and wrong answers on the examination, diagnosis, and treatment tasks, respectively. When the training samples increase, the number of experience and wrong answers all increase slowly. As shown in the figures, the curve of the experience is below the curve of the wrong answers. The reason is that our agent is unable to reflect experience from all failures. Besides, the diagnosis experience is easier to accumulate than the other tasks.

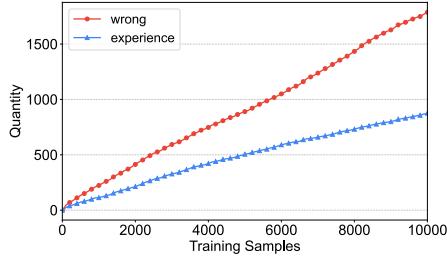


Fig. 11. The number of accumulated principles and wrong answers of examination with the increase of training samples.

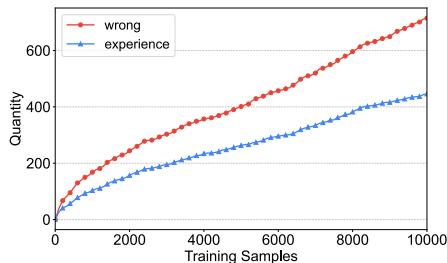


Fig. 12. The number of accumulated principles and wrong answers of diagnosis with the increase of training samples.

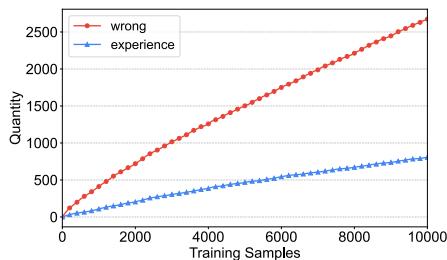


Fig. 13. The number of accumulated principles and wrong answers of treatment with the increase of training samples.

Note that all correct samples are added to the medical records library, so the number and trends are similar to Figure 6.

5.3.3 Case Study. To further analyze the usefulness of the medical record library and experience base, we conduct a case study here. Table 1 illustrates the performance of the experience base, medical record library separately, and MedAgent-Zero on the three tasks of a patient case. At first, all answers are incorrect without experience or medical records. When only the experience is added, the examination and diagnosis responses are correct, but the treatment response is incorrect. When only the relevant medical records are added, the examination response is incorrect, while diagnosis and treatment level responses are correct. Finally, when both the experience and records are added simultaneously, the doctor can provide correct responses for all three tasks.

Symptoms of the patient: The patient feels unwell with the following symptoms: discomfort in the throat, dry cough, dry mouth, these symptoms have lasted for 2 weeks.

Examination Task:

- **Retrieved experience:** In this case, the patient's persistent symptoms of discomfort in the throat, dry cough, and dry mouth for 1-2 weeks indicate the need for **thorough respiratory and blood-related examinations** to identify the cause of the illness. Each symptom may point to different potential conditions, so a thorough examination is necessary.
- **Retrieved records:** Top-3 relevant examples retrieved from record library all recommend blood test.

MedAgent-Zero: Thought: Given the persistent discomfort in the throat, dry cough, and dry mouth for 2 weeks, further examination focused on **respiratory function and blood-related testing** is necessary to identify the cause of the illness. Examination: blood test.

Diagnosis Task:

- **Retrieved experience:** The generated answer of Acute Nasopharyngitis was incorrect, as the symptoms indicated in the patient's examination are actually **more consistent with chronic bronchitis**. In this case, the blood test results pointed towards a bacterial infection, which would **align more with chronic bronchitis** than acute nasopharyngitis.
- **Retrieved records:** Top-3 relevant examples retrieved from record library all recommend chronic bronchitis.

MedAgent-Zero: Thought: **The patient's symptoms are consistent with chronic respiratory conditions**, as well as **the blood test results indicating a possible bacterial infection**. These factors strongly support the diagnosis of chronic bronchitis. Disease name: Chronic bronchitis.

Treatment Task:

- **Retrieved experience:** The symptoms described are discomfort in the throat, dry cough, and dry mouth, which **do not indicate severe respiratory issues**. It is essential to **match the level of treatment to the specific needs of the patient**, rather than assuming a higher level of treatment is necessary based solely on the presence of an infection.
- **Retrieved records:** Top-3 relevant examples retrieved from record library all recommend a mild treatment plan.

MedAgent-Zero: Thought: The symptoms described are discomfort in the throat, dry cough, and dry mouth, which **do not indicate severe respiratory issues**. The blood test suggests a **mild to moderate** bacterial infection. Treatment level: mild.

Fig. 14. Case study: Some key generated results of MedAgent-Zero on the three tasks. As shown in Table 1, only using the medical records library (retrieved records) or experience base (retrieved experience) will result in errors.

Figure 14 shows that with the addition of the experience and records, MedAgent-Zero answers correctly on all three tasks, with the retrieved experience and examples for each task contributing to the final answer. The experience base provides potential incorrect perspectives and key points that need special attention when answering. The records library provides the top- k most relevant reference answers based on the accumulated medical records. This case illustrates that both the experience base and medical record library are helpful for completing all three tasks, and they can complement each other to achieve better results.

Table 1. The prediction results of distinct strategies on the patient case in Figure 14.

Strategy	Examination	Diagnosis	Treatment
Only Experience Base	✓	✓	✗
Only Medical Records	✗	✓	✓
MedAgent-Zero	✓	✓	✓

6 REAL-WORLD EVALUATION

In Section 5, the effectiveness of the proposed MedAgent-Zero is verified in the simulated medical datasets. Although the medical records may be only helpful to the three types of medical tasks, we want to verify if the accumulated medical experience is useful to real-world medical datasets.

6.1 Task Definition

To evaluate the effectiveness of our evolution setting in realistic problems, we compare our method with other baselines on MedQA [7], a widely recognized and authoritative dataset for medical answering. MedQA includes questions in the format of multiple choice, mirroring the Medical Licensing Examination questions used to evaluate medical specialist competency. We chose the USMLE-style part consists of 1,273 questions, each offering four possible choice answers. Higher accuracy on this dataset demonstrates the better medical capability of a doctor agent.

Since we only generated diseases about the respiratory tract, such as COVID-19 and Influenza A, in our Agent Hospital at present, we only select related questions from the MedQA test set by GPT-3.5. Finally, there is a subset consisting of 72 questions, and our further experiments are conducted on this dataset.

6.2 Implementation Details

For inference, we adopt MedAgent-Zero by combining the medical record library from medical document learning (Section 4.2.2) and experience base from patient treatment in Agent Hospital (Section 4.2.1). The reason for this combination is that medical documents align the reality well so derived successful records contain accurate medical information, while experience from patient treatment is more specific and can be accumulated infinitely. The Experience Base includes experience concluded from the doctor agent’s wrong examination and diagnosis as these processes contain the most valuable medical knowledge.

We tune the number of experiences and records from the range between zero and ten, and the top- k of validated experience and most similar records chosen to form the inference prompt are both 2. For the accumulated Experience Base, we select the keyframe with every 2,000 patient cases diagnosed in Agent Hospital. We end up with 8,000 patient cases in total

Table 2. Main results on the respiratory disease dataset using GPT-3.5 and GPT-4 as the backbones.

Method	GPT-3.5	GPT-4
Vanilla	69.44	86.11
CoT [26]	73.61	86.11
MedAgents [24]	73.61	91.67
Medprompt [15]	81.94	90.28
MedAgent-Zero	84.72	93.06

as the performance is higher than all baselines. The LLM versions of the doctor agent are gpt-3.5-turbo-1106 and gpt-4-1106-preview.

6.3 Experimental Results

Experimental results are shown in Table 2. First, MedAgent-Zero achieves the best performance on the respiratory disease dataset, outperforming the SOTA method Medprompt [15] by 2.78% when using GPT-3.5 and outperforming the SOTA method MedAgents [24] by 1.39% when using GPT-4. The result validates that our model is helpful to agent evolution with only simulated and medical documents without any training samples from the MedQA, which effectively enhances the medical capacity of doctor agents. Second, The best performance of MedAgent-Zero is 93.06% based on GPT-4, which outperforms human experts in the MedQA dataset (around 87%) [14]. Third, GPT-4 based medical agents show prior performance than GPT-3.5 based in vanilla of any other methods, showing that GPT-4 is more powerful in the medical domain.

To summarize, our experimental results show that when evolved within the Agent Hospital by the MedAgent-Zero strategy, the medical agents can learn from simulated patients & medical documents and summarize helpful experience to achieve the best performance on real-world medical examinations, even without no manually labelled data.

6.4 Further Analyses

6.4.1 Ablation Studies. To further verify the effectiveness of the proposed two modules, we conducted an ablation study and the results are summarized in Table 3. First, MedAgent-Zero, where both the medical record library and experience base are utilized, achieves the best performance, showing that both modules are helpful. Second, the result with record library and experience base together is higher than use records or experience alone by 1.39% and 2.78% when inference with GPT-4, showing that the two parts have a synergistic effect and record library has a great influence on the final results. Third, we can find that inference with only the medical record library or experience base outperforms the CoT method by 8.33% and 2.78% when using GPT-3.5, respectively, which demonstrate the effectiveness of the two modules of MedAgent-Zero. It keeps the same trend when we use GPT-4 to infer.

6.4.2 Experience Accumulation Analysis. To demonstrate the effect of the experience accumulated in treating simulated patients on the Agent Hospital, we conduct experiments with experience that is summarized from different numbers of simulated patients. We keep the experiment setting using the top 2 experience and top 2 successful records to help answer the

Table 3. Ablation study of MedAgent-Zero, where ‘Record’ means using the medical record library and ‘Experience’ means using the Experience Base.

Model	Record	Experience	Accuracy
GPT-3.5	✓	-	81.94
	-	✓	79.39
	✓	✓	84.72
GPT-4	✓	-	91.67
	-	✓	90.28
	✓	✓	93.06

Table 4. The performance of MedAgent-Zero with accumulated experience from distinct numbers of simulated patients.

Model	0	2,000	4,000	6,000	8,000
GPT-3.5	73.61	79.16	77.78	81.94	84.72
GPT-4	86.11	91.67	88.89	90.28	93.06

MedQA tests. We choose some typical case points ranging from 0 to 8,000 cases. The results are shown in Table 4.

We can find that with the accumulation of patient cases to expand the experience base, the accuracy is getting higher in general. The performance with experience base accumulated using 8,000 cases is higher than those using 2,000/4,000/6,000 cases whether with GPT-3.5 or GPT-4. It is worth noting that the larger experience base is not always better, as we find that there is a marked drop between 2,000 cases and 4,000 cases. The reason may be that the distribution of cases from this period is significantly different from the respiratory disease dataset and some unhelpful experience is concluded. However, with more and more experience of high quality, the overall trend is getting better and better.

7 DISCUSSIONS

7.1 Main Findings

Firstly, our study verifies the possibility of self-evolution within Agent Hospital, providing a new approach for the study of medical LLMs/agents. This insight demonstrates a new way for agent evolution in simulation environments, where agents can improve their skills without human intervention. Secondly, the proposed MedAgent-Zero strategy offers a new method for parameter-free agent evolution without knowledge. By enabling agents to refine and expand their expertise through continuous interaction and feedback loops within the simulation, the strategy enhances their ability without any manually labelled data. Thirdly, Agent Hospital demonstrates promising scalability and interactivity, making it suitable for more complex medical simulation experiments. Its design allows for extensive customization and adjustment, enabling researchers to test a variety of scenarios and interactions within the healthcare domain.

7.2 Limitations

There are still some limitations to our work: 1) Only GPT-3.5 is adopted as the simulator for our Agent Hospital and evaluations. 2) Due to the interaction between agents and their

evolution involving API calls, the efficiency of our hospital is constrained by LLM generation. 3) Although the health records and examination results for each patient are generated to mimic real electronic health records without domain knowledge, there may be still some discrepancies to real-world records.

7.3 Future Work

Our future plans for the Agent Hospital mainly include follows: 1) Expanding the range of diseases covered in the simulation and extending into more medical departments, aiming to mirror the comprehensive services provided by a real hospital for further studies. 2) Enhancing the society simulation aspects of agents, such as incorporating a full promotion system for medical professionals, changing the distribution of the disease with time, and incorporating the history medical records of patients. These modifications will add depth to the interactions and decisions made by the agents. 3) Optimizing the selection and implementation of the base LLM model, aiming to achieve more efficient execution of the entire simulation process by leveraging powerful and open-source models.

8 CONCLUSIONS

In this paper, we construct a simulacrum of hospital for medical scenarios based on LLM and agent technology, which is named Agent Hospital. Agent Hospital not only includes two types of roles (medical professionals and patient agents) and dozens of specific agents, but also covers both in-hospital processes like triage, registration, consultation, examination, and treatment planning, as well as out-of-hospital stages such as illness and recovery. In this Agent Hospital, we propose the MedAgent-Zero strategy for the evolution of medical agents, which is parameter-free and knowledge-free, allowing for infinite agent training through simulated patients. This strategy primarily incorporates a medical record library and an experience base, enabling the accumulation of experience from correct and failed treatments as human doctors. On the simulated patient dataset, we observe that as the patient records increased, the accuracy of the doctor agents in examination, diagnosis, and treatment continuously improved. The doctor agent is able to complete the diagnosis and treatment of tens of thousands of patients within a few days, which would typically take at least two years for a human doctor. Furthermore, we find that the experience accumulated in Agent Hospital can significantly enhance the accuracy of doctor agents in a subset of the MedQA dataset, which even achieves state-of-the-art performance. Our study verifies that real-world simulation with a designed strategy can enhance the performance of LLM agents on specific tasks.

REFERENCES

- [1] Arian Askari, Roxana Petcu, Chuan Meng, Mohammad Aliannejadi, Amin Abolghasemi, Evangelos Kanoulas, and Suzan Verberne. 2024. Self-seeding and Multi-intent Self-instructing LLMs for Generating Intent-aware Information-Seeking dialogs. *arXiv preprint arXiv:2402.11633* (2024).
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [3] Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233* (2023).
- [4] Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023. S3: Social-network Simulation System with Large Language Model-Empowered Agents. *arXiv:2307.14984* [cs.SI]

- [5] Wenyue Hua, Lizhou Fan, Lingyao Li, Kai Mei, Jianchao Ji, Yingqiang Ge, Libby Hemphill, and Yongfeng Zhang. 2024. War and Peace (WarAgent): Large Language Model-based Multi-Agent Simulation of World Wars. arXiv:2311.17227 [cs.AI]
- [6] Shuyang Jiang, Yuhao Wang, and Yu Wang. 2023. Selfevolve: A code evolution framework via large language models. *arXiv preprint arXiv:2306.02907* (2023).
- [7] Di Jin, Eileen Pan, Nassim Oufattolle, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What Disease does this Patient Have? A Large-scale Open Domain Question Answering Dataset from Medical Exams. arXiv:2009.13081 [cs.CL]
- [8] Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu Chan, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. 2024. Adaptive Collaboration Strategy for LLMs in Medical Decision Making. *arXiv preprint arXiv:2404.15155* (2024).
- [9] Stefan Kirn, Christian Anhalt, Helmut Krcmar, and Andreas Schweiger. 2006. Agent. Hospital–health care applications of intelligent agents. *Multiagent Engineering: Theory and Applications in Enterprises* (2006), 199–220.
- [10] Lanjuan Li and Hong Ren. 2013. *Infectious Diseases* (8 ed.). People’s Medical Publishing House.
- [11] Nian Li, Chen Gao, Yong Li, and Qingmin Liao. 2023. Large Language Model-Empowered Agents for Simulating Macroeconomic Activities. arXiv:2310.10436 [cs.AI]
- [12] Yang Li, Yangyang Yu, Haohang Li, Zhi Chen, and Khaldoun Khashanah. 2023. TradingGPT: Multi-Agent System with Layered Memory and Distinct Characters for Enhanced Financial Trading Performance. arXiv:2309.03736 [q-fin.PM]
- [13] Yuan Li, Yixuan Zhang, and Lichao Sun. 2023. MetaAgents: Simulating Interactions of Human Behaviors for LLM-based Task-oriented Coordination via Collaborative Generative Agents. arXiv:2310.06500 [cs.AI]
- [14] Valentin Liévin, Christoffer Egeberg Hother, Andreas Geert Motzfeldt, and Ole Winther. 2023. Can large language models reason about medical questions? arXiv:2207.08143 [cs.CL]
- [15] Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, Renqian Luo, Scott Mayer McKinney, Robert Osazuwa Ness, Hoifung Poon, Tao Qin, Naoto Usuyama, Chris White, and Eric Horvitz. 2023. Can Generalist Foundation Models Outcompete Special-Purpose Tuning? Case Study in Medicine. arXiv:2311.16452 [cs.CL]
- [16] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [17] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–22.
- [18] Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. *arXiv preprint arXiv:2307.07924* (2023).
- [19] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research* 21, 140 (2020), 1–67.
- [20] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [21] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. 2017. Mastering the game of go without human knowledge. *nature* 550, 7676 (2017), 354–359.
- [22] Wolfgang Stammer, Felix Friedrich, David Steinmann, Hikaru Shindo, and Kristian Kersting. 2023. Learning by Self-Explaining. *arXiv preprint arXiv:2309.08395* (2023).
- [23] Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2024. Principle-driven self-alignment of language models from scratch with minimal human supervision. *Advances in Neural Information Processing Systems* 36 (2024).
- [24] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Yilun Zhao, Xingyao Zhang, Arman Cohan, and Mark Gerstein. 2023. Medagents: Large language models as collaborators for zero-shot medical reasoning. *arXiv preprint arXiv:2311.10537* (2023).
- [25] Lei Wang, Jingsen Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Recagent: A novel simulation paradigm for recommender systems. *arXiv preprint arXiv:2306.02552* (2023).
- [26] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances*

- in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., 24824–24837. https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf
- [27] Ross Williams, Niyousha Hosseinichimeh, Aritra Majumdar, and Navid Ghaffarzadegan. 2023. Epidemic Modeling with Generative Agents. *arXiv:2307.04986 [cs.AI]*
 - [28] Bushi Xiao, Ziyuan Yin, and Zixuan Shan. 2023. Simulating Public Administration Crisis: A Novel Generative Agent-Based Simulation System to Lower Technology Barriers in Social Science Research. *arXiv:2311.06957 [cs.CY]*
 - [29] Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. 2023. Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint arXiv:2309.04658* (2023).
 - [30] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems* 35 (2022), 15476–15488.
 - [31] An Zhang, Leheng Sheng, Yuxin Chen, Hao Li, Yang Deng, Xiang Wang, and Tat-Seng Chua. 2023. On Generative Agents in Recommendation. *arXiv:2310.10108 [cs.IR]*
 - [32] Junjie Zhang, Yupeng Hou, Ruobing Xie, Wenqi Sun, Julian McAuley, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. 2023. AgentCF: Collaborative Learning with Autonomous Language Agents for Recommender Systems. *arXiv:2310.09233 [cs.IR]*
 - [33] Tianjun Zhang, Aman Madaan, Luyu Gao, Steven Zheng, Swaroop Mishra, Yiming Yang, Niket Tandon, and Uri Alon. 2024. In-Context Principle Learning from Mistakes. *arXiv:2402.05403 [cs.CL]*
 - [34] Qinlin Zhao, Jindong Wang, Yixuan Zhang, Yiqiao Jin, Kaijie Zhu, Hao Chen, and Xing Xie. 2023. CompeteAI: Understanding the Competition Behaviors in Large Language Model-based Agents. *arXiv:2310.17512 [cs.AI]*
 - [35] Li Zhong, Zilong Wang, and Jingbo Shang. 2024. LDB: A Large Language Model Debugger via Verifying Runtime Execution Step-by-step. *arXiv preprint arXiv:2402.16906* (2024).

A CONTRIBUTIONS OF AUTHORS

Junkai Li is responsible for the design and implementation of the MedAgent-Zero, and authored the first version of Section 6. Siyu Wang is in charge of in-hospital experiments in Agent Hospital and wrote the first version of Section 5. Meng Zhang is responsible for building the Agent Hospital and contributed to part of Section 3. Weitao Li focused on the design and experimentation of the MedAgent-Zero experience accumulation mechanism. Yunghwei Lai is responsible for disease simulation and medical record generation, collaborating with Xinhui Kang to write the first version of Sections 2 and 3. Weizhi Ma managed the overall project progress, wrote the first version of the remaining sections, and revised the manuscript. Yang Liu is responsible for the overall research design and manuscript revisions.

B INVOLVED DISEASE KNOWLEDGE

B.1 Diseases

The information on COVID-19 depicted in Figure 15 is essential for constructing the health records of patients affected by COVID-19. We have currently compiled detailed medical information on eight diseases within our knowledge base to construct the simulated hospital environment. We would include a broader array of diseases to ensure the diversity and accuracy of patient information in the future.

COVID-19

Symptoms: dry throat, sore throat, fever, smell taste loss, runny nose, the central nervous system involvement, difficulty in breathing, hypoxemia, acute respiratory distress syndrome, sepsis shock, refractory metabolic acidosis, coagulopathy, and multiple organ failure.

Examination Results:

- **Blood Test:** In the early stage of the disease, the total number of peripheral blood white blood cells was normal or decreased, and the lymphocyte count was decreased. Some patients may have increased liver enzymes, lactate dehydrogenase, muscle enzymes, myoglobin, troponin and ferritin. In most patients, C-reactive protein (CRP) and erythrocyte sedimentation rate were increased, and procalcitonin was normal. In severe and critical patients, D-dimer was increased, peripheral blood lymphocytes were progressively decreased, and inflammatory factors were increased.
- **Chest X-ray Exam:** Chest imaging examination showed multiple small patchy shadows and interstitial changes in the early stage, which were obvious in the outer lung zone. Then it develops into multiple ground-glass opacities and infiltrations in both lungs. In severe cases, lung consolidation may occur, and pleural effusion is rare. In MIS-C, enlarged heart shadow and pulmonary edema are seen in patients with cardiac dysfunction.

Treatment Plan:

- **Mild:** rest in bed, strengthen supportive treatment, ensure adequate energy and protein intake, supplement vitamins, trace elements and other nutrients; Timely administration of ritonavir tablets or ambavir and romisivir injection.
- **Moderate:** timely physical cooling, drug antipyretic, prone position treatment, timely delivery of azvudine, monolavir capsule drug treatment.
- **Severe:** Treatment was given in the standard prone position for no less than 12 hours per day. Respiratory support, circulatory support, timely administration of intravenous human immunoglobulin for COVID-19.

Fig. 15. Utilized Medical Knowledge about the COVID-19.

B.2 Data Distribution

The patient case data in our study are evenly distributed among several diseases, including Acute Nasopharyngitis, Acute Rhinitis, Bronchial Asthma, Chronic Bronchitis, COVID-19, Influenza A, Influenza B, and Mycoplasma Infection. Detailed distributions of this data are illustrated in the subsequent table and figure.

Table 5. Distributions of the training set and the test set.

Disease	Train Dataset				Test Dataset			
	Mild	Moderate	Severe	Total	Mild	Moderate	Severe	Total
Acute Nasopharyngitis	417	418	418	1,253	20	20	20	60
Acute Rhinitis	416	415	417	1,248	21	21	21	63
Bronchial Asthma	417	417	417	1,251	22	22	22	66
Chronic bronchitis	417	417	416	1,250	21	21	21	63
COVID-19	417	417	416	1,250	21	21	22	64
Influenza A	416	417	417	1,250	21	21	21	63
Influenza B	417	416	417	1,250	20	20	20	60
Mycoplasma infection	415	418	415	1,248	20	20	21	61
Total	3,332	3,335	3,333	10,000	164	164	166	500

B.3 An Example of electronic Health Records of a Patient

Personal Information

Name: Kenneth Morgan

Age: 42

Gender: Male

Medical History: Diabetes, Chronic obstructive pulmonary disease

Disease Information

Disease: Acute Nasopharyngitis

Severity Level: Severe

Symptoms: Cough event, high fever, difficulty in breathing, acute muscle pain, complete loss of smell and taste, sore throat

Duration: Symptoms have been escalating rapidly over the past 48 hours

Examination Results

Blood Test:

ALT (Alanine Aminotransferase): 45 IU/L (7 – 40)
 AST (Aspartate Aminotransferase): 50 IU/L (13 – 35)
 Urea: 7.0 mmol/L (2.6 – 8.8)
 Creatinine (Cr): 95 µmol/L (41 – 81)
 Triglycerides (TG): 1.5 mmol/L (< 1.7)
 Total Cholesterol (TC): 6.0 mmol/L (< 5.18)
 Hepatitis B Surface Antigen (HBsAG): Negative
 HIV Antibody Test (anti-HIV): Negative
 Syphilis Test (RPR): Negative
 White Blood Cell Ct. (WBC): $3.0 \times 10^9/L$ (3.5 – 9.5)
 Red Blood Cell Ct. (RBC): $3.8 \times 10^{12}/L$ (3.5 – 5.5)
 Hematocrit (Hct): 35% (35 – 50)
 Hemoglobin (Hb): 110 g/L (115 – 150)
 Platelet Ct. (PLT): $200 \times 10^9/L$ (125 – 350)
 Lymphocyte Percentage (LYMPH%): 15% (20 – 50)
 Neutrophil Percentage (NEUT%): 80% (40 – 75)
 Lymphocyte Abs. Ct. (LYMPH#): $0.45 \times 10^9/L$ (1.3–3.2)
 Neutrophil Abs. Ct. (NEUT#): $2.4 \times 10^9/L$ (1.8 – 6.3)
 Monocyte Abs. Ct. (MONO#): $0.3 \times 10^9/L$ (0.2 – 1.0)
 Monocyte Percentage (MONO%): 10% (3 – 10)
 Eosinophil Abs. Ct. (EO#): $0.02 \times 10^9/L$ (0.02 – 0.52)
 Eosinophil Percentage (EO%): 0.7% (0.4 – 8.0)
 Basophil Abs. Ct. (BASO#): $0.01 \times 10^9/L$ (0 – 0.06)
 Basophil Percentage (BASO%): 0.3% (0 – 1)
 Mean Platelet Volume (MPV): 11 fL (9 – 13)
 Lactate Dehydrogenase: 250 U/L (135 – 225)
 Muscle Enzymes (CK): 200 U/L (22-198 for males)
 Myoglobin: 80 ng/mL (< 90 ng/mL)
 Troponin I: 20 ng/L (< 14 ng/L)
 Ferritin: 600 ng/mL (20-500 for males)
 CRP: 50 mg/L (< 3 mg/L)
 ESR: 40 mm/hr (0 – 20)
 Procalcitonin: 0.5 ng/mL (< 0.5)

D-dimer: 1.0 mg/L FEU (0 – 0.5)

Rh Type: Positive

ABO Group: O

Specific antigen: SARS-CoV-2 Nucleo-capsid

Blood Silver Level: 60 µg/L (50 – 150)

Chest X-ray Exam: Lung consolidation with bilateral pleural effusion

Chest Computerized Tomography: Multiple ground-glass opacities and infiltrations in both lungs

Serological Diagnosis: Normal

Viral Antigen Detection: Negative

Allergen Test: Normal

Bacterial Culture of Nasal Secretions: Normal

Respiratory Function Test: Severely impaired

Sputum Examination: Presence of viral particles

Nasopharyngeal Examination: Inflammation and edema

Serum Antibody Test: Positive for SARS-CoV-2 antibodies

Pulmonary Function Test: Impaired gas exchange

Nucleic Acid Amplification Test: Positive for SARS-CoV-2

Eosinophil Count in Sputum: Abnormal

Oral Pharyngeal Examination: Ulcerations and lesions

Nasal Endoscopy: Mucosal inflammation and congestion

Fig. 16. A complete generated medical record of the resident agent.