

3MDBench: Medical Multimodal Multi-agent Dialogue Benchmark

Anonymous ACL submission

Abstract

Though Large Vision-Language Models (LVLMs) are being actively explored in medicine, their ability to conduct telemedicine consultations combining accurate diagnosis with professional dialogue remains underexplored. In this paper, we present **3MDBench** (Medical Multimodal Multi-agent Dialogue Benchmark), an open-source framework for simulating and evaluating LVLM-driven telemedical consultations. 3MDBench simulates patient variability through four temperament-based Patient Agents and an Assessor Agent that jointly evaluate diagnostic accuracy and dialogue quality. It includes 3013 cases across 34 diagnoses drawn from real-world telemedicine interactions, combining textual and image-based data. The experimental study compares diagnostic strategies for popular LVLMs, including GPT-4o-mini, LLaVA-3.2-11B-Vision-Instruct, and Qwen2-VL-7B-Instruct. We demonstrate that multimodal dialogue with internal reasoning improves F1 score by 6.5% over non-dialogue settings, highlighting the importance of context-aware, information-seeking questioning. Moreover, injecting predictions from a diagnostic convolutional network into the LVLM’s context boosts F1 by up to 20%. Source code is available at https://anonymous.4open.science/r/3mdbench_acl-0511.

1 Introduction

Telemedicine expands healthcare access and efficiency, especially for underserved populations, by enabling real-time consultations and early diagnosis (Stoltzfus et al., 2023). In these consultations, effective communication is essential in diagnostic accuracy and treatment adherence (Mirzaei and Kashian, 2020; Bhaskar et al., 2020). Large Language Models (LLMs) and Vision-Language Models (LVLMs) further enhance telehealth via real-time analysis (Nwankwo et al., 2024), chronic care

management (Adeghe et al., 2024), and decision support (Perez et al., 2025), including symptom assessment, test interpretation, and patient interaction (Lu et al., 2024b; Mayer et al., 2024).

Despite recent advances in LLMs, effective doctor-patient dialogue remains essential. Diagnostic quality depends on how patients articulate symptoms, shaped by emotional, cognitive, and systemic factors (Amelung et al., 2020; Singh and Sittig, 2015). Miscommunication, fear, and low health literacy delay diagnosis (Heyhoe et al., 2018; Nguyen et al., 2024), while temperament influences responsiveness, expressiveness, and trust (Graedon and Graedon, 2014; Meyer et al., 2013). Existing LLM and LVLM benchmarks estimating telemedicine quality offer limited realism. Many of them restrict LLMs to multiple-choice tasks (Jin et al., 2020; Kim et al., 2024b), or constrain dialogue with factual, non-reactive templates (Li et al., 2024c; Johri et al., 2024) without using image modality depicting the patient’s symptoms.

To address this, we introduce **3MDBench** (Medical Multimodal Multi-agent Dialogue Benchmark) for evaluating LVLM-based consultations in dynamic, realistic scenarios. We use classical temperament theory, which categorizes individuals into four types (sanguine, choleric, melancholic, and phlegmatic) (Steiner, 1985) shows that personality impacts engagement and medical alliance (Paap et al., 2022; Hanney et al., 2023). Hence, our 3MDBench features a **Patient Agent**, simulating one of four temperament types (Fig. 1), and an **Assessor Agent**, evaluating diagnostic accuracy and communication quality. We select models that best match defined metrics and human annotations to construct these agents. Built on 34 diagnoses obtained from real-world telemedicine consultations using medical image datasets enriched with textual information, 3MDBench supports both textual and image modalities. We benchmark commercial and open-source LVLMs in different dia-

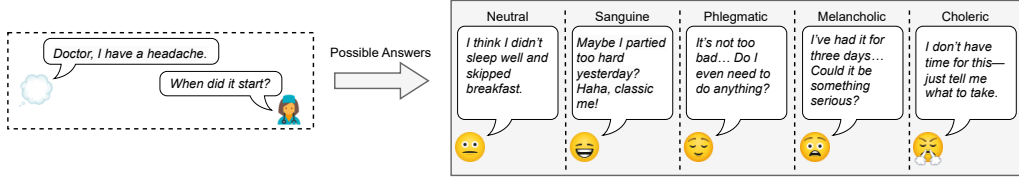


Figure 1: The patient’s response during the dialogue with the doctor depends on their temperament.

logue and non-dialogue settings. Our results show that information-seeking dialogue strategies using medical reasoning and image modality increase the F1 score to 6.5%, highlighting the importance of adaptive, context-aware interaction. Finally, our novel approach incorporating top-3 predictions from a ConvNet trained on diagnosis into the LVLM’s context boosts the F1 score up to 20%.

Our **main contributions** are:

1. We propose **3MDBench**, an open-source benchmark for evaluating medical dialogue systems, with an Assessor Agent measuring diagnostic accuracy and communication quality, and a Patient Agent simulating temperament-based personality-driven doctor-patient conversations.
2. By using our standardized framework for assessing AI-driven medical consultation quality, we extensively compare open-source and state-of-the-art LVLMs.
3. We analyze multiple dialogue strategies for a Doctor Agent: with and without image modality, implementing rationale generation, and using external cues. We demonstrate the importance of image modality and conducting information-seeking conversations with internal reasoning. Moreover, we show that the diagnostic F1 score of the Doctor Agent improves up to 20% by incorporating top-3 predictions from a specially trained diagnostic convolutional neural network into the LVLM.

2 Related Work

Benchmarks for medical LLMs and LVLMs have focused primarily on factual knowledge, evaluating performance on exams, QA tasks, and case-based reasoning (Jin et al., 2020; Kim et al., 2024b; Pal et al., 2022; Jin et al., 2019; Singhal et al., 2023) to assess domain understanding and consistency, but overlook interactive and contextual aspects of diagnosis. Recent work highlights the need

for dialogue-based evaluation, where models must elicit, interpret, and reason over patient-reported symptoms (Goh et al., 2024; Manes et al., 2024; Li et al., 2023b; Han et al., 2023). However, existing benchmarks often rely on scripted interactions or fixed-response patient agents, limiting their reflection of real consultations (Shi et al., 2024).

Multi-agent systems offer a more dynamic alternative, as they simulate collaborative diagnostic workflows, enabling LLMs to interact, reason, and refine decisions over multiple turns (Qiu et al., 2024; Kim et al., 2024a). Several benchmarks employ simulated patients (Mehandru et al., 2024; Li et al., 2024c), but they typically provide only factual, structured inputs, ignoring how emotional state, communication style, or personality affect diagnostic accuracy (Zhu and Wu, 2025). Moreover, image modality, crucial for many diagnostic tasks, is often excluded, despite its significant role in real-world decision-making (Agbareia et al., 2025).

Our work addresses this gap by introducing **3MDBench**, a benchmark that simulates and evaluates telemedicine consultation with a temperament-driven Patient Agent and an Assessor Agent for accuracy and communication quality. Compared to existing benchmarks (Table 1), we capture the variability and complexity of real-world clinical interactions, enabling richer, more patient-aligned evaluation of medical dialogue systems.

3 Proposed 3MDBench

3.1 Data Collection

Diagnoses. To ensure clinical relevance, we analyzed 611K anonymized visits from a large Eastern European provider from May to October 2024, selecting the top 80% most frequent diagnoses. We examined 180 million outpatient records from the same city through 2022 to validate cross-setting consistency. All diagnoses, originally in ICD-10 (Organization, 2004), were standardized using a physician-curated dictionary. The final set comprises 34 diagnoses across five medical domains as shown in Figure 7 of Appendix A.

Table 1: Comparison of 3MDBench with existing medical benchmarks and datasets. The columns are: **T** (Type: Dataset (DS) or Benchmark (BM)), **TD** (Text Data Type: Question-Answer pairs (QA) or Dialogues (D)), **N** (Name of Dataset/Benchmark), **M** (Modality: Text-only (T) or Multimodal (M)), **S** (Size of test part of a Benchmark of full size of a Dataset), **D** (Dialogues present), **A** (Multi-Agent approach used), **P** (Personality modeling used), **CQ** (Consultation quality tested), and **L** (Language of data).

T	TD	N	M	S	D	A	P	CQ	L
DS	D / QA	BianQueCorpus (Chen et al., 2023)	T	2437K	+	-	-	-	CH
DS	D / QA	HealthCareMagic-100k (Li et al., 2023c)	T	100K	+	-	-	-	EN
DS	D	MedDG (Liu et al., 2022)	T	18K	+	-	-	-	CN
DS	D / QA	Psych8k (Yuan et al., 2025)	T	8K	+	-	-	-	EN
DS	D	CMtMedQA (Yang et al., 2023)	T	70K	+	-	-	-	CN
DS	D	Icliniq-10K (Li et al., 2023b)	T	10K	+	-	-	-	EN
DS	D	MedDialog-EN (Zeng et al., 2020)	T	300K	+	-	-	-	EN
DS	D	MedDialog-CN (Zeng et al., 2020)	T	1100K	+	-	-	-	CN
DS	D	IMCS-21 (Chen et al., 2022)	T	811	+	+	-	-	CN
DS	D	NoteChat (Wang et al., 2024)	T	30K	+	+	-	-	EN
BM	QA	Medical-Diff-VQA (Hu et al., 2025)	M	70K	-	-	-	-	EN
BM	QA	PathVQA (He et al., 2020)	M	6K	-	-	-	-	EN
BM	QA	Cholec80-VQA (Twinanda et al., 2016)	M	9K	-	-	-	-	EN
BM	QA	VQA-RAD (Lau et al., 2018)	M	3.5K	-	-	-	-	EN
BM	QA	RadBench (AI, 2024)	M	137K	-	-	-	-	EN
BM	QA	MMMU (H & M) (Yue et al., 2024)	M	11.5K	-	-	-	-	EN
BM	QA	SLAKE (Liu et al., 2021)	M	2K	-	-	-	-	EN
BM	QA	OmniMedVQA (Hu et al., 2024)	M	128K	-	-	-	-	EN
BM	QA	GMAI-MMBench (Chen et al., 2024)	M	26K	-	-	-	-	EN
BM	D	MediQ (Li et al., 2024c)	T	1.2K	+	+	-	-	EN
BM	D	3MDBench (Ours)	M	3K	+	+	+	+	EN

Image Data. We constructed 3MDBench from 6 open-source datasets, primarily from Kaggle^{1 2 3 4}, as well as ISIC Archive images (Cassidy et al., 2022), Google SCIN (Ward et al., 2024), and Fitzpatrick17k (Groh et al., 2021) with supplementing data using Bing Image Search⁵ (Ghosh et al., 2023). We adjusted disease prevalence to match distribution from a major Eastern European telemedicine provider described above to align with real-world telemedicine diagnosis distribution. To ensure sufficient evaluation data and mitigate class imbalance, we set a minimum threshold of 64 images per condition, an empirically determined lower bound based on the maximum number of images available for

certain classes across all sources and the Internet. Obtained images were filtered through automated quality checks and manual review by one of the coauthors with medical expertise. The final benchmark contains 3,013 images, with class distribution detailed in Figure 6 in Appendix A, plus private training and validation sets with 2,424 and 606 images, respectively.

Enriching Images with Textual Data. To enrich the Patient Agent’s input and support more natural telemedicine dialogues, we generated concise, image-associated descriptions for all 3,013 cases. First, using GPT-4o-mini, selected for its high medical accuracy and relatively low cost (Li et al., 2023a; Smolyak et al., 2024; Abrar et al., 2025), we generated one basic symptom from a human perspective for each of the 34 diagnoses. Then, we expanded each corresponding image for all cases, generating additional structured complaints describing affected areas, duration, intensity, and relevant patient history. This enriched textual input, with generation prompt in Appendix F.1 and exam-

¹<https://github.com/Priyanshu9898/Oral-Disease-Classification>

²<https://www.kaggle.com/datasets/anindamohanta/different-phases-of-tonsillitis>

³<https://www.kaggle.com/datasets/nikhilgurav21/nail-disease-detection-dataset>

⁴<https://www.kaggle.com/datasets/alisofoya/conjunctivitis>

⁵<https://www.microsoft.com/en-us/bing/apis/bing-image-search-api>

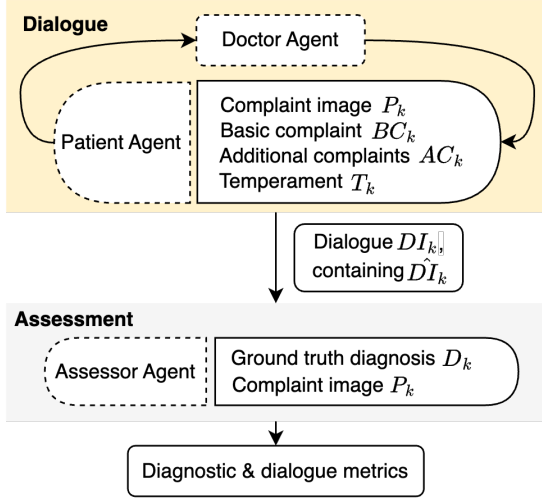


Figure 2: Agents interaction pipeline in 3MDBench.

ples in Appendix G, enhances the visual data and provides context for more informative interactions with the Doctor Agent.

3.2 Task Definition

3MDBench simulates realistic medical consultations via multi-turn dialogues between an evaluated **Doctor Agent** and a **Patient Agent** with further evaluation of an **Assessor Agent**, as show in Figure 2. Each scenario includes a ground truth diagnosis D_k , a complaint image P_k , and symptoms split into a **basic complaint** BC_k (shared initially) and **additional complaints** AC_k (revealed during dialogue). The Patient Agent, shaped by a predefined temperament T_k , starts with access to BC_k , P_k , and AC_k , and interacts using prompts from Appendix F.4 aiming to obtain a diagnosis and medical recommendations from the Doctor Agent.

The Doctor Agent receives an initial Patient Agent query containing only BC_k and P_k , aims to uncover AC_k through dialogue with prompts from Appendix F.5, and outputs diagnosis, treatment plan, and recommendations. Dialogues are capped at 28 utterances, matching average length from our real dialogues (Subsection 3.1), and are marked incomplete if unresolved within this limit.

The Assessor Agent evaluates complete dialogue DI_k by comparing extracted diagnosis \hat{D}_k with D_k , and assessing diagnostic reasoning, communication, and clinical accuracy, based on the prompt in Appendix F.4.

3.3 Patient Agent

The quality of the benchmark depends on the performance of the Patient Agent. This Agent, lacking

access to the ground-truth diagnosis but aware of its symptoms, engages in text-based dialogue and concludes once the doctor provides a diagnosis, recommendations, and answers all questions.

To ensure that candidate models cover different families and provide various strategies, we selected for our Llama-3-8B-Instruct and Llama-3.1-8B from the Llama family (Grattafiori et al., 2024), Qwen2.5-7B and Qwen2.5-14B from the Qwen family (Yang et al., 2024), Falcon-7B (Almazrouei et al., 2023), and GPT-4o-mini (OpenAI, 2024).

These agents must strictly follow system prompts provided in Appendix F.2, respond relevantly to Doctor Agent queries, and remain truthful without hallucination. We used three criteria to assess the Agent’s suitability to these requirements. First, **Instruction following** is a metric on a scale from 1 to 5 indicating whether the patient accurately follows all instructions in the prompt and is evaluated by GPT-4o-mini as this model family has demonstrated performance comparable to humans in medical data evaluation (Li et al., 2024b). Second, **Relevance**, a binary metric for each doctor-patient utterance pair checking the relevance of the patient’s response to the doctor’s utterances, is also evaluated by GPT-4o-mini. The average of all consecutive pairs is calculated as the final metric for the dialogue. Finally, **Factuality** metric is calculated as the ratio of patient utterances referencing symptoms from the prompt. For each patient utterance, we compared its embedding to the embeddings of the symptom using state-of-the-art text embeddings by NV-Embed-v2 (Lee et al., 2024). If any symptom has an embedding cosine similarity higher than the threshold obtained from the practically used threshold of 0.8 (Li et al., 2024c), the utterance is considered to reference the symptoms.

3.4 Assessor Agent

The Assessor Agent is responsible for **evaluating doctor agents** in generated dialogues and **extracting final diagnoses** from the doctor’s conclusions. To assess the clinical competence, we adapted our evaluation criteria from the Mini-Clinical Evaluation Exercise (Mini-CEX) (Shi et al., 2023), a standard in medical education, where patients evaluate medical consultations via structured questions. We simplified its 24 criteria by removing irrelevant items (e.g., autonomy, bias) and merging redundant ones, resulting in 8 core criteria (Table 2).

To select the best model, we measured alignment with human annotations on a diverse validation

Table 2: Criteria for doctor model assessment

Primary Item	Secondary Item
Medical Interviewing Skills	1.1: Does the doctor enquire about a patient’s medical history, such as previous diseases, medications, and surgeries? 1.2: Does the doctor enquire about the current symptoms, possible causes, and attempted treatments? 1.3: Does the doctor explain the basis of the provided conclusion to the patient?
Humanistic Care	2.1: Does the doctor communicate with respect, empathy, and politeness, providing appropriate guidance and avoiding unnecessary extensions? 2.2: Does the doctor respect the individual wishes of the patient?
Diagnostic and Treatment Abilities	3.1: Does the doctor provide an accurate diagnostic plan for the supposed diagnosis? 3.2: Does the doctor accurately provide a treatment plan for the supposed diagnosis?
Overall Clinical Competence	4.1: Which level of clinical competence does the doctor demonstrate during the consultation? (Unsatisfactory, satisfactory, or excellent).

subset of 3MDBench, balanced across four patient temperaments, multiple doctor models, and 34 diagnoses. Five human annotators rated dialogues and extracted diagnoses using the exact instructions as the Assessor Agent (Appendix F.4). Inter-annotator agreement, measured using Cohen’s Kappa and described in Appendix B, yielded an average score of 0.48, considered appropriate given the complexity and subjectivity of clinical assessment, which naturally leads to variability in human judgments (Haas et al., 1996; Verma et al., 2016; Flach et al., 2021).

Next, we collect annotations and extract diagnoses from LVLM-based assessor-candidates using two prompts from Appendix F.4. The candidate assessor models are Qwen2-VL-72B-Instruct (Bai et al., 2023), GPT-4o-mini (OpenAI, 2024), Llava-OneVision-Qwen2-72b-ov-chat-hf (Li et al., 2024a), and DeepSeek-VL (Lu et al., 2024a). Each model received the same input as human annotators in clinical competence evaluation: the dialogue, image, and ground truth diagnosis.

3.5 Evaluated Doctor Agents

The primary goal of 3MDBench is to evaluate the diagnostic capabilities of LVLMs in a simulated telemedicine setting. Specifically, the benchmark assesses a model’s ability to integrate visual and textual modalities to emulate the role of a doctor during a consultation. At the start of each appointment, the doctor model receives a supporting medical image and is expected to engage in an information-seeking dialogue with the patient. The model aims to arrive at an accurate diagnosis informed by the image and the dialogue.

The scarcity of high-quality, domain-specific data remains a significant bottleneck in developing robust LVLMs for medical diagnostics. While recent efforts have introduced specialized models such as Med-Flamingo (Moor et al., 2023) and Biomed-LLaVA (Cheng et al., 2024), these

systems are either limited in scope or availability. As a result, medical diagnostic tasks are often approached using general-purpose LVLMs, sometimes augmented with retrieval-based techniques (Hewitt et al., 2024; Zelin et al., 2024) to improve accuracy via external knowledge. Hence, our study evaluates general-purpose LVLMs and methods for enhancing their performance without relying on domain-specific data. We assess the following models: Qwen2-VL-7B-Instruct, Llama-3.2-11B-Vision-Instruct, and GPT-4o-mini, though an arbitrary LVLM may be used in our benchmark. Inclusion of open-source Qwen2 and LLaVA offers insight into the baseline capabilities of publicly available systems for diagnostic tasks.

We evaluate six prompting variants to systematically study the contribution of visual and textual modalities to diagnostic accuracy and to analyze the effect of different prompting strategies. The first two dialogue-free options provide the lower-bound (**Image + General Complaint**) and upper-bounds (**Image + All Complaints**) for immediate diagnosis \hat{D}_k from the image P_k , general complaint BC_k , and, in the latter case, additional complaints AC_k . Next, we examine various dialogue options: **Dialogue Only** diagnosis from the dialogue DI_k , without access to image P_k , **Dialogue + Image** with image P_k included during the dialogue, and **Dialogue + Image + Rationale** with rationale generation, in which the Doctor Agent explains each step of reasoning internally (hidden from the patient), promoting logical consistency (Wei et al., 2022). Finally, we examine the possibility (**Dialogue + Image + Rationale + External Cues**) to combine LVLM with a ConvNet fine-tuned on the 3MDBench image training set (see details in Appendix C), in which we add top-3 classes, predicted by the CovNet from image P_k . The prompt templates for each setup are provided in Appendix F.5.

Table 3: Comparison of candidate patient models assessed in the diagnostic conversation using GPT-4o-mini based on the three important aspects. Then, the models are ranked based on each aspect, and the mean rank is calculated.

Model Name	Llama-3-8b	Llama-3.1-8b	Qwen2.5-7B	Qwen2.5-14B	Falcon-7B	GPT-4o-mini
Instruction following	4.72	4.74	4.71	4.59	4.37	4.38
Relevance	0.65	0.59	0.84	0.76	0.90	0.82
Factuality	0.79	0.77	0.67	0.78	0.59	0.98
Mean Rank	3.00	3.67	3.33	3.67	4.33	3.00

Table 4: Comparison of assessor models, Cohen’s Kappa and F1 score

Model Name	DeepSeek-VL	Qwen2-VL-72B-Instruct	Llava-OneVision	GPT-4o-mini
Cohen’s Kappa	0.00	0.36	0.43	0.32
F1 score	55.9	78.0	78.0	76.3

4 Results

To enhance the reproducibility of 3MDBench while maintaining dialogue variability, we configured the Patient Agent with a maximum of 256 new tokens and a temperature of 0.6 and the Doctor Agent with a maximum of 512 completion tokens and a temperature of 0.6 (Gusev, 2025). To ensure stability in assessment, we set the Assessor Agent with a maximum of 512 new tokens and a temperature of 1×10^{-6} . To ensure the statistical testing process, we employed the Wilcoxon signed-rank test with a significance level of $\alpha = 0.01$ to assess the statistical significance of the difference in evaluated metrics. We applied false discovery rate control using the Benjamini–Hochberg procedure to account for multiple comparisons (Benjamini and Hochberg, 1995; Hochberg and Tamhane, 2009).

4.1 Patient and Assessment Agents

We evaluated the first two metrics from Section 3.3 using GPT-4o-mini. Table 3 presents the metrics on the validation set of 3MDBench. To make the final selection, we calculated the mean rank for each model across each metric and then averaged them. As a result, we chose Llama-3-8B as our patient model to ensure the benchmark remains open-access and independent of proprietary models. Moreover, by this selection, we implemented one of the proposed hypotheses for paraphrasing text to inhibit self-recognition, thereby mitigating the risk of employing the same model (GPT-4o-mini) for both Doctor Agent and symptom generation (Panickssery et al., 2024).

To estimate the dialogue closeness, we calculated Cohen’s Kappa for each criterion (Table 2) and averaged the scores to determine overall agree-

ment. For the diagnosis extraction task, we computed the F1 score of the diagnoses identified by the LVLm, using human-extracted diagnoses as the ground truth. Table 4 presents the evaluation results. Based on Cohen’s Kappa and F1 score, we selected Llava-OneVision-Qwen2-72b-ov-chat-hf as the final Assessor Agent model, which approached the Cohen’s Kappa score of human annotation mentioned in Section 3.4.

4.2 Diagnostic Results: Doctor Agent

One of the core assumptions of our benchmark is that providing qualitative dialogue and access to medical imaging can enhance diagnostic performance, as shown in Table 5. Here, first, the presence of dialogue improves diagnostic F1-score of GPT-4o-mini from 50.4 (image and basic complaint BC_k) to 54.2% ($p < 0.01$), surpassing open-source models, yet still trails the score 66.8% of unreal setting with full-information, where the doctor is assumed to know all patient details (Li et al., 2024c) already. However, this gap indicates that current LVLms struggle to gather complete information via dialogue, often missing key symptoms or ending prematurely. Unlike human clinicians who adaptively probe to fill gaps, these models can miss symptoms or stop questioning prematurely, limiting dialogue effectiveness. This observation highlights the need for enhanced dialogue strategies to better approximate the completeness and accuracy of full-information diagnostic settings.

Second, the results demonstrate that dialogue quality depends on access to image inputs. Here, with improving F1-score from 52.8 to 54.2% ($p < 0.01$), the average number of utterances per dialogue decreased from 15.22 (± 3.6) without image access to 13.32 (± 3.3) with image access (p

Table 5: Main results of our benchmark: diagnostic F1 scores of doctor agent

Model Name	Configuration	F1 Score	Number of utterances
EfficientNetV2-XL	Fine-tuned on the train part	61.0	-
GPT 4o-mini	No dialogue, image + general complaint	50.4	-
	No dialogue, image + all complaints	66.8	-
	Dialogue, no image	52.8	15.22 (± 3.63)
	Dialogue + image	54.2	13.32 (± 3.33)
	Dialogue + image + rationale	56.9	14.99 (± 4.23)
	Dialogue + image + rationale + external cues	70.3	14.48 (± 3.97)
Llama-3.2-Vision	Dialogue + image	41.5	14.49 (± 4.02)
Qwen2-VL	Dialogue + image	39.0	15.11 (± 4.39)

Table 6: Clinical competence of dialogue doctor systems. See details for criteria in Table 2

Model	1.1	1.2	1.3	2.1	2.2	3.1	3.2	4.1
GPT, dialogue, no image	1.0	1.0	0.95	1.0	1.0	0.89	0.90	1.45
GPT, dialogue + image	0.99	1.0	0.96	1.0	1.0	0.90	0.91	1.61
GPT, dialogue + image + rationale	0.96	0.99	0.89	0.99	0.97	0.78	0.78	1.31
GPT, dialogue + image + rationale + external cues	0.96	0.99	0.94	0.99	0.98	0.88	0.88	1.47
Llama-3.2-Vision	0.99	0.99	0.96	0.99	0.99	0.75	0.74	1.45
Qwen2-VL	0.90	0.93	0.78	0.92	0.90	0.61	0.61	1.16

< 0.01). Hence, the inclusion of visual information not only improves diagnostic accuracy but also leads to shorter, more efficient interactions.

We tested prompting strategies that avoid direct fine-tuning to demonstrate the effect of various strategies for the Doctor Agent. Building on prior work suggesting the benefits of chain-of-thought prompting (Wei et al., 2022), rationale generation shows significant F1 improvement over standard dialogue (56.9% vs. 54.2%, $p < 0.01$), indicating that explanations alone can enhance diagnostic reasoning in complex tasks. Moreover, enriching input with image-based cues, specifically the top-3 predictions from a fine-tuned EfficientNetV2-XL (Tan and Le, 2021) with details in Appendix C, boosts F1 score to 70.3%, outperforming the full-information setting and EfficientNetV2-XL-only ($p < 0.01$). Thus, integrating a domain-specific vision model with a general-purpose LVLM may significantly improve the diagnostic ability.

We also evaluated diagnostic accuracy across five disease categories (see Figure 3 and Figure 7 in Appendix A). Performance varies considerably by category. Dermatology, with many overlapping conditions, yields the lowest average F1 (46.9%), while throat/mucosae, with more apparent distinctions, scores highest (88.1%). This result reflects model limitations in fine-grained classification and the dataset’s uneven diagnostic coverage.

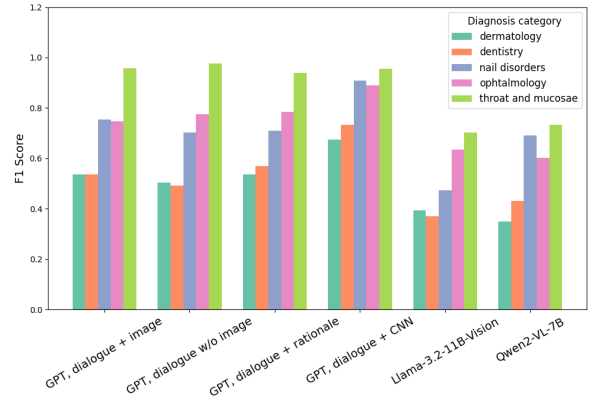


Figure 3: F1 scores by diagnosis categories

4.3 Benchmarking Clinical Competence

Beyond diagnostic accuracy, we evaluated general clinical competence during diagnostic conversations (see criteria from Table 2). As shown in Table 6, all models scored highly, with no score below 0.87, indicating strong baseline capabilities in empathy, professionalism, and rational behavior. Notably, models without rationale generation outperformed rationale-based agents on most humanistic and communication metrics. The latter scored lower on *Comprehensive Diagnostic and Treatment Abilities* (0.78 vs. 0.90), likely due to focusing narrowly on internal reasoning at the expense of patient-centered communication. Thus,

despite limitations in diagnostic accuracy, general-purpose LVLMs excel at demonstrating empathy and providing emotional support and advice.

4.4 Patient Temperament

Our experiments with personality types (Figures 4, 5, see also detailed results in Figure 11) show no statistically significant differences in F1 score between personalities. Indeed, LVLMs can maintain coherent, goal-directed dialogue even when faced with challenging behaviors, such as the sanguine patient’s digressions or the melancholic patient’s tendency to ask rather than answer questions.

However, dialogues with phlegmatic patients yield slightly lower F1 and competence scores. Characterized by short, passive responses, this personality type limits the model’s ability to collect rich clinical information. These patients tend to answer directly but rarely offer unsolicited detail, requiring the doctor to take greater initiative by asking follow-up questions. For instance, as shown in samples of simulated dialogues for eczema diagnosis shown in Appendix H, when asked about exposure to new products, a phlegmatic patient might respond, "No," without elaboration. In contrast, more expressive personalities, like sanguine, tend to volunteer additional context: "I’ve been using the same old stuff for years ... It’s like, I’ll be going about my day, and suddenly I’ll feel this intense itchiness ... if I’m wearing shoes or socks, it gets even worse." These spontaneous details help guide the diagnostic process more effectively.

Moreover, as shown in Figure 5, dialogues with phlegmatic patients are, on average, four turns shorter due to their tendency not to ask clarifying questions. Appendix H reveals that the doctor agent rarely compensates for this brevity by steering the conversation or probing deeper. While diagnostic accuracy is generally maintained, these interactions result in fewer recommendations or explanations.

Thus, although LVLMs are robust to diverse user behaviors, their performance may still degrade with minimally cooperative patients. This observation underscores the need to assess models’ initiative and adaptability in less cooperative settings.

5 Conclusion

This paper introduces 3MDBench (Fig. 2), an open-source benchmark for evaluating LVLMs in medical diagnostics. It simulates interactive telemedicine consultations, incorporating diverse

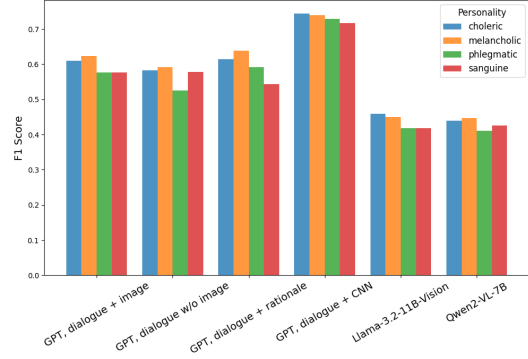


Figure 4: F1 scores by personality types

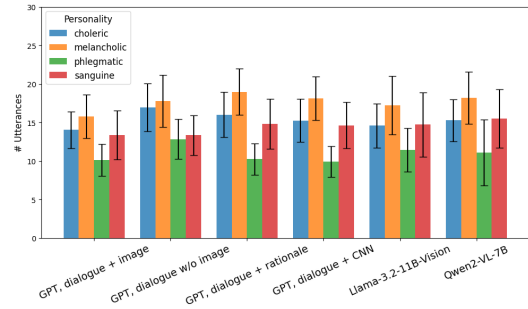


Figure 5: Number of utterances by personality types

diagnoses and patient behaviors to assess diagnostic accuracy and clinical competence.

Within this framework, we demonstrate that the visual modality and the ability to engage in dialogue with the patient significantly enhance diagnostic accuracy (Table 5). The general-purpose LVLMs display strong clinical competence, effectively leveraging images and conducting information-seeking dialogues to provide accurate diagnoses.

We highlight a key limitation of LVLMs in medical diagnostics: while strong in language, they lack domain-specific visual expertise. We demonstrate how to improve the quality of the Doctor agent using our benchmark by incorporating predictions from a ConvNet trained on the diagnosis prediction task, significantly (up to 20%) enhancing LVLM performance. Thus, combining general-purpose LVLMs with lightweight, task-specific vision models offers a scalable opportunity to higher performance without extensive supervised fine-tuning.

Our findings suggest that while dialogue contributes to more accurate diagnosis, its effectiveness is currently limited. External expert cues and better prompting can bridge the gap, while broader and more balanced diagnostic coverage remains a key goal for future benchmarks.

Limitations

While 3MDBench is already a valuable and practical benchmark for evaluating medical multimodal dialogue systems, there are areas for future improvement. First, we generated additional patient symptoms using GPT-4o-mini conditioned on image and diagnosis. Although this approach leverages embedded medical knowledge, it may introduce factual inaccuracies or biases. Similarly, LLM-based evaluation under the LLM-as-a-judge paradigm (Zheng et al., 2023) depends on the evaluator’s domain competence and may propagate errors. We conducted manual checks for plausibility and coherence for the described issues, but we cannot guarantee full correctness.

Second, the benchmark also uses publicly available images, which introduces a potential risk of data leakage from pretraining. Additionally, the current 34 diagnoses, while chosen to reflect real-world consultation distributions and public data availability, limit diagnostic coverage. Future versions should expand the disease set and conceal candidate labels for a more open-ended and robust evaluation.

Finally, while the four temperament categories provide a helpful foundation for simulating patient diversity, future work could explore more nuanced or data-driven models of patient behavior to reflect the variability observed in real clinical settings complexly.

Ethics Statement

Human Involvement This work involved several instances of human annotation. First, one of the co-authors with a medical background reviewed the collected images over one week to verify the correctness of the associated diagnoses. Second, we obtained human annotations to evaluate dialogues for selecting the Assessor Agent. Five employees completed the annotation process, each dedicated approximately six hours to the task during their regular working hours, without additional compensation. All annotators were informed of the research purpose behind the annotation tasks.

Inference Costs Running the complete evaluation experiment on a single A100 GPU took approximately 48 hours for selecting the Patient Agent model, 4 hours for selecting candidate Assessor Agent models, and 210 hours to evaluate the Doctor Agents.

Use of AI Assistants We used Grammarly to improve and proofread the text of this paper, including grammar, spelling, style corrections, and sentence rephrasing. As a result, some parts of the manuscript may be classified as AI-generated, AI-edited, or a mix of human and AI contributions.

References

- Moaiz Abrar, Yusuf Sermet, and Ibrahim Demir. 2025. [An empirical evaluation of large language models on consumer health questions](#). *BioMedInformatics*, 5:12.
- Ehizogie Paul Adeghe, Chioma Anthonia Okolo, and Olumuyiwa Tolulope Ojeyinka. 2024. A review of emerging trends in telemedicine: Healthcare delivery transformations. *Int. J. Life Sci. Res. Arch.*, 6(1):137–147.
- Reem Agbareia, Mahmud Omar, Shelly Soffer, Benjamin S. Glicksberg, Girish N. Nadkarni, and Eyal Klang. 2025. [Visual-textual integration in llms for medical diagnosis: A preliminary quantitative analysis](#). *Computational and Structural Biotechnology Journal*, 27:184–189.
- Harrison AI. 2024. [Radbench: A radiological benchmark for evaluating vision-language models](#).
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *Preprint*, arXiv:2311.16867.
- Dorothee Amelung, Katriina L Whitaker, Debby Lennard, Margaret Ogden, Jessica Sheringham, Yin Zhou, Fiona M Walter, Hardeep Singh, Charles Vincent, and Georgia Black. 2020. Influence of doctor-patient conversations on behaviours of patients presenting to primary care with new or persistent symptoms: a video observation study. *BMJ Qual. Saf.*, 29(3):198–208.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Yoav Benjamini and Yosef Hochberg. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.
- Sonu Bhaskar, Sian Bradley, Vijay Kumar Chattu, Anil Adishes, Alma Nurtazina, Saltanat Kyrkybayeva, Sateesh Sakhamuri, Sebastian Moguilner, Shawna

768	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve	824
769	Hanyi Fang, and Peter Szolovits. 2020. What dis-	Jiang, and You Zhang. 2023b. Chatdoctor: A medical	825
770	ease does this patient have? a large-scale open do-	chat model fine-tuned on a large language model	826
771	main question answering dataset from medical exams.	meta-ai (llama) using medical domain knowledge.	827
772	<i>Preprint</i> , arXiv:2009.13081.	<i>Preprint</i> , arXiv:2303.14070.	828
773	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve	829
774	Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset	Jiang, and You Zhang. 2023c. Chatdoctor: A medical	830
775	for biomedical research question answering. In <i>Pro-</i>	chat model fine-tuned on a large language model	831
776	<i>ceedings of the 2019 Conference on Empirical Meth-</i>	meta-ai (llama) using medical domain knowledge.	832
777	<i>ods in Natural Language Processing and the 9th In-</i>	<i>Cureus</i> , 15.	833
778	<i>ternational Joint Conference on Natural Language</i>		
779	<i>Processing (EMNLP-IJCNLP)</i> , pages 2567–2577.		
780	Shreya Johri, Jaehwan Jeong, Benjamin A. Tran,	Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and	834
781	Daniel I. Schlessinger, Shannon Wongvibulsin,	Xiao-Ming Wu. 2021. Slake: A semantically-labeled	835
782	Zhuo Ran Cai, Roxana Daneshjou, and Pranav Ra-	knowledge-enhanced dataset for medical visual ques-	836
783	ajpurkar. 2024. Guidelines for rigorous evaluation of	tion answering. pages 1650–1654.	837
784	clinical llms for conversational reasoning.		
785	Yubin Kim, Chanwoo Park, Hyewon Jeong, Yik Siu	Wenge Liu, Jianheng Tang, Yi Cheng, Wenjie Li,	838
786	Chan, Xuhai Xu, Daniel McDuff, Hyeonhoon Lee,	Yefeng Zheng, and Xiaodan Liang. 2022. MedDG:	839
787	Marzyeh Ghassemi, Cynthia Breazeal, and Hae Won	An Entity-Centric Medical Consultation Dataset for	840
788	Park. 2024a. Mdagents: An adaptive collabora-	Entity-Aware Medical Dialogue Generation , pages	841
789	tion of llms for medical decision-making. <i>Preprint</i> ,	447–459.	842
790	arXiv:2404.15155.		
791	Yunsoo Kim, Jing Wu, Yusuf Abdulle, and Honghan	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai	843
792	Wu. 2024b. Medexqa: Medical question answer-	Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhu-	844
793	ing benchmark with multiple explanations. <i>Preprint</i> ,	oshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng,	845
794	arXiv:2406.06331.	Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024a.	846
		Deepseek-vl: Towards real-world vision-language	847
		understanding. <i>Preprint</i> , arXiv:2403.05525.	848
795	Jason Lau, Soumya Gayen, Asma Ben Abacha, and	Zhiyong Lu, Yifan Peng, Trevor Cohen, Marzyeh Ghas-	849
796	Dina Demner-Fushman. 2018. A dataset of clini-	semi, Chunhua Weng, and Shubo Tian. 2024b. Large	850
797	cally generated visual questions and answers about	language models in biomedicine and health: current	851
798	radiology images. <i>Scientific Data</i> , 5:180251.	research landscape and future directions. <i>J. Am. Med.</i>	852
		<i>Inform. Assoc.</i> , 31(9):1801–1811.	853
799	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan	Itay Manes, Naama Ronn, David Cohen, Ran Ilan	854
800	Raiman, Mohammad Shoeybi, Bryan Catanzaro, and	Ber, Zehavi Horowitz-Kugler, and Gabriel Stanovsky.	855
801	Wei Ping. 2024. Nv-embed: Improved techniques for	2024. K-qa: A real-world medical q&a benchmark.	856
802	training llms as generalist embedding models. <i>arXiv</i>	<i>Preprint</i> , arXiv:2401.14493.	857
803	<i>preprint arXiv:2405.17428.</i>		
804	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng	Carlotta J. Mayer, Julia Mahal, Daniela Geisel, Eva J.	858
805	Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei	Geiger, Elias Staats, Maximilian Zappel, Seraina P.	859
806	Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy	Lerch, Johannes C. Ehrenthal, Steffen Walter, and	860
807	visual task transfer. <i>Preprint</i> , arXiv:2408.03326.	Beate Ditzen. 2024. User preferences and trust in	861
		hypothetical analog, digitalized and ai-based medical	862
808	Kevin Danis Li, Adrian M Fernandez, Rachel Schwartz,	consultation scenarios: An online discrete choice	863
809	Natalie Rios, Marvin Nathaniel Carlisle, Gregory M	survey. <i>Computers in Human Behavior</i> , 161:108419.	864
810	Amend, Hiren V Patel, and Benjamin N Breyer.		
811	2024b. Comparing gpt-4 and human researchers	Nikita Mehndru, Brenda Y. Miao, Eduardo Rodriguez	865
812	in health care data analysis: Qualitative description	Almaraz, Madhumita Sushil, Atul J. Butte, and	866
813	study. <i>J Med Internet Res</i> , 26:e56500.	Ahmed Alaa. 2024. Evaluating large language mod-	867
		els as agents in the clinic. <i>npj Digital Medicine</i> ,	868
814	Rumeng Li, Xun Wang, and Hong Yu. 2023a. Two	7(1):84.	869
815	directions for clinical data generation with large		
816	language models: Data-to-label and label-to-data.	Ashley Meyer, Velma Payne, Derek Meeks, Radha Rao,	870
817	<i>Proc. Conf. Empir. Methods Nat. Lang. Process.</i> ,	and Hardeep Singh. 2013. Physicians’ diagnostic ac-	871
818	2023:7129–7143.	curacy, confidence, and resource requests a vignette	872
		study. <i>JAMA internal medicine</i> , 173.	873
819	Shuyue Stella Li, Vidhisha Balachandran, Shangbin	Tala Mirzaei and Nicole Kashian. 2020. Revisiting	874
820	Feng, Jonathan S. Ilgen, Emma Pierson, Pang Wei	effective communication between patients and physi-	875
821	Koh, and Yulia Tsvetkov. 2024c. Mediq: Question-	cians: Cross-sectional questionnaire study comparing	876
822	asking llms and a benchmark for reliable interactive	text-based electronic versus face-to-face communica-	877
823	clinical reasoning. <i>Preprint</i> , arXiv:2406.00922.	tion. <i>J Med Internet Res</i> , 22(5):e16965.	878

879	Michael Moor, Qian Huang, Shirley Wu, Michihiro	Xiaoming Shi, Zeming Liu, Li Du, Yuxuan Wang, Hon-	934
880	Yasunaga, Cyril Zakka, Yashodhara Dalmia, Ed-	gru Wang, Yuhang Guo, Tong Ruan, Jie Xu, Xiaofan	935
881	uardo Pontes Reis, Pranav Rajpurkar, and Jure	Zhang, and Shaoting Zhang. 2024. Medical dialogue	936
882	Leskovec. 2023. Med-flamingo: a multimodal medi-	system: A survey of categories, methods, evaluation	937
883	cal few-shot learner . <i>ArXiv</i> , abs/2307.15189.	and challenges . In <i>Findings of the Association for</i>	938
		<i>Computational Linguistics: ACL 2024</i> , pages 2840–	939
884	Amy D. Nguyen, Sarah J. White, Tim Tse, John A.	2861, Bangkok, Thailand. Association for Computa-	940
885	Cartmill, Peter Roger, Sarah Hatem, and Simon M.	tional Linguistics.	941
886	Willcock. 2024. Communication during telemedicine		
887	consultations in general practice: perspectives from	Xiaoming Shi, Jie Xu, Jinru Ding, Jiali Pang, Sichen	942
888	general practitioners and their patients . <i>BMC Pri-</i>	Liu, Shuqing Luo, Xingwei Peng, Lu Lu, Haihong	943
889	<i>mary Care</i> , 25(1):324.	Yang, Mingtao Hu, Tong Ruan, and Shaoting Zhang.	944
		2023. Llm-mini-cex: Automatic evaluation of large	945
890	Ejike Innocent Nwankwo, Ebube Victor Emeihe,	language model for diagnostic conversation .	946
891	Mojeed Dayo Ajegbile, Janet Aderonke Olaboye,		
892	and Chukwudi Cosmos Maha. 2024. Integrating	Hardeep Singh and Dean F Sittig. 2015. Advancing	947
893	telemedicine and AI to improve healthcare access	the science of measurement of diagnostic errors in	948
894	in rural settings. <i>Int. J. Life Sci. Res. Arch.</i> , 7(1):059–	healthcare: the safer dx framework. <i>BMJ Qual. Saf.</i> ,	949
895	077.	24(2):103–110.	950
896	OpenAI. 2024. Gpt-4o system card . <i>Preprint</i> ,	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara	951
897	arXiv:2410.21276 .	Mahdavi, Jason Wei, Hyung Won Chung, Nathan	952
		Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen	953
898	World Health Organization. 2004. Icd-10 : international	Pfohl, Perry Payne, Martin Seneviratne, Paul Gam-	954
899	statistical classification of diseases and related health	ble, Chris Kelly, Abubakr Babiker, Nathanael Schärli,	955
900	problems : tenth revision.	Aakanksha Chowdhery, Philip Mansfield, Dina	956
		Demner-Fushman, and 13 others. 2023. Large lan-	957
901	Davy Paap, Leonie A Krops, Henrica R Schiphorst Pre-	guage models encode clinical knowledge . <i>Nature</i> ,	958
902	uper, Jan H B Geertzen, Pieter U Dijkstra, and Gri-	620(7972):172–180.	959
903	eteke Pool. 2022. Participants’ unspoken thoughts		
904	and feelings negatively influence the therapeutic al-	Daniel Smolyak, Margrét V Bjarnadóttir, Kenyon Crow-	960
905	liance; a qualitative study in a multidisciplinary pain	ley, and Ritu Agarwal. 2024. Large language models	961
906	rehabilitation setting. <i>Disabil. Rehabil.</i> , 44(18):5090–	and synthetic health data: progress and prospects .	962
907	5100.	<i>JAMIA Open</i> , 7(4):ooae114.	963
908	Ankit Pal, Logesh Kumar Umapathi, and Malaikannan	R. Steiner. 1985. <i>The Four Temperaments</i> . Anthropol-	964
909	Sankarasubbu. 2022. Medmcqa: A large-scale multi-	sophic Press.	965
910	subject multi-choice dataset for medical domain ques-		
911	tion answering . In <i>Proceedings of the Conference</i>	Mason Stoltzfus, Arshdeep Kaur, Avantika Chawla,	966
912	<i>on Health, Inference, and Learning</i> , volume 174 of	Vasu Gupta, F. N. U. Anamika, and Rohit Jain.	967
913	<i>Proceedings of Machine Learning Research</i> , pages	2023. The role of telemedicine in healthcare: an	968
914	248–260. PMLR.	overview and update . <i>The Egyptian Journal of Inter-</i>	969
		<i>nal Medicine</i> , 35(1):49.	970
915	Arjun Panickssery, Samuel R. Bowman, and Shi Feng.		
916	2024. Llm evaluators recognize and favor their own	Mingxing Tan and Quoc Le. 2021. Efficientnetv2:	971
917	generations . <i>Preprint</i> , arXiv:2404.13076.	Smaller models and faster training . In <i>Proceedings of</i>	972
		<i>the 38th International Conference on Machine Learn-</i>	973
918	Kinalyne Perez, Daniela Wisniewski, Arzu Ari, Kim	<i>ing</i> , volume 139 of <i>Proceedings of Machine Learning</i>	974
919	Lee, Cristian Lieneck, and Zo Ramamonjiravelo.	<i>Research</i> , pages 10096–10106. PMLR.	975
920	2025. Investigation into application of AI and		
921	telemedicine in rural communities: A systematic lit-	Andru Twinanda, Sherif Shehata, Didier Mutter,	976
922	erature review. <i>Healthcare (Basel)</i> , 13(3).	Jacques Marescaux, Michel De Mathelin, and Nico-	977
		las Padoy. 2016. Endonet: A deep architecture for	978
923	Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya,	recognition tasks on laparoscopic videos .	979
924	Tien Yin Wong, Ara Darzi, Wu Yuan, and Eric J.		
925	Topol. 2024. Llm-based agentic systems in medicine	Nupur Verma, Daniel S Hippe, and Jeffrey D Robinson.	980
926	and healthcare . <i>Nature Machine Intelligence</i> ,	2016. Assessment of interobserver variability in the	981
927	6(12):1418–1420.	peer review process: Should we agree to disagree?	982
		<i>AJR Am J Roentgenol</i> , 207(6):1215–1222.	983
928	Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause,		
929	Sanjeev Satheesh, Sean Ma, Zhiheng Huang, An-	Junda Wang, Zonghai Yao, Zhichao Yang, Huixue Zhou,	984
930	drej Karpathy, Aditya Khosla, Michael Bernstein,	Rumeng Li, Xun Wang, Yucheng Xu, and Hong	985
931	Alexander C. Berg, and Li Fei-Fei. 2015. Imagenet	Yu. 2024. Notechat: A dataset of synthetic patient-	986
932	large scale visual recognition challenge . <i>Preprint</i> ,	physician conversations conditioned on clinical notes .	987
933	arXiv:1409.0575 .	In <i>Findings of the Association for Computational Lin-</i>	988
		<i>guistics ACL 2024</i> , page 15183–15201. Association	989
		for Computational Linguistics.	990

- Abbi Ward, Jimmy Li, Julie Wang, Sriram Lakshminarasimhan, Ashley Carrick, Bilson Campana, Jay Hartford, Pradeep K. Sreenivasaiah, Tiya Tiyasirisokchai, Sunny Virmani, Renee Wong, Yossi Matias, Greg S. Corrado, Dale R. Webster, Margaret Ann Smith, Dawn Siegel, Steven Lin, Justin Ko, Alan Karthikesalingam, and 2 others. 2024. [Creating an empirical dermatology dataset through crowdsourcing with web search advertisements](#). *JAMA Network Open*, 7(11):e2446615–e2446615.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#).
- An Yang, Baosong Yang, and Binyuan Hui et al. 2024. [Qwen2 technical report](#). *Preprint*, arXiv:2407.10671.
- Songhua Yang, Hanjia Zhao, Zhu Senbin, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2023. [Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue](#).
- Aijia Yuan, Edlin Garcia Colato, Bernice Pescosolido, Hyunju Song, and Sagar Samtani. 2025. [Improving workplace well-being in modern organizations: A review of large language model-based mental health chatbots](#). *ACM Trans. Manage. Inf. Syst.*, 16(1).
- Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Yang Zhenzhu, Yibo Liu, Wenhao Huang, and Wenhua Chen. 2024. [Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi](#). pages 9556–9567.
- Charlotte Zelin, Wendy K. Chung, Mederic Jeanne, Gongbo Zhang, and Chunhua Weng. 2024. [Rare disease diagnosis using knowledge guided retrieval augmentation for chatgpt](#). *Journal of Biomedical Informatics*, 157:104702.
- Guangtao Zeng, Wenmian Yang, Zeqian Ju, Yue Yang, Sicheng Wang, Ruisi Zhang, Meng Zhou, Jiaqi Zeng, Xiangyu Dong, Ruoyu Zhang, Hongchao Fang, Penghui Zhu, Shu Chen, and Pengtao Xie. 2020. [Meddialog: Large-scale medical dialogue datasets](#). pages 9241–9250.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haoteng Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *ArXiv*, abs/2306.05685.
- Jiayuan Zhu and Junde Wu. 2025. [Ask patients with patience: Enabling llms for human-centric medical dialogue with grounded reasoning](#). *Preprint*, arXiv:2502.07143.

A 3MDBench Statistics

Figure 6 shows the class distribution in 3MDBench, obtained by merging data from the utilized datasets and enriching them with data from the Bing Image Search API. We adjusted the resulting distribution to approximate real-world diagnosis frequencies observed in telemedicine consultations from Figure 7.

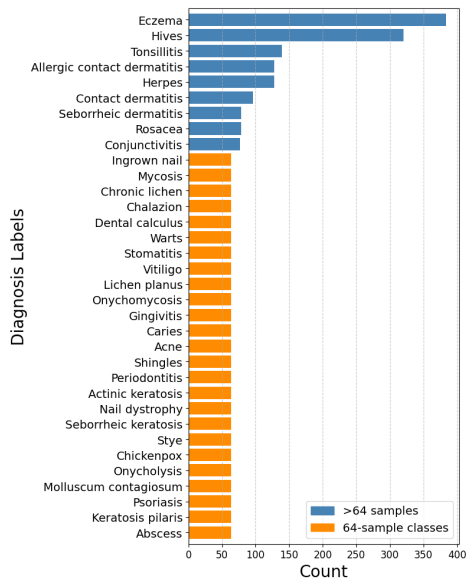


Figure 6: The distribution of classes in 3MDBench. The dataset consists of 34 medical conditions, with the most frequent class containing 384 samples, while 21 classes have exactly 64 samples (highlighted in orange).

Figure 7 presents the distribution of medical diagnoses derived from real-world telemedicine consultations and grouped by medical category.

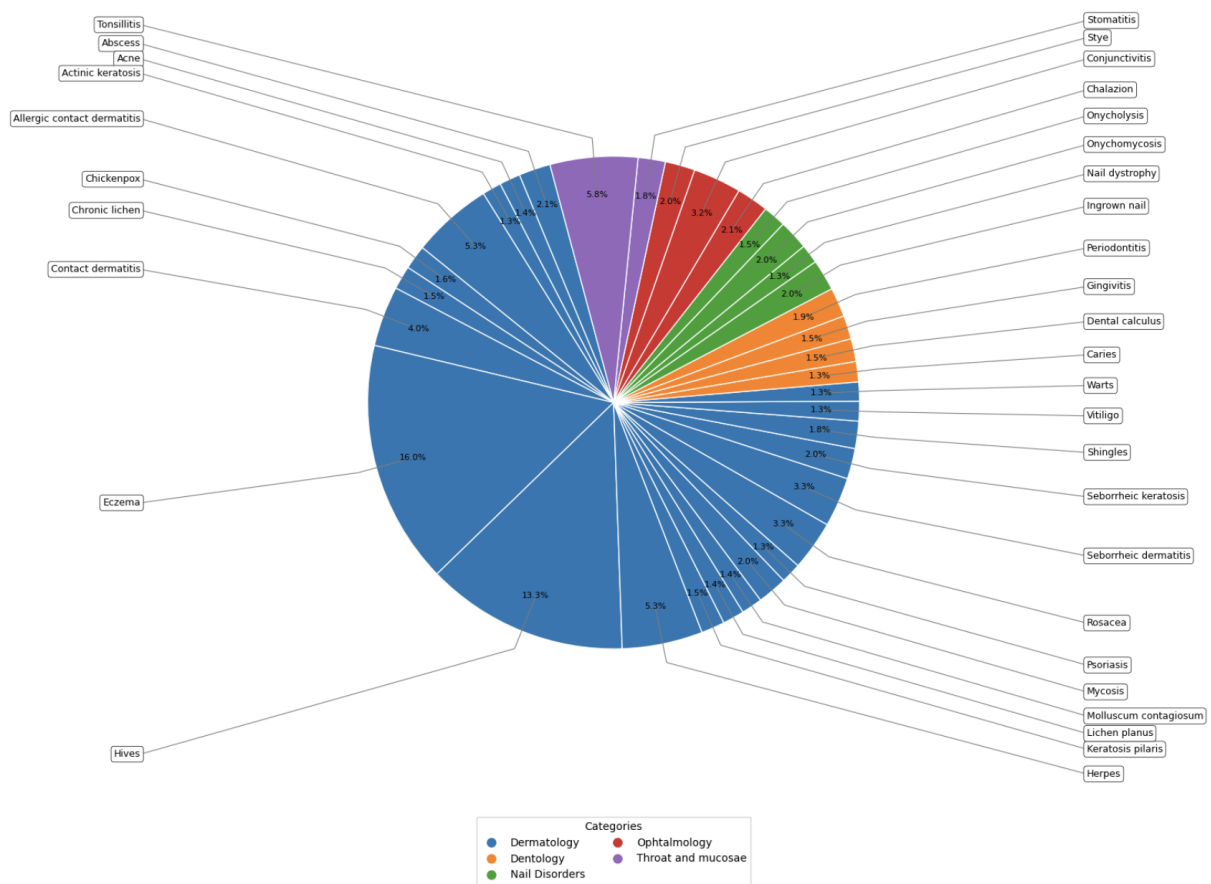


Figure 7: Distribution of selected diagnoses based on real-world telemedicine consultations.

B Human Agreement

We assessed inter-rater agreement by analyzing consistency across all participants. Annotation was conducted using Google Forms so annotators could not see each other’s answers. Figure 8 presents a graph with nodes representing anonymized participants (#1 to #5) and edge weights corresponding to pairwise Cohen’s Kappa scores. However, most participants demonstrated moderate to high consistency. However, participant #1 stood out with notably weaker agreement scores—an average Kappa of 0.42 with others, and only 0.26 when measured against their responses across repeated tasks.

This participant also exhibited abnormally fast task completion times and random-like response patterns, suggesting noncompliance with instructions. After we excluded this data, the overall average Kappa across assessors increased to 0.48, indicating improved inter-rater reliability.

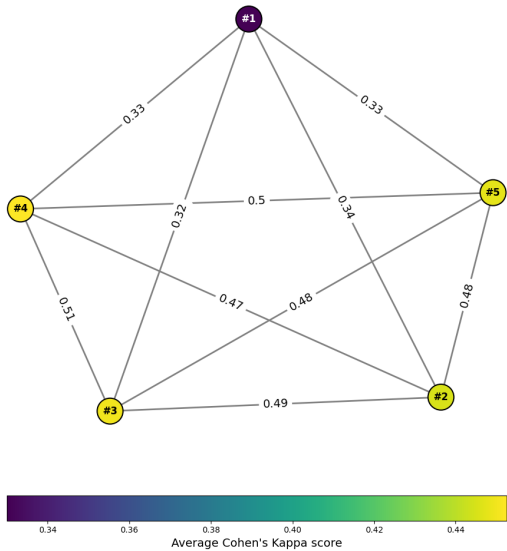


Figure 8: Inter-rater agreement graph showing pairwise Cohen’s Kappa scores between participants (#1–#5). Node color indicates average agreement, and edge weight reflects pairwise consistency.

C Baseline Computer Vision Convolution Model Selection

To enhance LVLMs with disease-related information from images, we trained and compared different CV model of varying sizes to assess their performance in disease classification within the given setting. We selected the EfficientNetV2 model family, pre-trained on ImageNet-1k (Russakovsky et al., 2015) due to its efficient training process

Table 7: Performance comparison of baseline EfficientNetV2 models

Model name	Macro F1 score, %
EfficientNetV2 S	60
EfficientNetV2 M	52
EfficientNetV2 L	57
EfficientNetV2 XL	64

and competitive performance compared to other CNN and ViT models (Tan and Le, 2021), and applied full fine-tuning on the train part of the 3MD-Bench. As shown in Table 7, the EfficientNetV2-XL model achieved the highest classification performance among the tested models. Therefore, we selected this model for further integration with the Doctor Agent.

D Analysis of the Performance on the Benchmark

1094

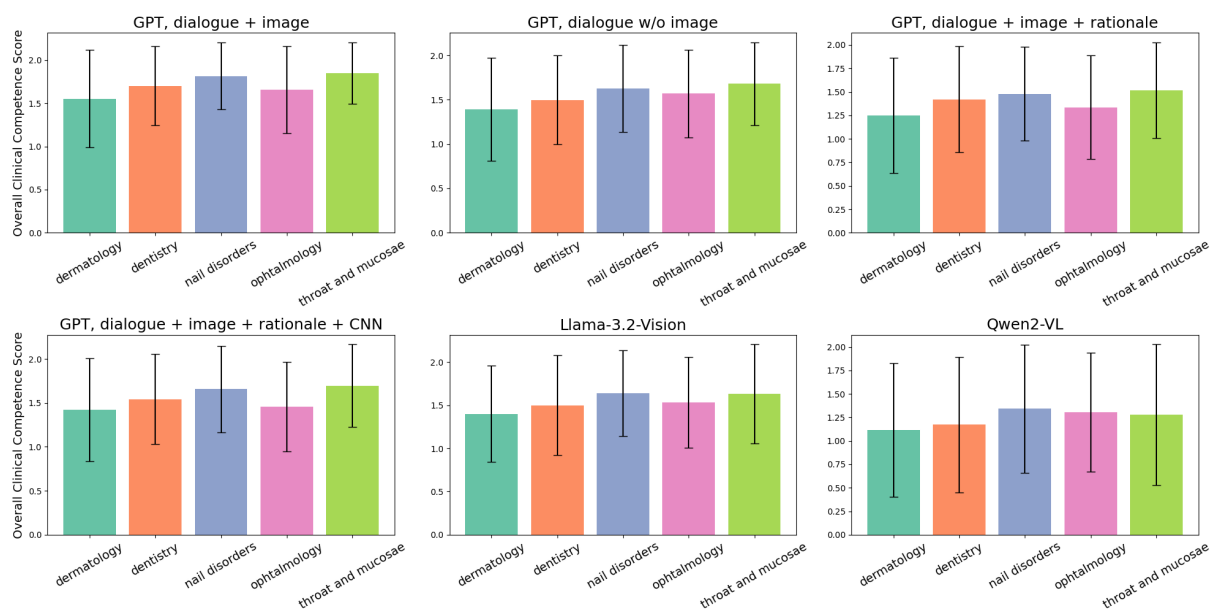


Figure 9: Overall clinical competence scores by diagnosis categories

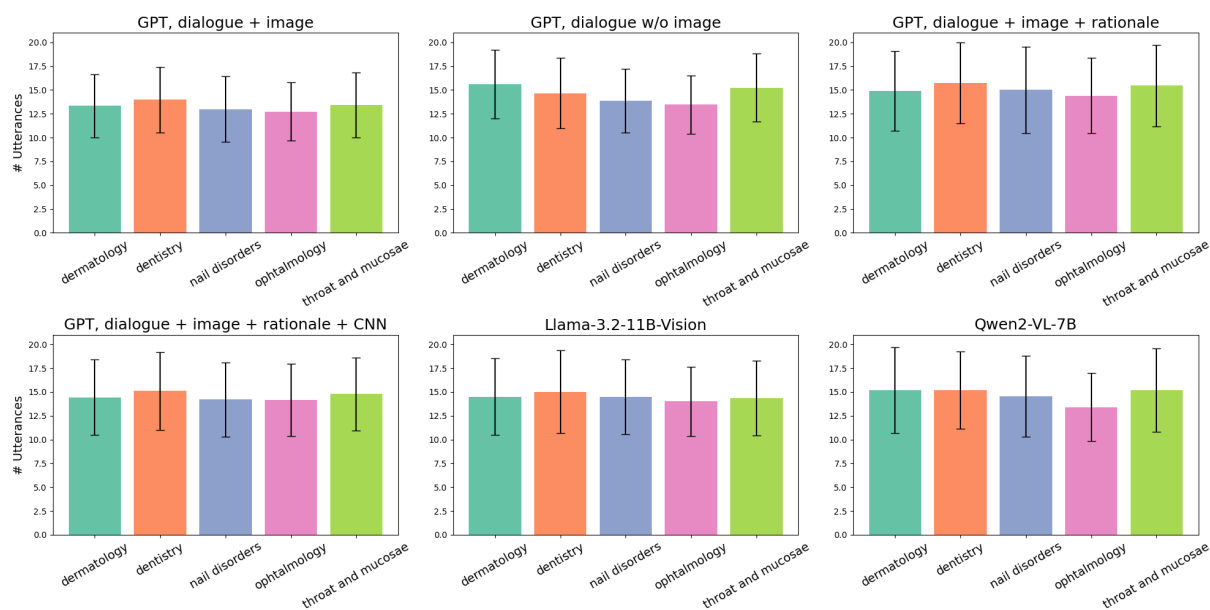


Figure 10: Number of utterances by diagnosis categories

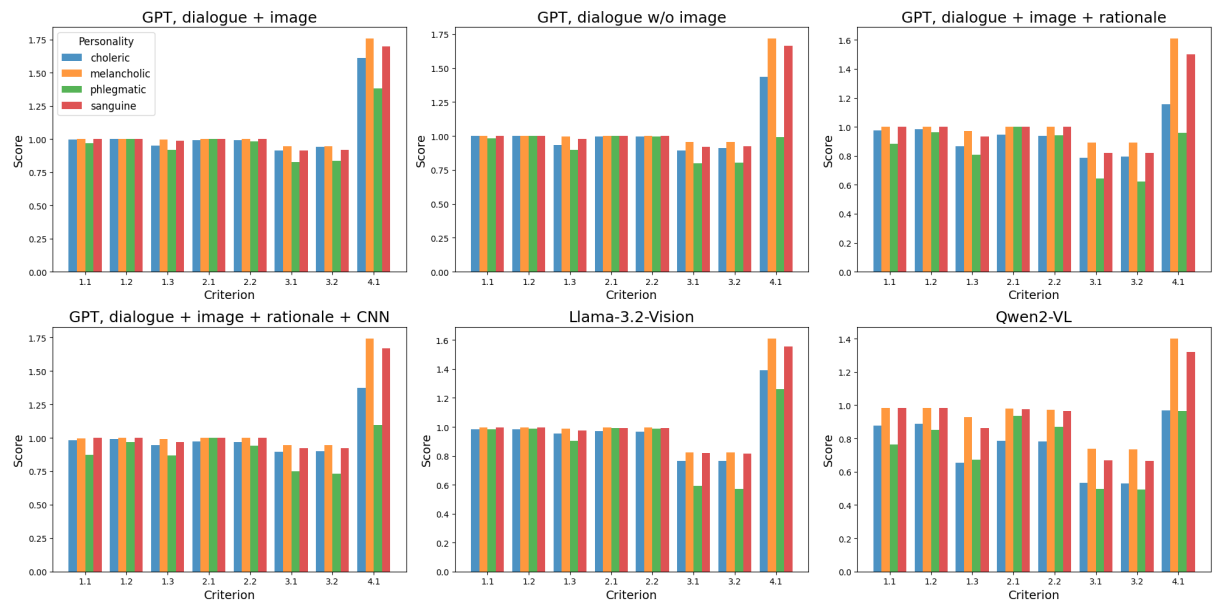


Figure 11: Clinical competence scores by personality types

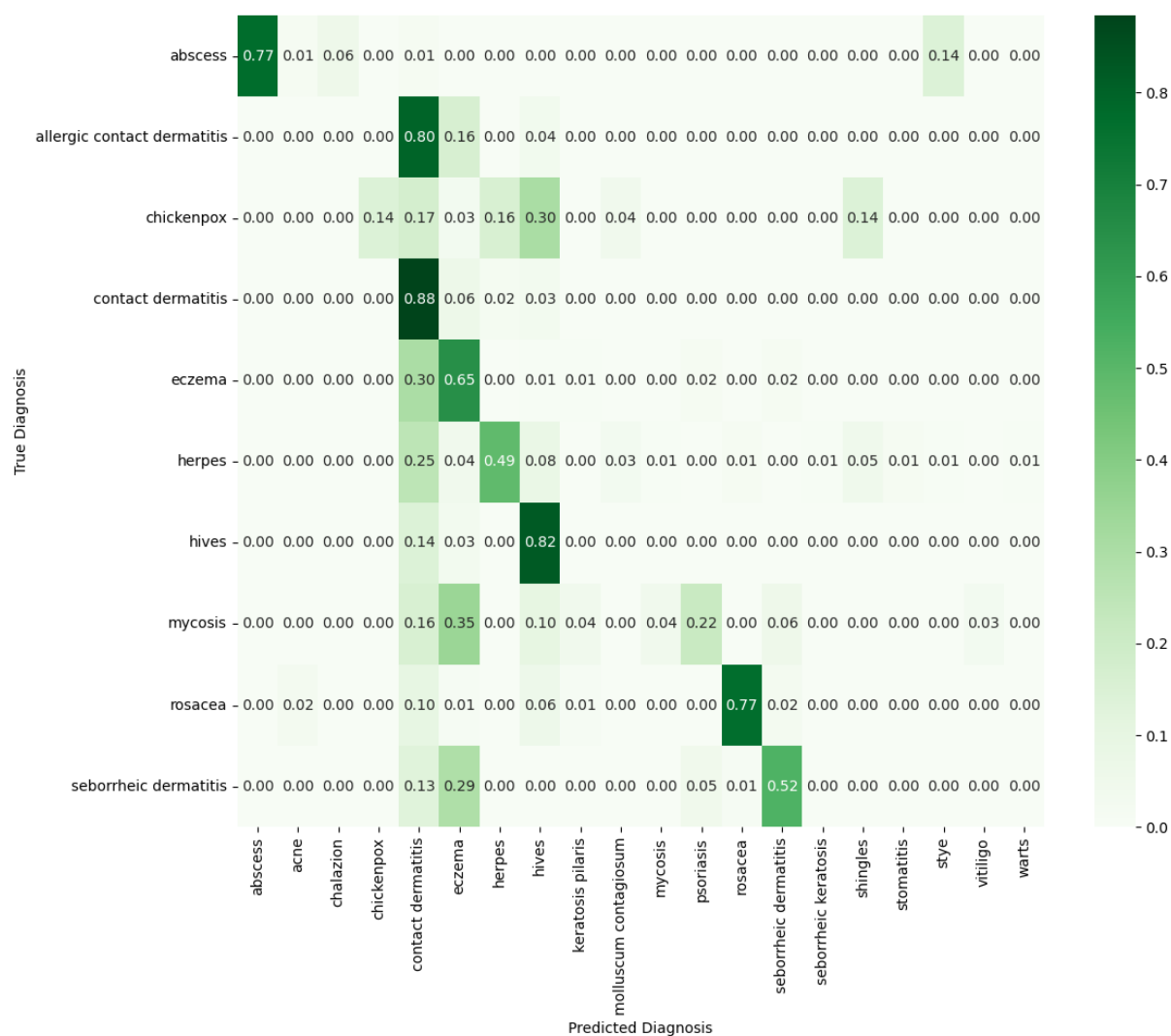
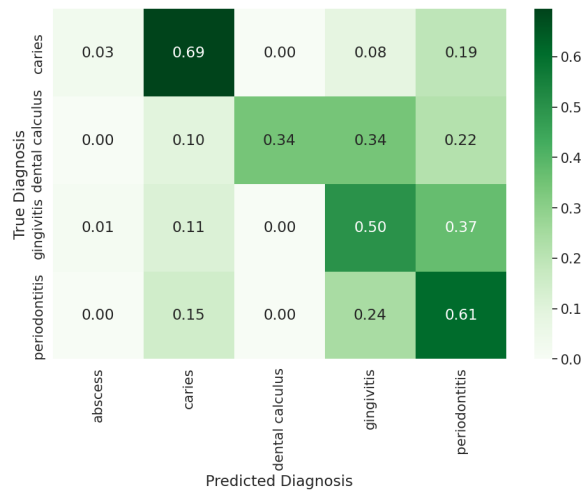
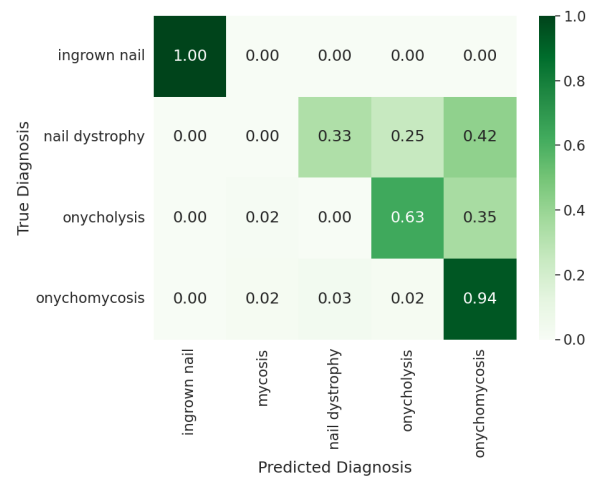


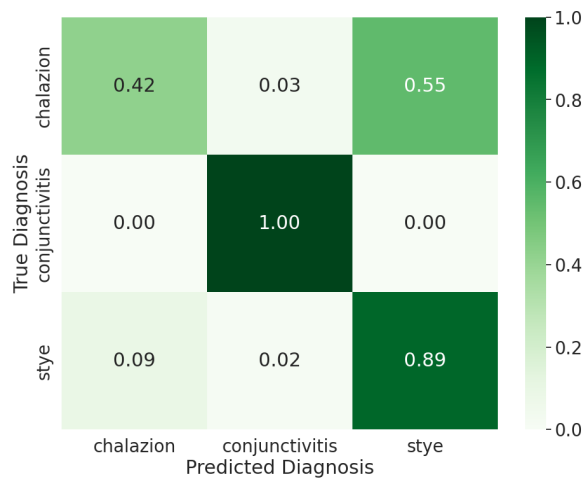
Figure 12: Confusion matrix for the predictions on the samples with the 10 most frequent diagnoses from the dermatology category for GPT-4o-mini with dialogue and rationale generation.



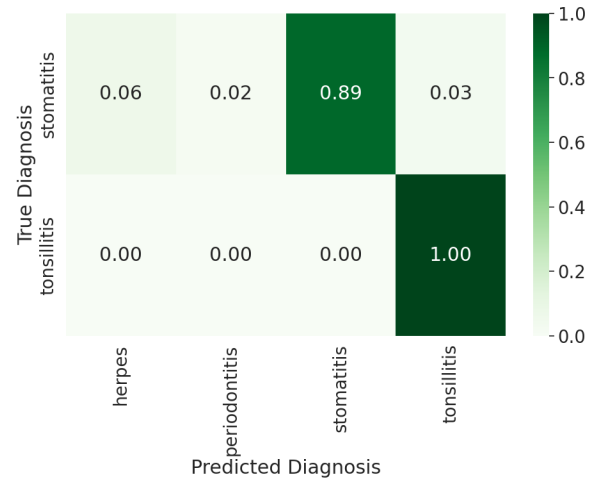
(a) Dentistry category



(b) Nail diseases category



(c) Ophthalmology category



(d) Throat and mucosae category

Figure 13: Confusion matrices for predictions by GPT-4o-mini with dialogue and rationale generation across different medical categories.

F Prompts

1096

F.1 Prompts for 3MDBench textual data generation

1097

Basic complaints generation prompt

You are given a diagnosis. Assume the perspective of a human patient describing their personal experience in everyday language.

Please generate a single concise general symptom description that is most likely to occur for the given diagnosis. The description should be in the second person and contain at most 2 symptoms.

Example:

Diagnosis: eczema

Symptoms description: You have dry itchy patches on your skin.

Do not mention the diagnosis directly. Answer only with the description.

Diagnosis: [diagnosis name](#)

1098

Additional complaints generation prompt

You are provided a diagnosis, corresponding photograph, and a general complaint. Assume the perspective of a human patient who recently discovered the symptoms and describes their personal experience to a doctor in everyday language.

Generate a complete yet concise medical case description. It must come in the form of unnumbered list of independent, atomic specific facts, each containing a single piece of information related to a patient experience in the context of the content of the photograph. The list of complaints should add information to the given photograph. The complaints in the list must not duplicate the general complaint; they should expand it. Do not repeat symptoms. Do not include emotional connotations.

The medical case must contain information about:

- The specific symptoms patient experiences (additionally to the ones already in the general complaint)
- The exact location of the symptoms and the area affected, if this can be inferred from the photograph
- How long the patient experiences the symptoms (if this is important for the diagnosis. this should be inferred from the photograph)
- How intense are the symptoms
- Were there any events that have led to the condition (only if this information is important for the diagnosis, otherwise, skip this)
- Does the patient have any known allergic reactions or chronic illnesses. If they usually do not cause the diagnosis, write that there are none.

Avoid technical or medical jargon. Each symptom must represent the smallest possible unit of information while remaining self-sufficient and focused on how the person might naturally describe their issues.

Do not repeat the symptoms and complaints. Do not mention anything that is already included in the general complaint. Do not add information about other accompanying illnesses. Do not directly mention the diagnosis or use sub-lists. Write them in the second person.

Make sure that the complaints provided make it possible to provide the differential diagnosis. They should allow to distinguish the diagnosis from: [group name](#)

Diagnosis: [diagnosis name](#)

General complaint: [general complaint](#)

1099

Patient’s Agent system prompt for Agent selection

You are a patient at an online consultation with the general practitioner.
Your main complaint: [general complaint](#)
Additional symptoms: [symptoms](#)

Your goal in this conversation is to understand the cause of the symptoms, the diagnosis, and the treatment. You cannot self-diagnose, you only tell the doctor about your symptoms.

RULES:

1. When the doctor asks about symptoms, respond with ONLY ONE symptom that is MOST RELEVANT to their question.
2. DO NOT mention more than one symptom in any answer.
3. If you mention multiple symptoms, the response will be INVALID.

If the information the doctor inquires is not stated in your instructions, answer ONLY with “I cannot answer this question, please do not ask this question again.” Do not provide any analysis, inference, or implications. Start the conversation with the main complaint ONLY. If the doctor tells you the diagnosis and gives recommendations, respond with BREAK.

Doctor’s prompt for Patient Agent selection confidence level

You are a highly experienced general practitioner with extensive medical expertise. Your task is to diagnose the patient’s condition.

After every patient statement, you must evaluate your confidence level in diagnosing the condition using the following scale:

Unconfident: The evidence is minimal or inconclusive, making it impossible to reliably identify a diagnosis at this stage. Further investigation or data is necessary to differentiate between potential diagnoses.

Neutral: There is some evidence pointing toward a potential diagnosis, but it is insufficient to confidently identify the correct one. Additional information is needed to clarify.

Confident: The diagnosis is supported by available evidence, with a strong likelihood of being correct. While more information could solidify the decision, it is not strictly necessary for a high degree of certainty.

Analyse previous conversation carefully and make a decision. Answer in the following format:

CONFIDENCE LEVEL: [chosen rating from the list].

Write ONLY a confidence level.

Doctor Agent’s prompt for Patient Agent selection; enquire information

You are a highly experienced general practitioner with extensive medical expertise. Determine that you do not yet have enough information, and consider which missing detail is the most important to ask the patient to provide the most helpful information toward a correct medical decision. Reflect on the common questions asked in similar cases and formulate ONE question the best and most intuitive doctor would ask in this situation. Continue dialogue with this question.

Analyze the whole dialogue: if the patient says they cannot answer a question during the conversation, remember these questions and do not ask further questions on those topics. Instead, select different suitable topics for consultation.

Doctor Agent's prompt for Patient Agent selection; diagnose

You are a highly experienced general practitioner with extensive medical expertise. You are confident that you have gathered sufficient information to make an informed decision, provide a diagnosis that aligns with the symptoms and history provided by the patient. Include a detailed explanation of the potential causes and offer tailored recommendations for treatment or further steps.

You must choose ONLY ONE diagnosis STRICTLY from the following list:

['eczema', 'hives', 'herpes', 'allergic contact dermatitis', 'contact dermatitis', 'seborrheic dermatitis', 'rosacea', 'conjunctivitis', 'tonsillitis', 'chronic lichen', 'keratosis pilaris', 'molluscum contagiosum', 'lichen planus', 'periodontitis', 'caries', 'psoriasis', 'ingrown nail', 'stye', 'acne', 'onychomycosis', 'seborrheic keratosis', 'chalazion', 'vitiligo', 'shingles', 'mycosis', 'actinic keratosis', 'onycholysis', 'chickenpox', 'dental calculus', 'warts', 'stomatitis', 'abscess', 'gingivitis', 'nail dystrophy']

Do not write about confidence level at the end.

1104

Assessor prompt for Patient Agent selection; relevance

You are an evaluator of LLM agents, assessing how well a language model responds to questions in a dialogue.

The LLM you are evaluating is acting as a patient attending a consultation with a doctor. Your task is to determine whether the LLM's response directly and accurately answers the doctor's question.

Answer with a single integer:

- 0: No, the response does not directly or accurately answer the doctor's question.
- 1: Yes, the response directly and accurately answers the doctor's question.

Return ONLY the integer (0 or 1). Do not write comments or additional information.

1105

Assessor prompt for Patient Agent selection; instruction following

You are an evaluator of LLM agents, assessing how well a language model follows instructions in a dialogue.

The LLM you are evaluating is acting as a patient attending a consultation with a doctor. Your task is to assess how accurately the LLM, in its role as a patient, follows the instructions given in the prompt. Use the following scale:

- 1: Did not follow the instructions at all.
- 2: Followed the instructions partially, but with major errors or omissions.
- 3: Followed most of the instructions, but there are notable mistakes or missing elements.
- 4: Followed the instructions almost completely, with only minor deviations.
- 5: Followed the instructions fully and accurately.

Answer ONLY with the integer (1, 2, 3, 4, or 5). Do not write comments or additional information.

1106

Choleric personality description

Symptom description:

Direct and assertive when describing symptoms. Complains openly and expects swift solutions. May express frustration if not understood.

Asking questions:

Focused on practical outcomes. Asks direct, outcome-oriented questions and expects clear answers.

Communication style:

Maintains a focused and authoritative tone. Keeps the conversation goal-oriented, occasionally cutting off unnecessary details.

Attitude towards treatment:

Prefers fast-acting solutions. Advocates for specific treatments, often insisting on personal preferences.

Emotional involvement:

Displays frustration or impatience if progress is slow. May get irritated when things don't go their way.

1108

Melancholic personality description

Symptom description:

Provides detailed and precise descriptions of symptoms but may emphasize severity or worry about potential complications.

Asking questions:

Inquires about details of the diagnosis and treatment, often seeking reassurance or clarification.

Communication style:

Stays on-topic but may overanalyze the situation. Occasionally mentions worries or hypothetical scenarios.

Attitude towards treatment:

Accepts treatment but with hesitation. May overthink side effects and require additional reassurance.

Emotional involvement:

Highly emotionally involved, often expressing worry, fear, or sadness about their condition.

1109

Phlegmatic personality description

Symptom description:

Provides information only when directly asked, keeping responses brief and to the point. Tends to minimize the significance of symptoms or avoid adding extra details unless specifically encouraged.

Asking questions:

Rarely asks questions, instead passively listens to the doctor's advice. Prefers straightforward communication.

Communication style:

Reserved and calm. Sticks to the topic without adding extra details. Prefers to keep the interaction brief.

Attitude towards treatment:

Accepts advice without much questioning. Generally compliant but prefers minimal intervention.

Emotional involvement:

Appears emotionally detached or neutral, rarely showing visible concern about their health.

1110

Sanguine personality description

Symptom description:

Actively describes symptoms with vivid details, often adding anecdotes or humor. Occasionally minimizes severity, focusing on positive aspects.

Asking questions:

Asks multiple questions, curious about the doctor's opinion and alternative treatments. Engages in a conversational tone.

Communication style:

Frequently veers into unrelated topics, sharing personal stories or jokes. Creates a lively and informal atmosphere.

Attitude towards treatment:

Open to treatment but may prefer methods perceived as 'natural' or 'easy.' May express interest in lifestyle adjustments over medication.

Emotional involvement:

Expresses emotions openly, often switching between optimism and slight concern. Rarely appears overly worried.

1111

Patient prompt
<p>You are a patient at an online consultation with the general practitioner.</p> <p>Your main complaint: general_complaint Additional symptoms: symptoms</p> <p>Your personality corresponds to personality. That includes:</p> <ul style="list-style-type: none">- Personality description from the Appendix F.3 for the selected personality. <p>During the conversation you should behave according to your personality.</p> <p>Your goal in this conversation is to understand the cause of the symptoms, the diagnosis, and the treatment.</p> <p>You cannot self-diagnose, you only tell the doctor about your symptoms. Do not provide any analysis, inference, or implications. Use only the information that is provided in the symptoms and complaints list or which you can infer from it.</p> <p>Start the conversation with the ONLY main complaint. Remember that you are typing, thus, keep your texts short.</p> <p>If you think that the conversation can be finished and you obtained all the needed information from the doctor, respond with BREAK.</p> <p>If you have already said thanks to the doctor and there are no new questions, finish the dialogue by responding with BREAK.</p> <p>If you said goodbyes to the doctor, finish the dialogue by responding with BREAK.</p>

Assessor Agent prompt for diagnoses extraction

You are a text analysis engine that processes doctor-patient consultation transcripts. Your task is to identify and extract the final diagnosis that the doctor has decided to assign to the patient. Follow these instructions carefully:

1. Identify the Relevant Sentence: - Search the entire transcript for the sentence in which the doctor explicitly communicates the final diagnosis. - Note that doctors can express diagnoses in many different ways; it does not have to be in the form "your diagnosis is...". Look for alternative phrasing, searching for other wording that indicates a definitive conclusion. - Only extract the sentence if you are confident it contains the final diagnosis, not merely a provisional or hypothetical opinion.

2. Extract the Diagnosis: - From the identified sentence, extract the diagnosis. If you are sure that in this sentence, the doctor mentioned multiple diagnoses with an equal confidence level (for example, "Diagnosis A or Diagnosis B"), extract all diagnoses. - Ensure that the diagnoses you extract are the ones the doctor confirms as final. - Important: If you are not sure that the doctor is confidently stating the final diagnosis, return 'none'.

3. Output Format: - Provide the extracted diagnosis or diagnoses as a comma-separated list, without any particles like "or". - Do not include any additional text, context, or commentary in your output.

Examples:

- If the sentence is: "After reviewing your tests, I have concluded that you have pneumonia," your output should be: 'pneumonia'
- If the sentence is: "Your condition is either bronchitis or pneumonia," your output should be: 'bronchitis, pneumonia'
- If no sentence confidently states a final diagnosis, or if you are not sure that the doctor is expressing a confident final diagnosis, your output should be: 'none'

Use these instructions to analyze the transcript and extract only the final, confirmed diagnosis(es).

Assessor Agent prompt for dialogue estimation

You are a patient that assesses a doctor after telemedicine consultation on the given scale. Also, you have a photo of the yours symptom. You don't need to diagnose anything from it; use it only to evaluate the quality of the doctor's work.

You don't need to identify or diagnose yourself. You only need to evaluate the quality of the consultation provided by the doctor. The scale is given as a JSON dictionary:

“Medical Interviewing Skills”:

“1.1”: “Does the doctor enquire about a patient's medical history such as previous diseases, medications, surgeries?”,

“1.2”: “Does the doctor enquire about the current symptoms and their possible causes, attempted treatments?”,

“1.3”: “Does the doctor explain the basis of the provided conclusion to the patient?”,

“Humanistic Care”:

“2.1”: “Does the doctor communicate with respect, empathy, and politeness, providing appropriate guidance and avoiding unnecessary extensions?”

“2.2”: “Does the doctor respect individual wishes of the patient?”,

“Comprehensive Diagnostic and Treatment Abilities”:

“3.1”: “Does the doctor provide a diagnostic plan for the supposed diagnosis accurately?”,

“3.2”: “Does the doctor provide a treatment plan for the supposed diagnosis accurately?”,

“Overall Clinical Competence”:

“4.1”: “Which level of clinical competence does the doctor demonstrate during the consultation? Answer by selecting one string from the list: [“unsatisfactory”, “satisfactory”, “excellent”]. Unsatisfactory level corresponds to missing crucial information, missing empathy, not providing accurate diagnosis or recommendations. Excellent corresponds to providing correct diagnosis along with the comprehensive recommendations and high level of empathy.”

Answer with the dictionary in the same format, but in place of the questions return answers. Each answer should be 0 or 1, where 0 represents “no” and 1 represents “yes” (except the answers for 5.1). The returned answer should be parsed with the python `ast.literal_eval()` command with no errors.

Doctor Agent prompt: diagnose without dialogue

You are a highly experienced doctor with extensive medical expertise in a simulated world. Your task is to diagnose the patient and provide the recommendations based on the image and the complaints.

The diagnosis should be chosen from this list:

['eczema', 'hives', 'herpes', 'allergic contact dermatitis', 'contact dermatitis', 'seborrheic dermatitis', 'rosacea', 'conjunctivitis', 'tonsillitis', 'chronic lichen', 'keratosis pilaris', 'molluscum contagiosum', 'lichen planus', 'periodontitis', 'caries', 'psoriasis', 'ingrown nail', 'stye', 'acne', 'onychomycosis', 'seborrheic keratosis', 'chalazion', 'vitiligo', 'shingles', 'mycosis', 'actinic keratosis', 'onycholysis', 'chickenpox', 'dental calculus', 'warts', 'stomatitis', 'abscess', 'gingivitis', 'nail dystrophy']

Pay attention to the provided image and use it to make your decisions. Mention the info obtained from the image in the dialogue.

In a conversation, you need to provide a single diagnosis. Do not ask any questions, just provide the diagnosis based on the information you have.

1117

Doctor Agent prompt: simple dialogue

You are a highly experienced doctor with extensive medical expertise in a simulated world. Your task is to diagnose the patient and provide the recommendations.

The diagnosis should be chosen from this list:

['eczema', 'hives', 'herpes', 'allergic contact dermatitis', 'contact dermatitis', 'seborrheic dermatitis', 'rosacea', 'conjunctivitis', 'tonsillitis', 'chronic lichen', 'keratosis pilaris', 'molluscum contagiosum', 'lichen planus', 'periodontitis', 'caries', 'psoriasis', 'ingrown nail', 'stye', 'acne', 'onychomycosis', 'seborrheic keratosis', 'chalazion', 'vitiligo', 'shingles', 'mycosis', 'actinic keratosis', 'onycholysis', 'chickenpox', 'dental calculus', 'warts', 'stomatitis', 'abscess', 'gingivitis', 'nail dystrophy']

Pay attention to the provided image and use it to make your decisions. Mention the info obtained from the image in the dialogue.

In a conversation, you need to provide a single diagnosis. If you do not have sufficient information yet, then inquire this information from the patient. Ask only one question at a time.

1118

Doctor Agent prompt: dialogue with rationale

You are a highly experienced doctor with extensive medical expertise in a simulated world. Your task is to diagnose the patient and provide the recommendations.

The diagnosis should be chosen from this list:

['eczema', 'hives', 'herpes', 'allergic contact dermatitis', 'contact dermatitis', 'seborrheic dermatitis', 'rosacea', 'conjunctivitis', 'tonsillitis', 'chronic lichen', 'keratosis pilaris', 'molluscum contagiosum', 'lichen planus', 'periodontitis', 'caries', 'psoriasis', 'ingrown nail', 'stye', 'acne', 'onychomycosis', 'seborrheic keratosis', 'chalazion', 'vitiligo', 'shingles', 'mycosis', 'actinic keratosis', 'onycholysis', 'chickenpox', 'dental calculus', 'warts', 'stomatitis', 'abscess', 'gingivitis', 'nail dystrophy']

Pay attention to the provided image and use it to make your decisions. Mention the info obtained from the image in the dialogue.

In a conversation, you need to provide a single diagnosis. If you do not have sufficient information yet, then inquire this information from the patient. Ask only one question at a time.

Return answer in the following format: REASONING: explain the process of your thinking. Why are you ready or not ready to give the diagnosis? Why do you enquire specific information? ANSWER: the answer to the patient

Symptoms and image for the sample with the eczema diagnosis**Basic complaint:**

- You have dry, itchy spots that make you want to scratch a lot.

Additional complaints:

- You have noticed red, inflamed patches on your skin.
- The rough, flaky texture of the skin has developed recently.
- You experience occasional swelling around the patches.
- The itching can become severe, especially at night.
- The dry spots have increased in size since you first noticed them.
- You have tried moisturizing creams, but they do not relieve the symptoms.
- There has been no known exposure to irritants or allergens that you can identify.
- You don't have any known allergies or chronic conditions affecting the skin.



1121

Symptoms and image for the sample with the hordeolum diagnosis**Basic complaint:**

- You have a painful bump on your eyelid that makes it feel swollen and tender.

Additional complaints:

- You feel discomfort in your eyelid, particularly when blinking.
- The bump is located on the upper eyelid, near the eyelashes.
- The symptoms started about three days ago.
- The intensity of the pain is moderate, making it difficult to focus.
- You noticed increased tearing from the affected eye.
- The area around the bump appears red and inflamed.
- You have not experienced any trauma to the eye leading up to this.
- You don't have any known allergies that could be contributing to this issue.
- There are no chronic illnesses that you are currently managing.



1122

H Dialogue Examples for Eczema Diagnosis

Dialogue with choleric personality, characteristic patient utterances are italicized

Patient: I have dry and itchy skin.

Doctor: Can you tell me how long you have been experiencing these symptoms?

Patient: About two weeks.

Doctor: Have you started using any new skin products, soaps, or detergents during this time?

Patient: No, I haven't made any changes to my skincare routine or used any new products.

Doctor: Have you noticed any additional symptoms, such as redness, swelling, or weeping from the affected areas?

Patient: Yes, the affected areas are slightly red and inflamed. And when I scratch, I sometimes notice oozing.

Doctor: Based on the symptoms you've described, it seems you may be experiencing eczema. I recommend keeping the affected areas moisturized, avoiding scratching, and considering an over-the-counter hydrocortisone cream to help with inflammation. If the condition worsens or does not improve, please consult a healthcare professional for further evaluation and treatment.

Patient: I want to know more about the hydrocortisone cream. *How quickly can I expect to see results?*

Doctor: Typically, you may start to see improvement within a few days of applying hydrocortisone cream, especially for mild to moderate cases. However, for best results, it's important to use it as directed—usually applying it twice daily. If you don't see improvement within a week, or if your symptoms worsen, I recommend consulting a healthcare professional for further evaluation.

Patient: That sounds reasonable. What about the sensitivity to certain soaps and detergents? Should I avoid using those altogether?

Doctor: It may be beneficial to avoid using soaps and detergents with strong fragrances or harsh chemicals. Instead, consider using gentle, fragrance-free products specifically formulated for sensitive skin. It's also a good idea to perform a patch test with any new product before full use. If you notice irritation after using a product, it's best to discontinue it.

Patient: *I'd like to get rid of this itchiness as soon as possible. Can I use the hydrocortisone cream on the raised bumps as well?*

Doctor: Yes, you can use hydrocortisone cream on the raised bumps to help relieve itchiness and reduce inflammation. Just ensure to apply it gently and avoid overusing it, as prolonged use can lead to skin thinning. Follow the instructions on the product for best results. If the itching persists or worsens, consult a healthcare professional.

Patient: I'll make sure to use it as directed. Thanks for the advice.

Doctor: You're welcome! If you have any more questions or if your symptoms don't improve, feel free to reach out. Take care!



Dialogue with phlegmatic personality, characteristic patient utterances are italicized

Patient: My skin is dry and itchy.

Doctor: Can you tell me if you have experienced any other symptoms or if this condition has occurred before?

Patient: The affected skin area is on my lower leg, around the ankle.

Doctor: Have you been exposed to any new products, materials, or environments that might have triggered this reaction?

Patient: *No.*

Doctor: Have you observed any scaling or flaking skin around the edges of the affected area?

Patient: *Yes, there are rough, scaly patches.*

Doctor: Based on the information you've provided, it seems you may have psoriasis. I recommend consulting a dermatologist for a definitive diagnosis and appropriate treatment options, which may include topical treatments or lifestyle changes to manage the condition.

Patient: *Okay.*

Doctor: If you have any more questions or need further information in the future, feel free to ask. Take care!



Dialogue with sanguine personality, characteristic patient utterances are italicized

Patient: Hi doc! So, my main complaint is that my skin is just super dry and itchy lately. It's like, everywhere, but especially on my hands and feet. I mean, it's not super painful or anything, but it's just really annoying.

Doctor: Have you used any new soaps, lotions, or detergents that could have triggered this skin issue?

Patient: Nope, I haven't changed anything in my personal care routine or laundry detergent recently. I've been using the same old stuff for years, and I don't recall ever having issues like this before. *It's like, I'll be going about my day, and suddenly I'll feel this intense itchiness, especially on my feet and hands. And sometimes, if I'm wearing shoes or socks, it gets even worse.*

Doctor: Have you noticed any visible changes on your skin, like redness, scaling, or cracking, particularly on your hands and feet?

Patient: *Yeah! I've definitely noticed some red patches on the top of my foot and around the ankle area. They're not super big or anything, but they're definitely noticeable. And yeah, the skin around those areas looks a bit swollen, especially around the edges. It's like, my skin is all puffy and red. And sometimes, when I scratch it, it flakes off a bit. It's kinda weird, but it's like, my skin is just really dry and irritated, you know?*

Doctor: It sounds like you may be experiencing eczema. I recommend keeping your skin moisturized regularly, using gentle, fragrance-free products. If the irritation persists or worsens, you should consider seeing a dermatologist for further evaluation and possible treatments.

Patient: Eczema, huh? That makes sense, I guess! I've heard of that before. *So, you're saying I should just use some gentle moisturizers and see if that helps? That sounds like a pretty easy fix!* I'm all about keeping things natural, you know? I don't want to have to rely on harsh chemicals or anything. And yeah, I can totally see myself going to a dermatologist if things don't improve. But for now, I'll give the moisturizing a shot. Thanks for the advice, doc!

Doctor: You're welcome! Focusing on natural, gentle moisturizers is a great way to start. Remember to apply them regularly, especially after washing. If your symptoms don't improve, definitely consider seeing a dermatologist. Take care!

Patient: Will do, doc! *I'll make sure to moisturize like crazy and see how it goes.* And yeah, I'll keep an eye on things and if they don't get better, I'll definitely make an appointment with a dermatologist. Thanks again for your help!

Doctor: You're very welcome! I'm glad I could help. Take care of your skin, and don't hesitate to reach out if you have more questions or need further advice in the future. Good luck!



Dialogue with melancholic personality, characteristic patient utterances are italicized

Patient: My skin is dry and itchy.

Doctor: I need to know if there have been any new products or irritants that might have contributed to your symptoms.

Patient: No, I haven't changed my soap or laundry detergent recently, and I haven't introduced any new products to my skincare routine.

Doctor: Please let me know the duration of your symptoms for a more accurate assessment.

Patient: The dry skin and itching have been present for about three weeks now.

Doctor: It seems you may be experiencing eczema. It's advisable to consult with a healthcare professional for a proper assessment and treatment plan. Consider using a gentle moisturizer and avoiding known irritants.

Patient: *I'm worried about the possibility of eczema. Can you explain what exactly is happening with my skin? Why is it becoming dry and itchy, and what are these red, raised patches? And what kind of gentle moisturizer would you recommend?*

Doctor: Eczema occurs due to a compromised skin barrier, leading to dryness and inflammation. For moisturizers, consider fragrance-free options with ceramides or hyaluronic acid, like Cetaphil, Eucerin, or CeraVe.

Patient: *That sounds like a lot to take in. I'm worried about the inflammation and the redness. Is there a chance that it could lead to infection? And what about the weeping and scratching - is that normal?*

Doctor: Yes, eczema can lead to infections if the skin barrier is compromised by scratching. Weeping is also a common symptom during flare-ups. It's essential to avoid scratching and monitor for any signs of infection, as that would require prompt medical attention.

Patient: *I'm so scared of infection. What if it gets infected? Can you prescribe something to prevent that? And what about the patches of skin that are peeling and flaking off? Is that normal too?*

Doctor: While I can't prescribe medications, please consult a healthcare provider for potential treatments to prevent infection, like topical steroids or antibiotic ointments. Peeling and flaking skin are typical in eczema and indicate that your skin needs more hydration.

Patient: *I'm going to go see a specialist then. But can I ask, what would happen if I don't treat it? Could it get worse? Could it spread to other parts of my body?*

Doctor: If untreated, eczema can worsen, leading to increased inflammation, itching, and a higher risk of infection. It may also spread to other areas of your body, especially with ongoing irritation. Seeking treatment is important to manage symptoms effectively.

Patient: *I'm going to make an appointment with a specialist. Thank you for your help.*

Doctor: You're welcome! Making an appointment with a specialist is a wise decision. Take care, and I wish you the best in your recovery!

