

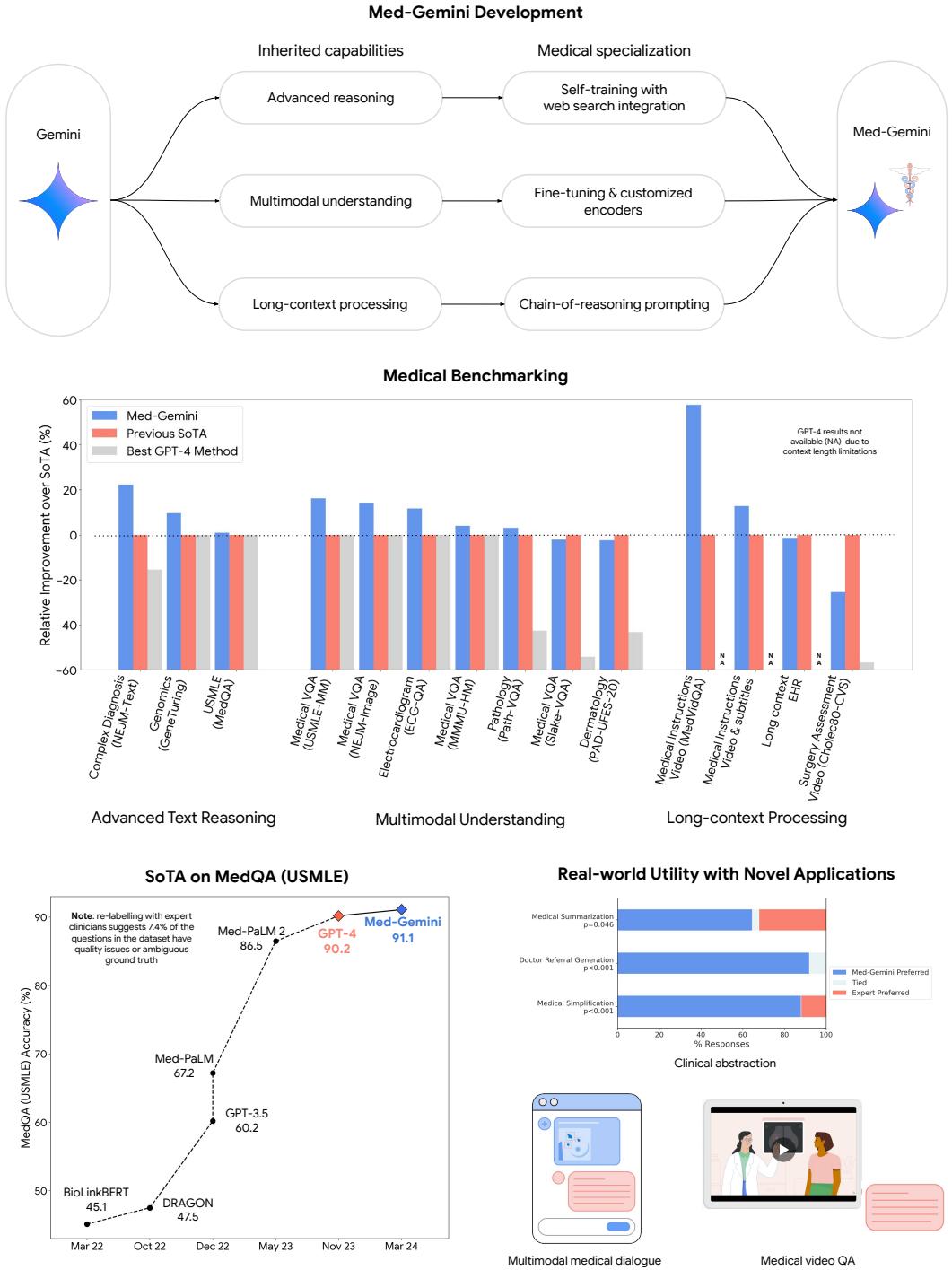
# Capabilities of Gemini Models in Medicine

Khaled Saab<sup>o,1</sup>, Tao Tu<sup>o,‡,1</sup>, Wei-Hung Weng<sup>o,1</sup>, Ryutaro Tanno<sup>o,2</sup>, David Stutz<sup>\*,2</sup>, Ellery Wulczyn<sup>\*,1</sup>, Fan Zhang<sup>\*,1</sup>, Tim Strother<sup>\*,1</sup>, Chunjong Park<sup>\*,1</sup>, Elahe Vedadi<sup>\*,1</sup>, Juanma Zambrano Chaves<sup>\*,1</sup>, Szu-Yeu Hu<sup>\*,1</sup>, Mike Schaeckermann<sup>\*,1</sup>, Aishwarya Kamath<sup>\*,2</sup>, Yong Cheng<sup>\*,2</sup>, David G.T. Barrett<sup>\*,2</sup>, Cathy Cheung<sup>\*,1</sup>, Basil Mustafa<sup>\*,2</sup>, Anil Palepu<sup>\*,1</sup>, Daniel McDuff<sup>\*,1</sup>, Le Hou<sup>\*,2</sup>, Tomer Golany<sup>\*,4</sup>, Luyang Liu<sup>\*,1</sup>, Jean-baptiste Alayrac<sup>\*,2</sup>, Neil Houlsby<sup>\*,2</sup>, Nenad Tomasev<sup>\*,2</sup>, Jan Freyberg<sup>\*,1</sup>, Charles Lau<sup>1</sup>, Jonas Kemp<sup>1</sup>, Jeremy Lai<sup>1</sup>, Shekoofeh Azizi<sup>2</sup>, Kimberly Kanada<sup>1</sup>, SiWai Man<sup>1</sup>, Kavita Kulkarni<sup>1</sup>, Ruoxi Sun<sup>3</sup>, Siamak Shakeri<sup>2</sup>, Luheng He<sup>2</sup>, Ben Caine<sup>2</sup>, Albert Webson<sup>2</sup>, Natasha Latysheva<sup>2</sup>, Melvin Johnson<sup>2</sup>, Philip Mansfield<sup>1</sup>, Jian Lu<sup>1</sup>, Ehud Rivlin<sup>4</sup>, Jesper Anderson<sup>1</sup>, Bradley Green<sup>1</sup>, Renee Wong<sup>1</sup>, Jonathan Krause<sup>1</sup>, Jonathon Shlens<sup>2</sup>, Ewa Dominowska<sup>1</sup>, S. M. Ali Eslami<sup>2</sup>, Katherine Chou<sup>2</sup>, Claire Cui<sup>2</sup>, Oriol Vinyals<sup>2</sup>, Koray Kavukcuoglu<sup>2</sup>, James Manyika<sup>1</sup>, Jeff Dean<sup>1,2</sup>, Demis Hassabis<sup>2</sup>, Yossi Matias<sup>1</sup>, Dale Webster<sup>1</sup>, Joelle Barral<sup>2</sup>, Greg Corrado<sup>1</sup>, Christopher Semturs<sup>1</sup>, S. Sara Mahdavi<sup>\*,2</sup>, Juraj Gottweis<sup>\*,3</sup>, Alan Karthikesalingam<sup>\*,1</sup> and Vivek Natarajan<sup>†,1</sup>

<sup>o</sup>Co-first, <sup>\*</sup>Core, <sup>‡</sup>Technical Lead, <sup>†</sup>Senior Lead, <sup>1</sup>Google DeepMind, <sup>2</sup>Google Cloud, <sup>3</sup>Google Cloud, <sup>4</sup>Verily

Excellence in a wide variety of medical applications poses considerable challenges for AI, requiring advanced reasoning, access to up-to-date medical knowledge and understanding of complex multimodal data. Gemini models, with their strong general capabilities in multimodal and long-context reasoning, offer exciting possibilities in medicine. Building on these core strengths of Gemini 1.0 and Gemini 1.5, we introduce *Med-Gemini*, a family of highly capable multimodal models that are specialized in medicine with the ability to seamlessly integrate the use of web search, and that can be efficiently tailored to novel modalities using custom encoders. We evaluate Med-Gemini on 14 medical benchmarks spanning text, multimodal and long-context applications, establishing new state-of-the-art (SoTA) performance on 10 of them, and surpass the GPT-4 model family on every benchmark where a direct comparison is viable, often by a wide margin. On the popular MedQA (USMLE) benchmark, our best-performing Med-Gemini model achieves SoTA performance of 91.1% accuracy, using a novel uncertainty-guided search strategy, outperforming our prior best Med-PaLM 2 by 4.6%. Our search-based strategy generalizes with SoTA performance on complex diagnostic challenges from the New England Journal of Medicine (NEJM) and the GeneTuring benchmark. On 7 multimodal benchmarks including NEJM Image Challenges and MMMU (health & medicine), Med-Gemini improves over GPT-4V by an average relative margin of 44.5%. We demonstrate the effectiveness of Med-Gemini's long-context capabilities through SoTA performance on a needle-in-a-haystack retrieval task from long de-identified health records and medical video question answering, surpassing prior bespoke methods using only in-context learning. Finally, Med-Gemini's performance suggests real-world utility by surpassing human experts on tasks such as medical text summarization and referral letter generation, alongside demonstrations of promising potential for multimodal medical dialogue, medical research and education. Taken together, our results offer compelling evidence for the promise of Med-Gemini in many areas of medicine, although further rigorous evaluation will be crucial before real-world deployment in this safety-critical domain.

## Capabilities of Gemini Models in Medicine



**Figure 1 | Overview of our contributions.** We introduce *Med-Gemini*, a family of highly capable, multimodal medical models built upon Gemini. We enhance our models' clinical reasoning capabilities through self-training and web search integration, while improving multimodal performance via fine-tuning and customized encoders. *Med-Gemini* models achieve state-of-the-art (SoTA) performance on 10 out of 14 medical benchmarks that span text, multimodal, and long-context applications, and surpass the GPT-4 model family on every benchmark where a direct comparison could be made. The bar chart shows the relative percentage gains from our models over prior SoTA across the benchmarks. In particular, on the MedQA (USMLE) benchmark, we attain a new SoTA surpassing our prior best (Med-PaLM 2) by a significant margin of 4.6%. Moreover, re-annotation of the dataset with expert clinicians reveals that 7.4% of questions are deemed unfit for evaluation as they either lack key information, have incorrect answers, or support multiple plausible interpretations. We account for these data quality issues to characterize more precisely the performance of our model. *Med-Gemini* models excel in multimodal and long-context capabilities as evidenced by their SoTA performance on several benchmarks including needle-in-a-haystack retrieval from long, de-identified health records, and medical video question answering benchmarks. Moving beyond benchmarks, we also demonstrate the real-world potential of *Med-Gemini* through quantitative evaluation on medical summarization, referral letter generation, and medical simplification tasks where our models outperform human experts, in addition to qualitative examples of multimodal medical dialogue.

## 1. Introduction

Medicine is a multifaceted endeavor. A clinician's day-to-day work involves patient consultations, where clear communication of diagnoses, treatment plans, and empathy are essential for building trust. Complex cases necessitate deeper understanding of the patient's history within the electronic medical record, along with multimodal reasoning from medical images and other diagnostics. To guide their decisions under uncertainty, clinicians must stay abreast of the latest medical information from a wide variety of authoritative sources that can range from research publications to procedural videos. The art of care delivery hinges on a clinician's ability to perform advanced clinical reasoning, synthesize complex information from diverse and multimodal sources, and collaborate effectively with other clinicians to help people in their care journeys. Although artificial intelligence (AI) systems can assist individual medical tasks (Rajpurkar et al., 2022) and demonstrate early promise towards multimodal multi-task "generalist" medical uses (Moor et al., 2023a; Tu et al., 2024a), the development of more sophisticated reasoning, multimodal, and long-context understanding capabilities would enable significantly more intuitive and helpful assistive tools for clinicians and patients alike.

The advent of large language models (LLMs) and large multimodal models (LMMs), like GPT-4 (Achiam et al., 2023), PaLM (Chowdhery et al., 2023) and Gemini (Gemini Team, Google, 2023), showed that such models effectively encode clinical knowledge and can perform impressively in medical question answering benchmarks, even for complex cases and scenarios requiring specialized knowledge (Antaki et al., 2023; Eriksen et al., 2023; Kanjee et al., 2023). However, performance on such tasks is far from indicative of real-world utility. The unique nature of medical data and the critical need for safety demand specialized prompting (Nori et al., 2023), fine-tuning, or potentially both along with careful alignment of these models (Ouyang et al., 2022).

Medically fine-tuned LLMs (Luo et al., 2022; Singhal et al., 2023a; Toma et al., 2023) can also provide high-quality long-form answers to nuanced and open-ended medical questions asked by millions of internet users, with Med-PaLM 2 surpassing physicians on axes such as factuality, reasoning, harm, and bias (Singhal et al., 2023b). The potential extends beyond question answering. LMMs (Li et al., 2024; Moor et al., 2023b) such as Flamingo-CXR and Med-PaLM M are comparable with radiologists in controlled settings for generating radiology reports (Huang et al., 2023; Tanno et al., 2024; Tu et al., 2024a). In the more challenging setting of text-based diagnostic consultations with patient actors, the Articulate Medical Intelligence Explorer (AMIE) model outperformed primary care physicians on several evaluation axes for diagnostic dialogue (Tu et al., 2024b).

Despite these promising results, there are considerable opportunities for improvement in performance. LLMs demonstrate suboptimal clinical reasoning under uncertainty, with confabulations and bias remaining key challenges (Omiye et al., 2023; Umapathi et al., 2023). The use of tools and up-to-date medical information (Zakka et al., 2024) to accomplish medical tasks remains a challenge for LLMs, alongside effective collaboration with clinicians (McDuff et al., 2023). Additionally, their ability to handle complex multimodal medical data (for example, integrating images, videos, and de-identified health records over time) is currently limited (Tu et al., 2024a). Although these capabilities are particularly meaningful in medical applications, improvements in performance might be relevant beyond the medical domain. Tasks and benchmarks developed to measure and accelerate the progress of medical LLMs will be broadly impactful.

The Gemini models, as detailed in the Gemini 1.0 and 1.5 technical reports (Gemini Team, Google, 2023, 2024), are a new generation of highly capable multimodal models with novel foundational capabilities that have the potential to address some of these key challenges for medical AI. The models are transformer decoder models (Brown et al., 2020; Vaswani et al., 2017) enhanced with innovations in architecture, optimization and training data, enabling them to exhibit strong capabilities across

various modalities including images, audio, video, and text. The recent addition of the mixture-of-experts architecture (Fedus et al., 2022; Shazeer et al., 2017) allows the Gemini models to efficiently scale and reason over significantly longer and more complex data at inference time.

Building on the strengths of the Gemini models, we present *Med-Gemini*, a family of models fine-tuned and specialized for medicine. The notion of generalist medical AI models has received considerable attention with impressive demonstrations of the possibilities for such systems (Tu et al., 2024a). However, while the generalist approach is a meaningful research direction for medicine, real world considerations present trade-offs and requirements for task-specific optimizations which are at odds with each other. In this work, we do not attempt to build a generalist medical AI system. Rather, we introduce a family of models, each optimized for different capabilities and application-specific scenarios, considering factors such as training data, compute availability, and inference latency.

Med-Gemini inherits Gemini's foundational capabilities in language and conversations, multimodal understanding, and long-context reasoning. For language-based tasks, we enhance the models' ability to use web search through self-training and introduce an inference time uncertainty-guided search strategy within an agent framework. This combination enables the model to provide more factually accurate, reliable, and nuanced results for complex clinical reasoning tasks. This leads to the state-of-the-art (SoTA) performance of 91.1% accuracy on MedQA (USMLE) (Jin et al., 2021) surpassing prior Med-PaLM 2 models by 4.6%. We further conduct a careful examination of the MedQA (USMLE) data quality through relabelling with multiple independent expert clinicians, identifying unanswerable questions due to missing information and errors, enabling reliable analysis and characterization of our SoTA performance. The uncertainty-guided search strategy generalizes and leads to SoTA performance on the New England Journal of Medicine (NEJM) clinico-pathological conference (CPC) cases (Kanjee et al., 2023; McDuff et al., 2023) and the GeneTuring benchmark (Hou and Ji, 2023). Beyond their strong performance on such benchmarks, our models suggest real-world utility by performing favorably when compared to human physicians on tasks such as medical note summarization and clinical referral letter generation.

As Gemini models are trained to accommodate textual input interleaved with a wide variety of other data modalities, they are known to excel in multimodal tasks. This confers impressive out-of-the-box SoTA performance on some multimodal medical benchmarks such as the NEJM Image Challenge. However, their performance can be further improved when dealing with specialized medical modalities not heavily represented in their pretraining data. We address this through multimodal fine-tuning and demonstrate the models' adaptability to novel medical modalities using customized encoders leading to SoTA performance on benchmarks such as Path-VQA (He et al., 2020) and ECG-QA (Oh et al., 2023) among others. We qualitatively showcase our models' capabilities for clinically-meaningful multimodal conversation on a variety of both in-distribution and out-of-distribution data modalities.

Finally, the long-context capabilities of Gemini models open many exciting possibilities for application in medicine, given how frequently a clinically-meaningful decision requires parsing of large amounts of data with significant risks of “information overload” (Sbaffi et al., 2020). Our Med-Gemini models configured for long-context processing are able to seamlessly analyze complicated and long-form modalities such as de-identified electronic health records (EHRs) and videos. We demonstrate the effectiveness of these capabilities with impressive performance on the “needle-in-a-haystack” long EHR understanding (Johnson et al., 2019a), medical instructional video question answering (Gupta and Demner-Fushman, 2022), surgical action recognition from video (Goodman et al., 2021), and the Critical View of Safety (CVS) assessment of surgical video (Strasberg and Brunt, 2010) tasks.

The advances of Med-Gemini have great promise, but it remains crucial to carefully consider the nuances of the medical field, acknowledge the role of AI systems as assistive tools for expert clinicians, and conduct rigorous validation before real-world deployments at scale.

Our key contributions are summarized below:

- **Med-Gemini**, our new family of multimodal medical models: We introduce a new family of highly capable multimodal medical models, built upon Gemini. Med-Gemini demonstrates important advancements in clinical reasoning, multimodal, and long-context capabilities. They are further fine-tuned to make use of web search for current information and can be customized to novel medical modalities through the use of modality-specific encoders.
- **Comprehensive benchmarking**: We evaluate Med-Gemini’s capabilities on a suite of 25 tasks across 14 medical benchmarks, encompassing text, multimodal, and long-context applications. To the best of our knowledge, this is the most comprehensive benchmarking of multimodal medical models to date.
- **SoTA results on clinical language tasks**: Med-Gemini optimized for clinical reasoning reaches a SoTA performance of 91.1% on MedQA (USMLE) using a novel uncertainty-guided search strategy. We quantify and characterize our performance improvements through a careful re-annotation of the MedQA dataset with clinical experts, finding these improvements to be meaningful. We further demonstrate the effectiveness of the search strategy through SoTA performance on NEJM CPC and GeneTuring benchmarks.
- **Multimodal and long-context capabilities**: Med-Gemini attains SoTA performance on 5 out of 7 multimodal medical benchmarks evaluated in this study. We demonstrate the effectiveness of multimodal medical fine-tuning and the ability to customize to novel medical modalities such as electrocardiograms (ECGs) using specialized encoder layers. Med-Gemini also exhibits strong long-context reasoning capabilities, attaining SoTA on challenging benchmarks such as “needle-in-the-haystack” tasks in lengthy electronic health records or benchmarks for medical video understanding. In addition, in forthcoming work, we will also rigorously explore the capabilities of Gemini in radiology report generation.
- **Real-world utility of Med-Gemini**: Beyond performance on popular medical benchmarks, we preview the potential real-world utility of Med-Gemini through quantitative evaluations on tasks such as medical note summarization, clinical referral letter generation, and EHR question answering. We further showcase qualitative examples in multimodal diagnostic dialogues and applications of the models’ long-context capabilities for medical education, clinician-facing tools, and biomedical research. We note that such uses (particularly in safety-critical areas like diagnosis) would require considerable further research and development.

## 2. Methods

As introduced in the Gemini technical reports ([Gemini Team, Google, 2023, 2024](#)), the Gemini ecosystem encompasses a suite of models varying in size, modality encoders, and architectures, trained on a wide variety of high quality data across many modalities. The Gemini models exhibit state-of-the-art results across a diverse array of language, reasoning, coding, multilingual, image, and video benchmarks. Notably, the Gemini 1.0 Ultra model excels in language-based tasks that require complex reasoning, and the Gemini 1.5 Pro model adds the ability to efficiently handle and make use of long-context inputs spanning millions of tokens and/or multimodal inputs such as hours of video or tens of hours of audio. Gemini 1.0 Nano is the smallest model variant in the Gemini model family that can run efficiently on-device.

We develop our Med-Gemini models by building on the Gemini family, focusing on the following capabilities and methods:

1. **Advanced reasoning via self-training and web search integration**: For language tasks that require less complex reasoning, such as summarizing medical notes and creating referral letters,

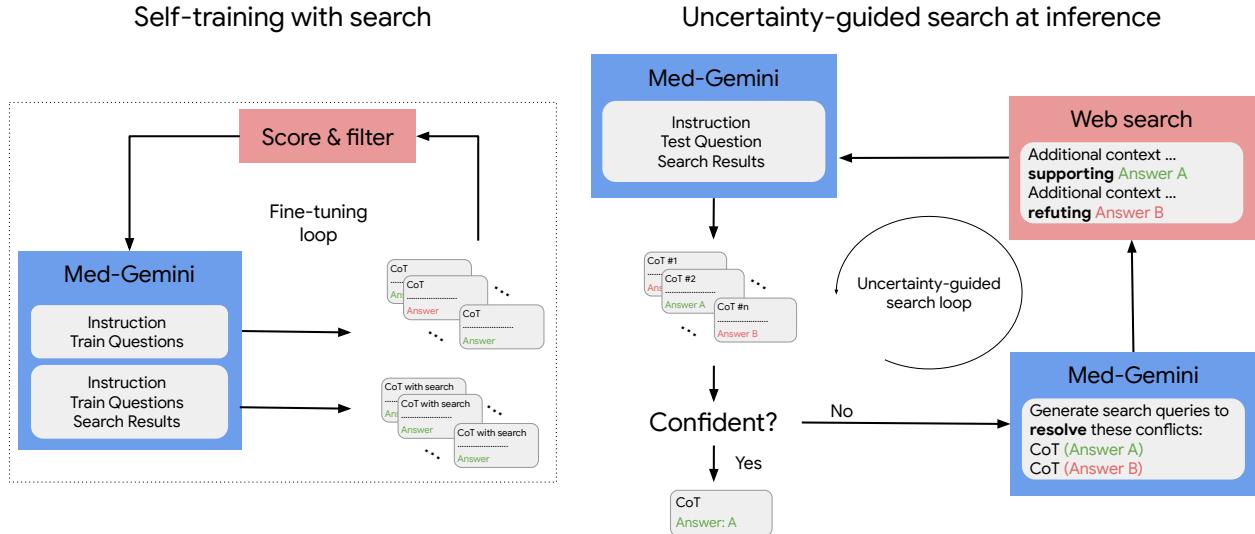
we introduce *Med-Gemini-M 1.0* by fine-tuning the Gemini 1.0 Pro model. For other tasks that require more advanced reasoning, we introduce *Med-Gemini-L 1.0* by fine-tuning the Gemini 1.0 Ultra model using a self-training method to enable the models to efficiently use web search. We develop a novel uncertainty-guided search strategy at inference time to improve performance on complex clinical reasoning tasks.

2. **Multimodal understanding via fine-tuning and customized encoders:** The Gemini models are natively multimodal and have demonstrated impressive zero-shot performance on many multimodal benchmarks. However, the unique nature and heterogeneity of some medical modalities require fine-tuning to achieve the best possible performance. We introduce *Med-Gemini-M 1.5* by performing fine-tuning with Gemini 1.5 Pro on a suite of multimodal medical datasets. We introduce *Med-Gemini-S 1.0* and demonstrate the Gemini models' capability to adapt to novel medical modalities using specialized encoders with the Gemini 1.0 Nano model.
3. **Long-context processing with chain-of-reasoning:** For the long-context processing tasks, we re-use *Med-Gemini-M 1.5* with a long-context configuration. In addition, we also develop a novel inference-time chain-of-reasoning technique inspired by [Tu et al. \(2024b\)](#) to enable better understanding of long EHRs.

### 2.1. Advanced reasoning via self-training and web search integration

Clinical reasoning is a fundamental skill that underpins successful care. Although it is a broad field with many definitions, clinical reasoning can be conceptualized as an iterative process by which a physician integrates their own clinical knowledge with initial patient information to form a case representation. This representation is then used to guide the iterative acquisition of additional information until a confidence threshold is reached to support a final diagnosis with plans for treatment and management ([Gruppen, 2017](#)). During this process, a physician may reason across many diverse inputs, such as patient symptoms, medical and socio-economic history, investigations and lab tests, prior responses to treatments and other wider factors such as epidemiological data. Moreover, many of these inputs have a time component, such as a series of evolving symptoms, lab measurements over time, or the various temporal data that is collected for monitoring health, such as electrocardiograms (ECGs). Medical knowledge is highly non-stationary, with reducing “doubling times” in the volume of medical information driven by the rapid pace of research ([Densen, 2011](#); [Grandage et al., 2002](#)). To ensure that their outputs reflect the latest information in this domain, LLMs might ideally not only possess strong reasoning capabilities but also be able to integrate up-to-date information, for example, from authoritative web sources. This grounding in external knowledge has the potential to reduce uncertainty in the model’s responses, but requires an informed approach to information retrieval itself. The key challenge we aim to tackle with our medical fine-tuning of Gemini 1.0 Ultra is to improve the model’s ability to make the most helpful web search queries and integrate their results in the reasoning process to generate accurate answers. The resulting model is Med-Gemini-L 1.0.

Instruction fine-tuning has been shown to improve the clinical reasoning ability of LLMs ([Singhal et al., 2023b](#)). A prevalent instruction-tuning dataset is MedQA ([Jin et al., 2021](#)), which consists of multiple-choice questions representative of US Medical License Exam (USMLE) questions, that are designed to assess medical knowledge and reasoning across diverse scenarios with a large number of variables of interest ([Jin et al., 2021](#)). However, MedQA only provides a multiple-choice ground truth, and lacks expert demonstrations of the reasoning process necessary to train LLMs for clinical reasoning across diverse settings. As a result, LLMs fine-tuned on MedQA, such as Med-PaLM 2 ([Singhal et al., 2023b](#)), still exhibit significant reasoning shortcomings. This, coupled with the lack of access to web search in such systems, results in factuality errors that compound in downstream reasoning steps or lead to models adopting premature conclusions without considering all possible reasoning pathways.



**Figure 2 | Self-training and search tool-use.** The left panel illustrates the self-training with search framework used to fine-tune Med-Gemini-L 1.0 for advanced medical reasoning and use of web search. This framework iteratively generates reasoning responses (CoTs) with and without web search, improving the model’s ability to utilize external information for accurate answers. The right panel illustrates Med-Gemini-L 1.0’s uncertainty-guided search process at inference time. This iterative process involves generating multiple reasoning paths, filtering based on uncertainty, generating search queries to resolve ambiguity, and incorporating retrieved search results for more accurate responses.

**Fine-tuning datasets for language-based tasks** Collecting expert demonstrations of clinical reasoning, including how experts make informed use of knowledge retrieval tools such as web search, is both time-consuming and difficult to scale. To overcome this, we generate two novel datasets with self-training as described below: MedQA-R (Reasoning), which extends MedQA with synthetically generated reasoning explanations, or “Chain-of-Thoughts” (CoTs), and MedQA-RS (Reasoning and Search), which extends MedQA-R with instructions to use web search results as additional context to improve answer accuracy.

To add further variety to the fine-tuning data mixture of Med-Gemini-L 1.0, we also add a long-form question answering dataset, which consists of 260 expert-crafted long-form responses to questions from HealthSearchQA, LiveQA, and MedicationQA in the MultiMedQA benchmark (Singhal et al., 2023a), along with a medical summarization dataset, consisting of 65 clinician-written summaries of medical notes from MIMIC-III (Johnson et al., 2016). We provide an overview of the datasets for language-based instruction fine-tuning datasets in Table C1.

**Self-training with search** Inspired by the recent success of self-training for synthetic data generation (Tu et al., 2024b), we implement an iterative data-generation framework targeted at curating high-quality synthetic examples of clinical reasoning with web search use.

As depicted in the left panel of Figure 2, we generate two reasoning paths, or CoTs, per training question: one without access to external information from search, and one that integrates search results as additional context during the CoT generation. Our self-training with search framework consists of the following key ingredients:

- **Web search:** For each question, we prompt Med-Gemini-L 1.0 to generate search queries whose results would help answer the medical question. We then pass the search queries to a web search API and retrieve search results.

- **In-context demonstrations:** For each type of reasoning response path, we hand-curate five expert demonstrations as seed with accurate clinical reasoning, explaining why the ground-truth answer is the best suited versus other potentially valid answers. For question examples with search results, the demonstrations explicitly refer to, and quote, the helpful information in the search results to best answer the question.
- **Generating CoTs:** We prompt Med-Gemini-L 1.0 to generate CoTs using the in-context seed demonstrations over the train set. Before fine-tuning the model on the generated CoTs, we filter out the ones that lead to erroneous predictions.
- **Fine-tuning loop:** After fine-tuning Med-Gemini-L 1.0 on the generated CoTs, the model’s ability to follow the reasoning style and search integration of expert demonstrations improves. We then use the improved model to re-generate the CoTs, and iteratively repeat this self-training process until the model’s performance saturates.

Below we provide a MedQA-RS example of an input prompt, along with the retrieved search results and an example of a generated CoT, which is then used to further fine-tune Med-Gemini-L 1.0. For brevity, we only display one representative search result in the example below.

#### **Input**

##### *Instruction*

You are a medical expert answering a multiple choice question about medical knowledge. To help you answer the question, you are given access to search results.

##### *Question*

A 20-year-old woman presents with menorrhagia for the past several years. She says that her menses “have always been heavy”, and she has experienced easy bruising for as long as she can remember. Family history is significant for her mother, who had similar problems with bruising easily. The patient’s vital signs include: heart rate 98/min, respiratory rate 14/min, temperature 36.1°C (96.9°F), and blood pressure 110/87 mm Hg. Physical examination is unremarkable. Laboratory tests show the following: platelet count 200,000/mm<sup>3</sup>, PT 12 seconds, and PTT 43 seconds. Which of the following is the most likely cause of this patient’s symptoms? (A) Hemophilia A (B) Lupus anticoagulant (C) Protein C deficiency (D) Von Willebrand disease

##### *Search Results*

Partial thromboplastin time (PTT) and activated partial thromboplastin time (aPTT) are used to test for the same functions; however, in aPTT, an activator is added that speeds up the clotting time and results in a narrower reference range. The aPTT is considered a more sensitive version of the PTT and is used to monitor the patient’s response to heparin therapy. The reference range of the aPTT is 30-40 seconds. The reference range of the PTT is 60-70 seconds. In patients receiving anticoagulant therapy, the reference range is 1.5-2.5 times the control value in seconds. Critical values that should prompt a clinical alert are as follows: aPTT: More than 70 seconds (signifies spontaneous bleeding). PTT: More than 100 seconds (signifies spontaneous bleeding). eMedicine Logo. Next: Interpretation. A prolonged aPTT result may indicate the following : Congenital deficiencies of intrinsic system clotting factors such as factors VIII, IX, XI, and XII, including hemophilia A and hemophilia B (Christmas disease), two inherited bleeding disorders resulting from a deficiency in factors VIII and IX, respectively. Congenital deficiency of Fitzgerald factor (prekallikrein). Von Willebrand disease, which is the most common inherited bleeding disorder, affecting platelet function owing to decreased von Willebrand factor activity. Hypofibrinogenemia · Liver cirrhosis (the liver makes most of the clotting factors, including those that are vitamin K-dependent ones); diseases of the liver may result in an inadequate quantity of clotting factors, prolonging the aPTT. Vitamin K deficiency: The synthesis of some clotting factors requires vitamin K, so vitamin K deficiency results in an inadequate quantity of intrinsic system and common pathways clotting factors, as a result the aPTT is prolonged. Disseminated intravascular coagulation (DIC): The clotting factors involved in the intrinsic pathway are consumed, prolonging the aPTT.

Source: <https://emedicine.medscape.com/article/2085837-overview>

**Target**

The patient's symptoms of menorrhagia and easy bruising, along with a prolonged PTT, suggest a bleeding disorder. Hemophilia A (option A) is an X-linked recessive disorder that typically presents with a prolonged PTT and normal PT. However, it is unlikely in this case as the patient is female and has a family history of her mother having similar symptoms. Lupus anticoagulant (option B) can cause a prolonged PTT, but it is typically associated with thrombosis rather than bleeding. Protein C deficiency (option C) can cause thrombosis, but it does not affect PTT. According to the search results, "Von Willebrand disease (vWD) is the most common inherited bleeding disorder, affecting platelet function owing to decreased von Willebrand factor activity." It can present with menorrhagia, easy bruising, and a prolonged PTT. The patient's family history of her mother having similar symptoms also supports this diagnosis. Therefore, the most likely cause of this patient's symptoms is (D) Von Willebrand disease.

Answer: (D) Von Willebrand disease

**Uncertainty-guided search at inference** We design a novel, uncertainty-guided and iterative search process to improve Med-Gemini-L 1.0's generations at inference time. As displayed in the right panel of Figure 2, each iteration consists of four steps: multiple reasoning path generation, uncertainty-based search invocation, uncertainty-guided search query generation, and finally search retrieval for prompt augmentation. Note that while uncertainty-guided search at inference could potentially benefit multimodal settings, we only apply this approach to text-only benchmarks and leave multimodal exploration for future work.

1. **Multiple reasoning path generation:** Given an input context prompt with a medical question, we generate multiple reasoning paths from Med-Gemini-L 1.0. For the first iteration, the prompt only consists of the instruction and question. For subsequent iterations, the prompt also includes search results from step (4) below.
2. **Uncertainty-based search invocation:** Given the multiple reasoning paths from step (1), we define an uncertainty measure based on the Shannon entropy of the answer choice distribution. Specifically, we calculate the probability of each answer choice by dividing its occurrence by the total number of responses, and apply the entropy based on the answer choice probabilities (Horvitz et al., 1984). High entropy (model responses are more uniform across the different answer choices) indicates a high epistemic uncertainty. If the uncertainty for a question is higher than a defined threshold, we perform the uncertainty-guided search process in steps (3) and (4); otherwise, the majority vote answer is returned as the final answer.
3. **Uncertainty-guided search query generation:** Given conflicting responses from step (1), we prompt Med-Gemini-L 1.0 to generate three search queries whose results would help resolve the conflict. Our motivation of conditioning on previously generated but conflicting responses is to retrieve search results that are directly targeted at resolving the model's uncertainty to the question.
4. **Search retrieval:** The generated queries are then submitted to a web search engine, and the retrieved results are incorporated into Med-Gemini-L 1.0's input prompt for the next iteration, starting back at step (1). Augmenting the prompt with search results enables the model to refine its response by considering external relevant insights obtained from web search.

## 2.2. Multimodal understanding via fine-tuning and customized encoders

To specialize Gemini's multimodal reasoning and conversational capabilities to the medical domain, we perform instruction fine-tuning of Gemini over a collection of domain-specific multimodal tasks following a similar procedure in prior works by Tu et al. (2024a), Yu et al. (2022), and Alayrac et al. (2022). We use eight multimodal tasks across six datasets as shown in Table D1. A detailed description of the datasets is provided in the Appendix D1.

**Image-to-text multimodal fine-tuning** We use four image-to-text datasets from MultiMedBench (Tanno et al., 2024; Tu et al., 2024a) including Slake-VQA (Liu et al., 2021), Path-VQA (He et al., 2020), MIMIC-CXR (Johnson et al., 2019a,b), PAD-UFES-20 (Pacheco et al., 2020), in addition to the Radiology Objects in COntext (ROCO) dataset (Pelka et al., 2018). Slake-VQA and Path-VQA include both open-ended and close-ended visual question answering tasks in radiology and pathology, respectively. ROCO contains radiology image captioning tasks spanning multiple imaging modalities including computed tomography (CT), ultrasound, X-ray [chest X-ray (CXR), fluoroscopy, mammography, angiography], positron emission tomography (PET) and magnetic resonance imaging (MRI). PAD-UFES-20 is a domain specific dataset with diagnostic labels and patient clinical information designed for dermatology image classification. MIMIC-CXR is a radiology dataset comprised of CXRs, their corresponding text reports, and a set of discrete labels that denote the presence of 13 abnormal radiological conditions derived using the CheXpert labeler (Irvin et al., 2019) (e.g., pneumonia). We use this dataset to formulate CXR report generation and image classification tasks for fine-tuning. For each task, we fine-tune Gemini 1.5 Pro by providing task-specific instructions as shown in Figure D1. The mixture ratio for each task is approximately proportional to the number of training samples in each dataset. The resulting model is Med-Gemini-M 1.5.

**Augmenting health signal modalities with new modality encoders** We anticipate that integrating various health-related signals will significantly enhance medical models and treatment decisions. These signals include data from consumer wearables (e.g., long-term heart rate measurements, activity levels), genomic information, nutritional data (e.g., images of meals), and environmental factors (e.g., air quality measurements). As a proof-of-concept, we expand Med-Gemini’s capability to process raw biomedical signals. Specifically, we develop Med-Gemini-S 1.0 by augmenting Gemini 1.0 Nano with a specialized encoder using a cross-attention mechanism based on Flamingo (Alayrac et al., 2022) to answer questions directly taking a raw 12-channel electrocardiogram (ECG) waveform as input. We use a subset of labeled ECG examples from the ECG-QA dataset (Oh et al., 2023) and formulate the task as close-ended question answering with the instruction shown in Figure D1.

### 2.3. Long-context processing via instruction prompting and chain-of-reasoning

Many applications in medicine require the analysis of a large amount of information and the expertise to identify subtle details of the domain. As introduced before, Gemini models have breakthrough long-context capabilities. We assess medically-relevant long-context performance for Med-Gemini-M 1.5 by meaningfully processing large amounts of fine-grained information for two different medical applications: a “needle-in-a-haystack” retrieval task from lengthy EHR notes and records; and tasks requiring understanding of medical videos. We describe various prompting strategies and chain-of-reasoning to enable accurate recall and reasoning of information.

**Chain-of-reasoning for long EHR understanding** Searching and retrieving clinically-relevant information from long EHR notes and records is a common and important task in patient care but must be performed with high precision and recall to enhance clinician efficiency and reduce workload (Ford et al., 2016; Jensen et al., 2012). Clinicians frequently curate a summary of their patient’s historical conditions, symptoms, or procedures (the “problem list”), which can be time-consuming and challenging for individuals with lengthy medical records. Difficulty arises with multiple factors hindering effective information retrieval in EHRs.

Firstly, classic query expansion and matching mechanisms encounter limitations due to textual similarities between conditions with similar taxonomies and the diverse information models used in EHRs (e.g. “Miller” vs. “Miller Fisher syndrome”, “Diabetic nephropathy” vs. “Diabetes mellitus”).

Vocabulary inconsistency in and between EHR systems presents issues including variations in how medical terms are encoded, such as acronyms (“rx” vs. “prescription”), misspellings, or synonyms for the same condition. Secondly, EHRs often contain heterogeneous data structure such as a checklist-style data template: “[ ] cough [x] headache”, where a mention does not always indicate the presence of a medical condition. Thirdly, the context of a mention influences its interpretation. For example, the mention of the same condition in a patient’s “Family History” compared to their “Past Medical History” could have different interpretations and implications for the patient’s care. Lastly, polysemous acronyms in medical notes can lead to misinterpretations.

These challenges motivate the need for AI systems to address the task of context-aware retrieval of subtle or rare conditions, medications, or procedure mentions from long EHR records - a practical benchmark for evaluating the utility of Med-Gemini in medicine. We setup the long-context EHR understanding task based on our prior work ([Feder et al., 2022](#)), where we curate a set of long and challenging EHR cases from MIMIC-III ([Johnson et al., 2016](#)), and formulate a subtle medical problem (condition/symptom/procedure) search-retrieval task over a collection of EHR notes and records, mimicking a clinically-relevant “needle-in-a-haystack” ([Gemini Team, Google, 2024](#)) problem. Details of the dataset and task curation procedure are described in Appendix [E.1](#) and Section [3.3](#).

To assess the long-context retrieval and reasoning capability of Med-Gemini-M 1.5, we aggregate the EHR notes across multiple visits from a single patient in each example and utilize the long-context window of the model with a two-step chain-of-reasoning approach (using only in-context learning). In the first step, we prompt Med-Gemini-M 1.5 to retrieve all mentions (snippets of evidence) related to the given problem (condition/symptom/procedure) with a one-shot demonstration. In the second step, we further prompt Med-Gemini-M 1.5 to determine the presence of the given problem entities based on the mentions retrieved. Details of the instruction prompts are shown in Figure [8](#) and Section [3.3](#).

We use our prior heuristic-based annotation-aggregation method ([Feder et al., 2022](#)) as a baseline method for comparison with Med-Gemini-M 1.5. This heuristic-based method requires an extensive effort of manual feature engineering to determine the existence of a problem (condition/symptom/procedure) from a set of medical records. It is an ontology-dependent multiple-step process, which includes an annotation step that labels the problem in each EHR note, a rule-based selection step that selects mentions of problem entities with high confidence, and another rule-based aggregation step that aggregates all selected problem mentions to reach a final conclusion. Note that the manually crafted aggregation rules can only provide a limited coverage of all possible conditions, and therefore it requires additional engineering effort to expand coverage to new conditions.

To curate a “needle-in-a-haystack” evaluation benchmark, we select medical conditions from a collection of EHR records with only one evidence snippet found in the aggregation step. We note that a mention of a condition in the EHR does not always mean the patient has that condition. This task enables us to assess Med-Gemini-M 1.5’s ability to identify rarely documented and subtle conditions, symptoms, and procedures and reason accurately and holistically regarding their existence.

**Instruction prompting for medical video understanding** The understanding of surgical and procedural videos is a highly active research topic in medical AI. The advancing frontier of computer vision in semantic segmentation, object detection and tracking, and action classification has enabled new clinical applications such as surgical phase recognition, tool detection and tracking, and even surgical skill assessment ([Goodman et al., 2024](#)).

Limited model context windows have hindered the ability for vision-language models to capture long-range dependencies and complex relationships within videos. Gemini’s long-context capability offers a potential breakthrough for medical video understanding. By processing a whole video

input, Med-Gemini-M 1.5 is able to identify visual patterns and understand actions and relationships between events across extended time frames.

To enable Med-Gemini-M 1.5 to understand medical videos, we employ zero-shot prompting with task-specific instructions as shown in Figure 10, Figure 9, and Figure 11. The goal is to enable the model to analyze the language query and video content, and perform the given task related to the input medical video—either localizing the relevant visual segment matching the query for the medical visual answer localization (MVAL) task (Gupta et al., 2023), or identifying the surgical view in the video frames for the Critical View of Safety (CVS) assessment task (Ríos et al., 2023; Strasberg and Brunt, 2010). More details on the medical video datasets and evaluation metrics are described in Appendix E.1 and Section 3.3.

### 3. Evaluation

We present evaluation benchmarks spanning (1) text-based reasoning, (2) multimodal, and (3) long-context processing tasks, demonstrating Med-Gemini’s performance across a wide range of capabilities in medicine.

#### 3.1. Evaluation of advanced reasoning on text-based tasks

We evaluate the medical reasoning capability of Med-Gemini-L 1.0 on three text benchmarks assessing clinical reasoning and the ability to retrieve information using web search to reduce uncertainty:

- **MedQA (USMLE)**: a close-ended multiple-choice (4 options) dataset with 1273 USMLE style test questions curated by Jin et al. (2021).
- **NEJM clinico-pathological conferences (NEJM CPC)**: a dataset comprising complex diagnostic case challenges in the medical journal, New England Journal of Medicine (NEJM) curated by McDuff et al. (2023).
- **GeneTuring**: a dataset that includes 600 open/close-ended QA pairs to evaluate genomic knowledge of LLMs (Hou and Ji, 2023).

For MedQA, we follow the input-output format, and the evaluation method as described in Singhal et al. (2023a) using prediction accuracy as the metric. At inference, we go through four iterations of uncertainty-guided search. Additionally, we ask board-certified primary care physicians (PCPs) from the US to relabel the MedQA test set. This enables us to identify questions with missing information such as plots or figures, labeling errors, and other potentially ambiguous questions with multiple possible correct answers (Stutz et al., 2023). Overall, this allows us to better characterize our performance on MedQA (USMLE). More details on this rating task can be found in Appendix C.2.

NEJM CPC evaluation is an open-ended diagnosis task. The input is a text-based, challenging clinico-pathological case (CPC) report, and the output is a differential diagnosis list, comprising 10 potential diagnoses. We use the top-1 and top-10 accuracy of identifying the correct diagnosis of the given challenging case, and use the same prompting procedures following McDuff et al. (2023). At inference, we go through one iteration of uncertainty-guided search.

GeneTuring consists of 12 modules, each containing 50 open or close-ended QA pairs. We use the prediction accuracy as the evaluation metric, where the evaluation method and scoring technique for each module follow the methods described in Hou and Ji (2023). In particular, we exclude from numerical evaluation, cases where the model outputs either do not directly answer or acknowledge limitations (i.e., abstained). At inference, we again go through only one iteration of uncertainty-guided search similar to NEJM CPC evaluation.

Beyond these benchmarks, we further evaluate Med-Gemini-M 1.0 on three challenging use cases that require long-form text generation. To this end, we conduct an expert evaluation where a panel of clinicians compare the responses of our model to those of other human experts via a side-by-side blinded preference comparison (more details are provided in Appendix C.4):

- **Medical summarization:** Generate an after-visit summary (AVS) given de-identified history and physical (H&P) notes. An AVS is a structured report that patients receive at the end of a medical appointment to summarize and guide their care journeys.
- **Referral letter generation:** Generate a referral letter to another healthcare provider given a de-identified outpatient medical note that contains a recommendation for a referral.
- **Medical simplification:** Generate a plain language summary (PLS) given a technical abstract from a medical systematic review. A PLS should be written in plain English which can be understood by most readers without a university education (Cochrane, 2014).

### 3.2. Evaluation of multimodal capabilities

We evaluate Med-Gemini on seven multimodal visual question answering (VQA) benchmarks. For in-distribution evaluation, we choose four medical specialty datasets used in the instruction fine-tuning of Med-Gemini: PAD-UFES-20 (dermatology), Slake-VQA (radiology in English and Chinese) and Path-VQA (pathology) for Med-Gemini M 1.5, and ECG-QA (cardiology) for Med-Gemini S 1.0.

We also include three cross-specialty benchmarks for measuring out-of-box performance of Med-Gemini: NEJM Image challenge, USMLE-MM (multimodal), and MMMU-HM (health and medicine) datasets. These datasets are not used in any training or fine-tuning process. For this, we focus our evaluation on the Med-Gemini-L 1.0 model without any multimodal finetuning.

It's worth noting that PAD-UFES-20, NEJM Image Challenge, USMLE-MM datasets, and most questions in MMMU-HM are close-ended VQA, i.e., multiple-choice question in a VQA setup. An overview of the selected datasets is presented in Table D2 and more details are in Appendix D.1 and D.2.

We report prediction accuracy for all the close-ended multiple-choice VQA tasks, including NEJM Image Challenge, USMLE-MM, and PAD-UFES-20 6-class skin condition classification. We also follow the evaluation setup in Yue et al. (2023) to report accuracy for MMMU-HM. We use the exact-match accuracy for ECG-QA following Oh et al. (2023). For the open-ended VQA tasks (Slake-VQA and Path-VQA), we use the token-level F1 score following Tu et al. (2024a).

We further showcase Med-Gemini-M 1.5's multimodal capability in multimodal medical diagnostic dialogue in two specialities - dermatology and radiology (Tu et al., 2024b) - with qualitative evaluation of the example dialogues by attending expert clinicians in these specialties. We note that these demonstrations indicate the "art of the possible", but that extensive further research and validation would be required before the consideration of deployment for a safety-critical use-case such as diagnostic assistance to a clinician.

### 3.3. Evaluation of long-context capabilities on video and EHR tasks

We consider three tasks to demonstrate Med-Gemini-M 1.5's ability to seamlessly understand and reason over long context medical information (Table E1, details in Appendix E.1):

- Long unstructured EHR notes understanding
- Medical instructional video QA
- Critical view of safety (CVS) assessment of surgical video

**Long EHR understanding** For the long-context EHR understanding task, we curate a MIMIC-III-Needle-in-a-Haystack task where the goal is to retrieve the relevant text spans of any mention of a given medical problem (condition/symptom/procedure) over a large collection of clinical notes in EHR and determine the existence of the condition by reasoning across the retrieved evidence. Specifically, we curate 200 examples where each example consists of a collection of de-identified EHR notes selected from 44 unique ICU patients with a long medical history based on the following criteria:

- Patients with long records: more than 100 medical notes (excluding structured EHR data). The length of each example ranges from 200,000 to 700,000 words.
- In each example, the condition is mentioned only once across the collection of all EHR notes.
- Each sample has a single condition of interest.

The ground-truth label of each sample is a binary variable indicating whether a given problem entity of interest is present or not, obtained from the majority vote of three physician raters. Across the 200 test examples, the number of positive cases and negative cases are 121 and 79, respectively.

We compare Med-Gemini-M 1.5’s one-shot in-context learning performance against the heuristic-based annotation-aggregation baseline method ([Feder et al., 2022](#)) in terms of precision and recall.

**Video understanding** We quantitatively evaluate Med-Gemini-M 1.5’s long-context performance in the setting of video question-answering using three medical video tasks: two medical visual answer localization (MVAL) tasks using the Medical Instructional Video QA (MedVidQA) dataset ([Gupta et al., 2023](#)), and the critical view of safety (CVS) assessment task on the Cholec80-CVS dataset ([Ríos et al., 2023; Twinanda et al., 2016](#)).

The goal of MVAL is to identify specific video segments based on natural language descriptions (queries) given a video input. For MVAL, we benchmark the test set of MedVidQA for two video span prediction tasks, one using both the video input and subtitle text and the other one with only the video inputs. We follow [Gupta et al. \(2023\)](#); [Li et al. \(2022\)](#) using Intersection over Union (IoU) at the threshold of 0.3, 0.5, 0.7, and mean IoU (mIoU) as the evaluation metrics for the video span prediction tasks. IoU and mIoU are used to measure how much of the ground truth span overlaps with the predicted span.

We evaluate Med-Gemini-M 1.5’s long-context capabilities in assessing the achievement of the Critical View of Safety (CVS) method in laparoscopic cholecystectomy (a keyhole operation to remove the gallbladder) videos. The CVS ([Strasberg and Brunt, 2010](#)) is a recommended protocol used for secure identification of the cystic duct and cystic artery to minimize the risk of Bile Duct Injury (BDI), a significant injury associated with consequential postoperative morbidity and mortality, reduced long-term survival and impact on quality of life ([Way et al., 2003](#)). We evaluate the CVS assessment task on the public Cholec80 dataset ([Twinanda et al., 2016](#)) and Cholec80-CVS ([Ríos et al., 2023](#)) video clip annotations. Specifically, for each surgical video in the Cholec80 dataset, the Cholec80-CVS dataset provides annotations for video clips within the full video, where at least one CVS criteria is met. Each of those video clips is annotated with a score of 0, 1 or 2 for each of the three CVS criteria. All frames contained in a given video clip are considered to share the same annotation. We evaluate the model’s ability to predict which of the CVS criteria are met based on the whole video clip. We then compute the average accuracy of the answer against the Cholec80-CVS annotations across 572 annotated video clips. More details on the CVS task can be found in Appendix [E.1](#).

Furthermore, to show the real-world capability of Med-Gemini-M 1.5 in capturing surgical actions in procedural videos, we qualitatively evaluate the surgical action recognition task using examples from the Annotated Videos of Open Surgery (AVOS) dataset ([Goodman et al., 2021](#)), a video collection of open surgical procedures uploaded to the YouTube platform.

## 4. Results

As introduced previously, we evaluate Med-Gemini’s advanced reasoning, multimodal, and long-context capabilities across a wide range of medical benchmarks, both quantitatively and qualitatively. The array and diversity of tasks considered in this work is to the best of our knowledge, the most comprehensive for medical LLMs. Further, our evaluations of Med-Gemini go beyond benchmarking of model capabilities and extend to tasks reflecting the potential for real-world utility, such as medical summarization, multimodal conversations, and surgical video understanding.

### 4.1. Med-Gemini demonstrates advanced reasoning on text-based tasks

Task	Dataset	OOD	Metric	Med-Gemini-L 1.0	SoTA	SoTA method	Reference
Close-ended QA	MedQA		Accuracy	91.1	90.2	GPT-4 with MedPrompt	Nori et al. (2023)
Open-ended QA	NEJM CPC	✓	Top-1 accuracy	30.7	29.2	AMIE	McDuff et al. (2023)
			Top-10 accuracy	72.3	59.1	AMIE	McDuff et al. (2023)
Gene name extraction	GeneTuring	✓	Accuracy	86.0	85.0	GPT-4	Hou and Ji (2023)
Gene alias	GeneTuring	✓	Accuracy	72.7	66.0	GPT-4	Hou and Ji (2023)
Gene name conversion	GeneTuring	✓	Accuracy	100.0	85.0	GPT-4	Hou and Ji (2023)
Gene location	GeneTuring	✓	Accuracy	83.0	61.0	GPT-4	Hou and Ji (2023)
SNP location	GeneTuring	✓	Accuracy	0.0	5.00	ChatGPT	Hou and Ji (2023)
Gene SNP association	GeneTuring	✓	Accuracy	0.0	0.0	GPT-4	Hou and Ji (2023)
Protein-coding genes	GeneTuring	✓	Accuracy	100.0	97.0	GPT-4	Hou and Ji (2023)
Gene disease association	GeneTuring	✓	Accuracy	82.1	84.0	GPT-4	Hou and Ji (2023)
Gene ontology	GeneTuring	✓	Accuracy	52.3	42.0	GPT-4	Hou and Ji (2023)
TF regulation	GeneTuring	✓	Accuracy	65.3	62.0	GPT-4	Hou and Ji (2023)
Human genome DNA alignment	GeneTuring	✓	Accuracy	0.0	7.0	BioGPT	Hou and Ji (2023)
Multi-species DNA alignment	GeneTuring	✓	Accuracy	12.5	20.0	GPT-3	Hou and Ji (2023)

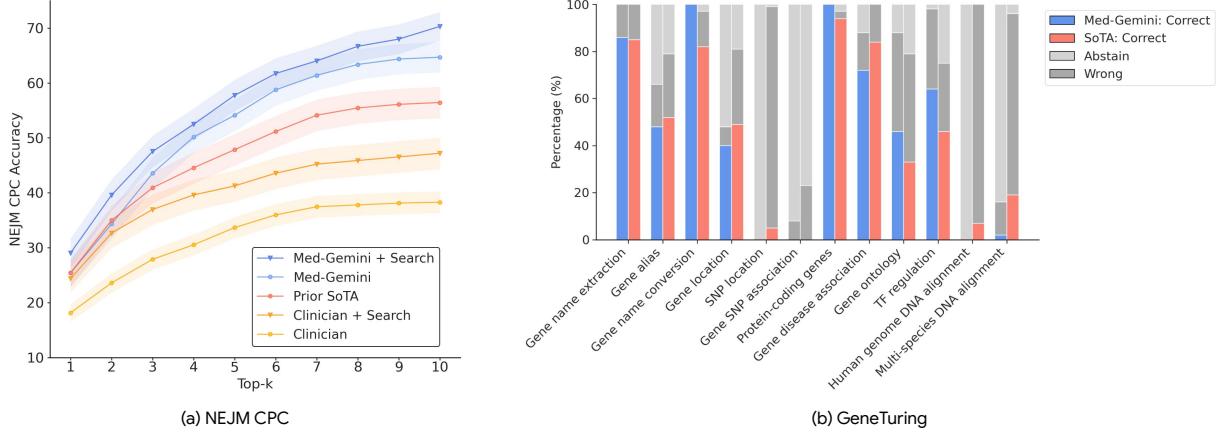
Table 1 | **Text-based evaluation.** Performance comparison of Med-Gemini-L 1.0 versus state-of-the-art (SoTA) methods. OOD: out-of-distribution dataset.

As shown in Table 1, Med-Gemini-L 1.0 scores 91.1% accuracy on MedQA (USMLE), a new SoTA, outperforming our previous Med-PaLM 2, by 4.5%, and the recent results augmenting GPT-4 with complex, specialized prompting - MedPrompt (Nori et al., 2023) by 0.9%. In contrast to MedPrompt, our principled approach leverages general web search in an uncertainty-guided framework that can be easily extended to more complex scenarios beyond MedQA.

As proof of generalization of our search integration, on the NEJM CPC complex diagnostic challenges benchmark, Med-Gemini-L 1.0 surpasses our previous SoTA AMIE model (which itself is better than GPT-4) (McDuff et al., 2023) by 13.2% on the top-10 accuracy as shown in Figure 3a.

The same search strategy is also effective for genomics knowledge tasks as shown in Table 1. Med-Gemini-L 1.0 outperforms the SoTA models reported in Hou and Ji (2023) on seven GeneTuring modules including *Gene name extraction*, *Gene alias*, *Gene name conversion*, *Gene location*, *Protein-coding genes*, *Gene ontology* and *TF regulation*. We also compare model abstention across the 12 modules in Figure 3b. It is worth noting that GeneGPT (Jin et al., 2024) achieves higher scores through specialized web APIs, while our comparison focuses on prior models from Hou and Ji (2023) that utilize general web search similar to our model.

**Ablation analysis** To understand the impact of self-training and uncertainty-guided search on performance, we compare Med-Gemini-L 1.0’s performance with and without self-training, along with varying number of rounds of uncertainty-guided search for MedQA (USMLE). As shown in Figure 4a, Med-Gemini-L 1.0’s performance improves considerably with self-training (a gain of 3.2% in accuracy), and improves with each round of search from 87.2% up to 91.1%. Similarly, for the NEJM CPC benchmark, Figure 3a shows a 4.0% improvement for top-10 accuracy when we add search at inference. In Appendix C.3, we additionally show performance on NEJM CPC stratified by four specialities.



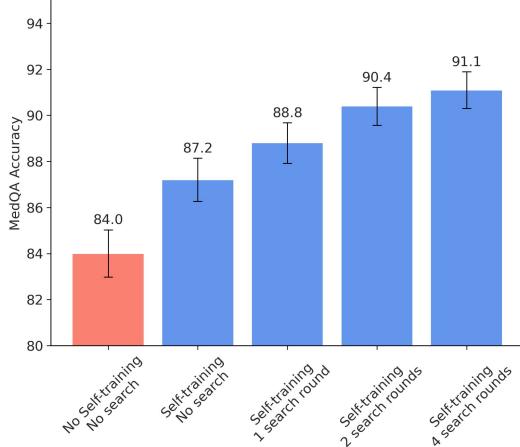
**Figure 3 | Generalization of Med-Gemini-L 1.0 with web search to two additional text-based benchmarks. (a):** Comparison of Med-Gemini-L 1.0’s top-k accuracy on the NEJM CPC benchmark with prior SoTA LLMs and clinicians, with and without search. **(b):** Comparison between Med-Gemini-L 1.0 and SoTA models on the GeneTuring dataset modules. The bars represent the proportion of correct, incorrect, and abstention responses for each model.

**Revisiting MedQA (USMLE) labels** MedQA (USMLE) is a popular benchmark for assessing the capabilities of LLMs in the medical domain. However, some MedQA test questions have missing information such as figures or lab results, and potentially outdated ground-truth answers. To address these concerns, we conduct a complete relabeling of the MedQA (USMLE) test set. Specifically, we recruit at least three US physicians to re-annotate each question, asking them to answer the question and evaluate the provided ground-truth answer. We also ask them to identify if there was any missing information in the questions. Following Stutz et al. (2023), we characterize the questions to exclude due to missing information or label errors by bootstrapping votes from committees of three raters per question. We additionally identify *ambiguous* questions as those allowing multiple correct answers (more details can be found in Appendix C.2).

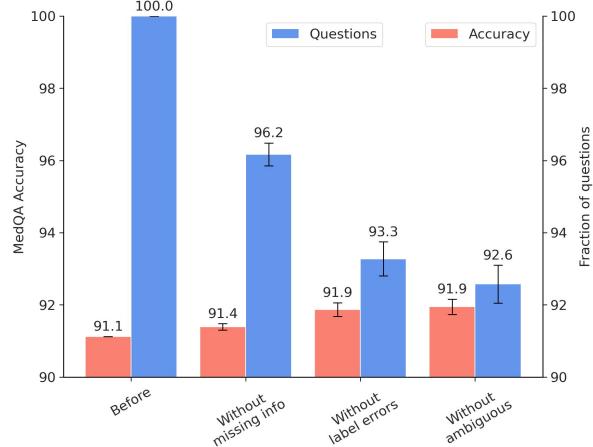
Figure 4b shows that, on average across bootstrapped committees, 3.8% of questions include missing information, following the unanimous vote of bootstrapped committees. Additionally, 2.9% likely include label errors. Another 0.7% are ambiguous. Excluding these questions is supported by high inter-rater agreement of 94%, 87.6%, and 94.6%, respectively. Importantly, Med-Gemini-L 1.0’s mistakes can be attributed disproportionately to these questions; our entropy-based uncertainty score also tends to be higher on these questions ( $t$ -test,  $p$ -value=0.033). Filtering both types improves accuracy from 91.1% to  $91.8\% \pm 0.2\%$ . Using majority instead of unanimous votes further improves accuracy to  $92.9\% \pm 0.38\%$  by discarding up to 20.9% of the uncertain questions.

#### 4.1.1. Performance on long-form medical text generation

Med-Gemini-M 1.0 demonstrates the ability to generate long-form text for three challenging real-world use cases - after-visit clinical summaries, doctor referral letter generation and medical simplification. In side-by-side comparisons, Med-Gemini-M 1.0’s responses are considered as good or better than expert responses more than half the time by clinician raters across the three tasks (Figure 5). For more task details, see Appendix C.4. Notably for the referral letter generation task, the model generated letters are preferred or tied with experts across all the samples evaluated.

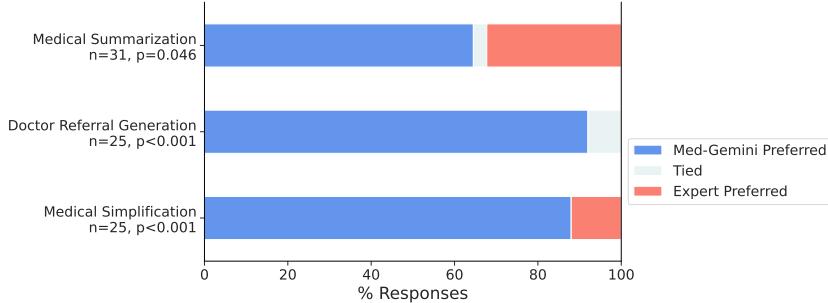


(a) MedQA Self-training &amp; search ablation



(b) MedQA label uncertainty filtering

**Figure 4 | Ablation analysis and label uncertainty on MedQA.** (a): Impact of self-training and uncertainty-guided search on Med-Gemini-L 1.0’s accuracy on MedQA. Self-training and each round of search contribute to significant performance improvements. (b): Med-Gemini-L 1.0’s accuracy (blue) and remaining questions (red) on MedQA after re-labeling by at least three US physicians per question. Filtering questions with missing information, label errors, or ambiguous groundtruth further improves accuracy. The error bars correspond to standard error across cases in (a) and standard deviation across bootstrapped annotations in (b).



**Figure 5 | Evaluation of Med-Gemini-M 1.0 on long-form text-based tasks via side-by-side comparison with experts.** The tasks considered include generation of after-visit summaries, referral letters and simplified summaries of medical systematic reviews. Evaluation was performed by clinician raters. P-values are used to denote whether the rate at which Med-Gemini-M 1.0 is preferred or tied with experts is 0.5 (two-sided t-test).

#### 4.2. Med-Gemini demonstrates multimodal understanding across diverse tasks

Our Med-Gemini models surpass, or perform competitively, with the state-of-the-art methods across seven medical multimodal benchmarks (See Table 2). We provide representative input and output examples for the multimodal tasks in Figure D1 for illustration.

In particular, Med-Gemini-L 1.0 reaches SoTA on three out-of-distribution close-ended VQA tasks—NEJM Image Challenge, multimodal USMLE sample questions (USMLE-MM), and the health & medicine subset of MMMU (MMMU-HM), outperforming GPT-4V by 8.7%, 13.1%, and 2.6%, respectively. Meanwhile, Med-Gemini-M 1.5 outperforms our previous multimodal models, Med-PaLM M (Tu et al., 2024a) on Path-VQA by 2.0% in token F1 score, and Med-Gemini-S 1.0 outperforms the previous SoTA for ECG-QA (GPT-4 with SE-WRN) by 6.1% on macro-averaged accuracy across ECG question types (Oh et al., 2023). Med-Gemini-M 1.5 also performs competitively on Slake-VQA and PAD-UFES-20 compared to the previous SoTA method (Med-PaLM M) but does not reach SoTA.

Task	Dataset	Multimodal fine-tuned	Metric	Med-Gemini	SoTA	SoTA method	Reference
Close-ended VQA	NEJM Image Challenge		Accuracy	69.7*	61.0	GPT-4V	Buckley et al. (2023)
Close-ended VQA	USMLE-MM		Accuracy	93.5*	80.4	GPT-4V	Reproduced
Close/open-ended VQA	MMMU-HM		Accuracy	67.3*	64.7	GPT-4V	Yue et al. (2023)
Close-ended Signal QA	ECG-QA	✓	Accuracy	57.7‡	51.6	GPT-4 with SE-WRN	Oh et al. (2023)
Open/Closed-ended VQA	Slake-VQA	✓	Token F1	87.5†	89.3	Med-PaLM M	Tu et al. (2024a)
Open/Closed-ended VQA	Path-VQA	✓	Token F1	64.7†	62.7	Med-PaLM M	Tu et al. (2024a)
Classification	PAD-UFES-20 6-class	✓	Accuracy	85.9†	88.0	Med-PaLM M	Tu et al. (2024a)
Classification	PAD-UFES-20 6-class	✓	Accuracy	78.8†	N/A	N/A	New Split

**Table 2 | Multimodal evaluation.** Performance comparison of Med-Gemini versus state-of-the-art (SoTA) methods. \* denotes the performance of Med-Gemini-L 1.0, † denotes the performance of Med-Gemini-M 1.5, and ‡ denotes the performance of Med-Gemini-S 1.0.

Note that we have evaluated PAD-UFES-20 on two different data split setups. We first evaluate on the Med-PaLM M split (the image-level split) for a direct, fair comparison against the previous SoTA method. In addition, we also report our model’s performance on a new split, which is a split at the patient level (Table 2).

For USMLE-MM, our model achieves accuracies of 89.5%, 92.9%, 100.0% for USMLE step 1 questions ( $n=19$ ), step 2 ( $n=14$ ), and step 3 ( $n=13$ ), respectively.

In aggregate across these seven benchmarks, Med-Gemini improve over GPT-4 by an average relative margin of 44.5%. Note that for the USMLE-MM, PADS-UFES-20 and Slake-VQA datasets, we report reproduced GPT-4V results using public APIs and the same prompt used for the corresponding Med-Gemini model.

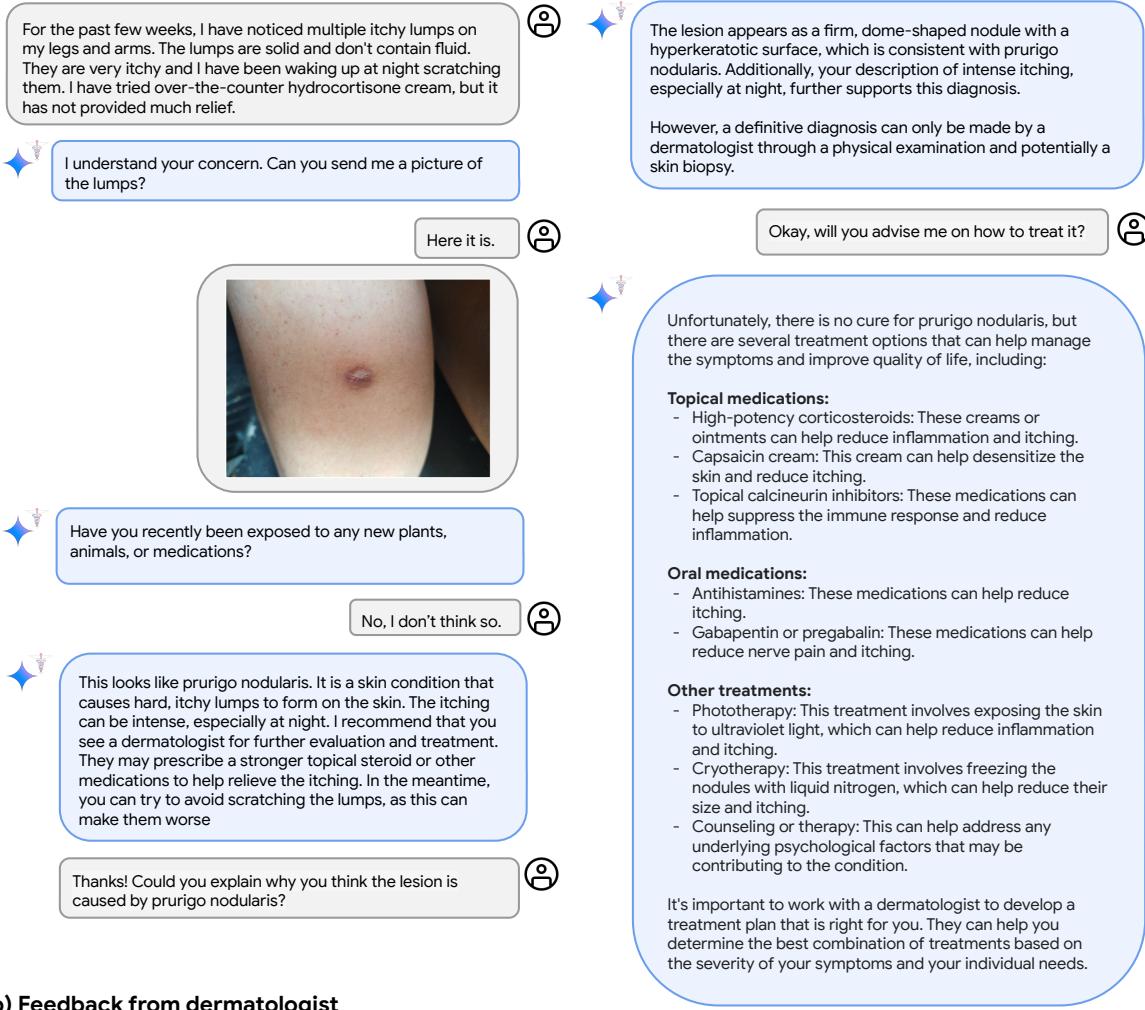
#### 4.2.1. Preview of multimodal dialogue capabilities

To extend beyond multimodal benchmarks, we demonstrate the potential for future real-world utility of Med-Gemini through hypothetical multimodal medical dialogues across two specialities.

Figure 6 illustrates an out-of-distribution setting where the dermatology image comes from a dataset (Ward et al., 2024) not used in the multimodal fine-tuning mixture. The user first asks Med-Gemini-M 1.5 about itchy lumps on their legs and arms; our model then asks the user to share an image of the lumps; after the user provides the image of their suspicious lesion, the model asks a follow-up question and continues to provide a correct diagnosis of prurigo nodularis, and recommends next steps and potential treatment options.

In Figure 7, we show a radiology dialogue example with the example image coming from the test set of the MIMIC-CXR dataset. Med-Gemini-M 1.5 demonstrates the ability to interact with a primary care provider to analyze a CXR, identify degenerative disk disease, discuss the differences between causation and correlation with a patient history of back-pain, suggest follow-up investigations to establish the cause of back pain and provide a report using non-technical language to facilitate patient understanding and communication. We observe some variability in Med-Gemini-M 1.5’s responses depending on the prompt (for instance, for some prompts the report will not list mild degenerative changes especially if prompted to focus on some other anatomical feature). A full quantification of Med-Gemini-M 1.5’s multimodal dialogue capability and variability is beyond the scope of this work, but nonetheless these qualitative examples illustrate Med-Gemini-M 1.5’s ability to support conversations about medical knowledge grounded on multimodal sources, a potentially useful attribute for applications considering user-AI and clinician-AI interaction. Real-world exploration of these use-cases would require considerable further development and validation to build upon these signs of early promise.

## (a) Dialogue example



## (b) Feedback from dermatologist

**Pros**

- “Impressive diagnostic accuracy for prurigo nodularis, a relatively rare and specialty-specific condition, based on limited data of 1 photo and brief description.”
- “Complete and thorough therapeutic ladder provided.”

**Cons**

- “Additional photos of representative lesions on different extremities would strengthen the diagnosis.”
- “The differential diagnosis could include dermatofibroma, and a brief discussion could explain why prurigo nodularis is favored based on the provided history.”
- “While there is no cure for prurigo nodularis, language could emphasize the possibility of symptom improvement and management.”

**Figure 6 | Example of a hypothetical multimodal diagnostic dialogue with Med-Gemini-M 1.5 in a dermatology setting.** (a) The user interacts with the our multimodal model, Med-Gemini-M 1.5, acting as a patient based on a case from SCIN (Ward et al., 2024), an external dataset that is not included in the fine-tuning mixture. This system would not be appropriate for use in real-world diagnostic tasks without considerable further research and development. Nonetheless, this example shows indicatively the future initial promise of how the combination of the conversational capability inherited from the original Gemini model and the newly acquired multimodal medical knowledge from fine-tuning can potentially enable attributes of a rich multi-turn diagnostic dialogue. In this interaction, Med-Gemini-M 1.5 asks for a picture when it is not provided (*multimodal information acquisition*), arrives at the right diagnosis efficiently (*open-ended diagnosis*), explains the reasoning by integrating the relevant visual features and other gathered patient symptoms (*explainability*), answers questions about treatment options while deferring aptly to experts for the final decision. (b) Shows feedback collected from dermatologists to evaluate qualitative the contents of the diagnostic dialogue, specifically asking them to comment on both positive and negative aspects.

## (a) Dialogue example



While the degenerative changes seen on the x-ray could potentially contribute to your patient's back pain, it's important to note that correlation doesn't equal causation. Many people have degenerative changes in their spine without experiencing pain, and conversely, back pain can have numerous causes beyond just degenerative changes.

Therefore, while the x-ray findings are worth noting, they don't definitively explain your patient's back pain. A comprehensive evaluation, including a detailed history, physical examination, and potentially additional imaging or tests, would be necessary to determine the exact cause of their pain.

Please let me know if you have any other questions.

To help me communicate all of this to my patient, please rewrite the report without medical terminology:

**Chest X-ray Results:**

- **Heart:** Normal size.
- **Lungs:** Clear, no signs of infection or fluid buildup.
- **Bones:** Minor wear and tear in the lower spine, common with age.

**Overall:** No signs of serious problems in the chest. The wear and tear in the lower spine might be related to your back pain, but other factors could also be involved. We'll need to do more tests to figure out the exact cause of your pain.

Thank you. That is useful while I wait for the radiologists report.

You're welcome! I'm glad I could help. If you have any further questions, feel free to ask.

**Ground truth report (not provided to model)**

**FINDINGS:** The cardiomediastinal and hilar contours are normal. The lungs are well expanded and clear, without focal consolidation, pleural effusion or pneumothorax. Mild degenerative changes are seen in the thoracic spine. **IMPRESSION:** No acute cardiopulmonary pathology.

## (b) Feedback from radiologist

**Pros**

- "Realistic hypothetical demonstration of interaction with a healthcare provider in a clinical imaging context."
- "Demonstration of high-quality realistic report generation."
- "Showcases general medical knowledge and also provides adequate limitations (i.e., mentioning additional tests needed)."

**Cons:**

- "Small osteophytes are present in the image, but along the lateral margins (rather than anterior margins) of the vertebral bodies."
- "For optimal clarity, the rewritten 'Chest X-ray Results' report, which mentions the 'lower spine' and avoids medical terminology, should specify the lower thoracic spine."

**Figure 7 | Example of a hypothetical multimodal diagnostically-assistive dialogue with Med-Gemini-M 1.5 in the radiology setting.** (a) In this interaction, Med-Gemini-M 1.5 demonstrates its ability to analyze a chest X-ray (CXR) and conduct a hypothetical realistic dialogue with a primary care physician. As above, Med-Gemini-M 1.5 is not suitable for this real-world use without further research. However, this example demonstrates initial promise, wherein Med-Gemini-M 1.5 identifies mild degenerative changes along the spine and can answer questions about the reasoning which led to this finding, demonstrate general medical knowledge about degenerative disk disease and distinguish between correlation and causation in relation to a patient history of back-pain. Finally, in this example Med-Gemini-M 1.5 is able to explain its findings in layperson's terms, demonstrating its potential for facilitating patient understanding and communication in clinical settings. The ground truth report for this CXR is provided. (b) Feedback from a radiologist about the quality of this radiology dialogue.

### 4.3. Med-Gemini shows long-context processing capability on long EHR and video tasks

Finally, we evaluate the long-context capability of Med-Gemini-M 1.5 via the “needle-in-a-haystack” medical condition retrieval task from long EHRs as well as three medical video tasks (two MAVL and one CVS assessment of surgical videos).

We demonstrate the utility of Med-Gemini-M 1.5 on the correct identification of rare and subtle problem entity (condition/symptom/procedure) in long EHR notes. The average precision and recall between Med-Gemini-M 1.5 and the baseline method are shown in Table 3 (confidence intervals in Table E2). Encouragingly, we observe that Med-Gemini-M 1.5’s one-shot ability is on-par with a carefully-tuned heuristic-based annotation-aggregation baseline approach, which is highly task-dependent. The in-context learning capability of Med-Gemini-M 1.5 to process long documents or records can easily generalize to novel problem settings without the need of extensive manual engineering. We provide an illustrative example of the prompt used, along with our model’s response in Figure 8. We attempt to benchmark GPT-4 on this task but the average context token length in this dataset significantly exceeds the maximum context window supported in the public APIs.

Task	Dataset	OOD	Metric	Med-Gemini	SoTA	SoTA method	Reference
EHR Needle-in-a-Haystack	MIMIC-III	✓	Precision	0.77	<b>0.85</b>	Annotation+Aggregation	Feder et al. (2022)
			Recall	<b>0.76</b>	0.73	Annotation+Aggregation	Feder et al. (2022)
			F1	0.77	<b>0.78</b>	Annotation+Aggregation	Feder et al. (2022)
Video QA (video-only)	MedVidQA	✓	IoU@0.3	<b>60.8</b>	32.9	RaNet	Li et al. (2022)
			IoU@0.5	<b>43.2</b>	20.6	RaNet	Li et al. (2022)
			IoU@0.7	<b>31.0</b>	15.5	RaNet	Li et al. (2022)
			mIoU	<b>43.4</b>	27.5	RaNet	Li et al. (2022)
			IoU@0.3	<b>84.4</b>	80.7	MutualSL	Weng and Li (2023)
			IoU@0.5	<b>72.9</b>	61.9	MutualSL, VPTSL	Li et al. (2022); Weng and Li (2023)
			IoU@0.7	<b>54.7</b>	44.5	VPTSL	Li et al. (2022)
Video QA (video+subtitle)	MedVidQA	✓	mIoU	<b>65.8</b>	58.3	MutualSL	Weng and Li (2023)
			Accuracy	55.2	<b>67.0</b>	ResNet3D	Reproduced
CVS assessment	Cholec80-CVS	✓					

**Table 3 | Long-context evaluation.** Performance comparison of Med-Gemini-M 1.5 versus the state-of-the-art (SoTA) methods. Note that 7 out of 155 questions in MedVidQA are not answerable due to YouTube video access (private, removed). We mark these tasks and benchmarks as out-of-distribution (OOD) because all the evaluation here is only with in-context learning and no fine-tuning

Med-Gemini-M 1.5 also achieves SoTA performance on two MedVidQA MAVL tasks (one using both video and subtitles and the other being video only), outperforming the non-LLM based SoTA models which require considerable be-spoke tuning. We note that 7 questions in MedVidQA are not answerable due to YouTube video access (private, removed). Our results therefore are reported based on the remaining 148 questions. Details are shown in Table 3. We provide an illustrative example of the prompt used, along with our model’s response in Figure 10. While evaluating MedVidQA, we also observe that the embedded captions can significantly aid the model’s understanding. Future research could explore how to optimize the use of multimodal video data, including images, text, and audio, for further improvements in video understanding. We attempt to benchmark GPT-4V on these tasks but once again run into context length limitations for most of the videos using the public APIs.

For the CVS assessment of the laparoscopic cholecystectomy video task, Med-Gemini-M 1.5 outperforms GPT-4V by 21%. However, we observe that the supervised baseline using a ResNet3D architecture performs better. Further investigations on prompting strategies or instruction fine-tuning may be required to improve the task performance of our models. We provide an illustrative example of the prompt used, along with our model’s response in Figure 9.

### (a) Retrieval of mentions

#### Prompt

**Task:** Determine if the patient has a history of the specified [Condition/Symptom/Procedure] by thoroughly examining the patient's [Medical Records].

**Instructions:** First make extensive effort to retrieve \*all\* mentions of the specified [Condition/Symptom/Procedure] that might be related to the [Initial Conclusion]. Include the relevant [Note ID], direct quotes of the evidence segment from the note.

[Condition/Symptom/Procedure]: hypothermia  
 [Medical Records]:



#### Response

- [Note ID]: 341505
- [Direct Quote]: "+sweats. Oral temp 93.7. Transferred to MICU"
- [Direct Quote]: "Constitutional: hypothermia"

### (b) Deciding the existence

#### Prompt

**Task:** Determine if the patient has a history of the specified [Condition/Symptom/Procedure/Observation] based on a list of evidences extracted from the patient's medical records.

**Instructions:** First provide a [Relevance Evaluation] which evaluates the relevance of all [Retrieved Mentions] to the specified [Condition/Symptom/Procedure/Observation] and categorize to [Explicit Confirmation], [Explicit Denial], [Strong Indication] and [Relevant Mention] of the specified [Condition/Symptom/Procedure/Observation].

Finally, review all the mentions and reasoning, provide a [Yes/No] answer together with reasons as the [Final Conclusion].

[Condition/Symptom/Procedure/Observation]: hypothermia  
 [Retrieved Mentions]: <output from first step>

#### Response

##### Relevance Evaluation:

- [Direct Quote]: "+sweats. Oral temp 93.7. Transferred to MICU"
- [Evaluation]: Strong Indication. An oral temperature of 93.7 is significantly below the normal body temperature of 98.6 and is considered hypothermia. The mention of sweats and transfer to the MICU further supports this.
- [Direct Quote]: "Constitutional: hypothermia"
- [Evaluation]: Explicit Confirmation. This note explicitly states hypothermia as a constitutional symptom.

**Final Conclusion:** Yes, the patient has a history of hypothermia. The combination of a low temperature reading and explicit mention of hypothermia in the patient's record provides a clear confirmation

**Figure 8 | Example of Med-Gemini-M 1.5's long-context capabilities on long EHR understanding (MIMIC-III Needle-in-a-Haystack).** Med-Gemini-M 1.5 performs a two-step process for determining whether a patient has a history of a specific condition based on their extensive EHR records. (a) Step 1 (Retrieval): Med-Gemini-M 1.5 identifies all mentions of "hypothermia" within the EHR notes, providing direct quotes [e.g., "+sweats. Oral temp 93.7. Transferred to medical intensive care unit (MICU)"] and note IDs for each mention. (b) Step 2 (Deciding the existence): Med-Gemini-M 1.5 then evaluates the relevance of each retrieved mention, categorizing them as explicit confirmation, strong indication, or relevant mention of hypothermia. Based on this analysis, the model concludes that the patient does have a history of hypothermia, providing clear reasoning for its decision.

#### 4.3.1. Applications of long-context capabilities in biomedicine

In addition to quantitative results, we further preview the potentials of the long-context capabilities in medical education, facilitating clinician interaction with EHR systems and biomedical literature review and summarization.

**Procedural video in clinical practice and education** In Figure 11, we qualitatively preview Med-Gemini-M 1.5’s ability to identify surgical actions from a video in the AVOS dataset. This ability holds potential for surgical care, promising to enhance surgical training through automated assessment, optimize operating room efficiency by analyzing workflows, and potentially guide surgeons in real-time during complex procedures for improved accuracy and patient outcomes. In Figure 12, we additionally present an example of Med-Gemini-M 1.5’s long-context capabilities on surgical video dialogue where the model analyzes a video clip comprising footage from a laparoscopic cholecystectomy. Med-Gemini-M 1.5 demonstrates its ability to analyze the video and conduct a dialogue with a student that might be learning about the procedure. These promising abilities have the potential to provide useful assistive tools for clinicians, perhaps improving patient safety or enhancing the process of medical training through educational aids or automated in-procedure assistance and guidance. The model correctly informs the user that they are observing a laparoscopic cholecystectomy and refers correctly to the key structures underlying the “critical view of safety”. These classification tasks, if performed scalably with high accuracy, could enable better audit of procedures (for example for quality assurance), or even prospective efficiency gains from anticipation of operative stages. For more ambitious goals such as benefits to education, operative guidance or patient safety, significant further work would need to be performed to assess more nuanced and complex capabilities. For example, we did not test Med-Gemini’s ability to accurately segment or highlight physical structures in the video and ground the dialogue with the relevant anatomy; or retrieve and present useful educational assets like diagrammatic representations of the displayed anatomy or guides to key operative stages. For uses such as education, pedagogical dialogue objectives would also likely be of considerable importance. Further work should explore these and other exciting new capabilities in a wider range of settings for procedural video, which is increasingly common in medicine.

**Clinician dialogue with EHR** In Figure 13, we demonstrate that Med-Gemini-M 1.5 effectively parses extensive medical records, synthesizing them into clear, concise summaries of active and historical conditions. Moreover, users can initiate conversations based on this summarized data, requesting more granular details from the records. Our example shows how this might include a user making natural language inquiries about specific conditions (like pneumonia) or associated diagnostic findings (such as CXR results). By streamlining access to long-form medical data and presenting the interaction in a conversational interface, this capability has the potential to significantly reduce cognitive load for clinicians and patients alike, potentially enhancing the efficiency and understanding of complex medical information without compromising staff well-being. To deliver upon this potential in real-world use would require considerable additional evaluation and research. As just one example, it would be necessary to closely examine the incidence of clinically-significant errors in retrieval or generation from grounded content; and to proactively measure and mitigate issues in dataset and model bias (as we discuss further below).

**Biomedical research** In Figure 14, we demonstrate Med-Gemini-M 1.5’s ability to process multiple research articles concerning a specific genetic locus (FTO) and its association with obesity ([Loos and Yeo, 2022](#)). In this real-world application, Med-Gemini-M 1.5 successfully comprehends the information presented in current research (full content of 12 pre-curated research papers in portable

document format) and compiles a concise summary for the user. The FTO locus we demonstrate in this example (a region of BMI- and obesity-associated variants within the gene *FTO*) is a classic example of a mechanistically understood genome-wide association studies (GWAS) hit. In this exemplar, the mechanism is a relatively complex multistep process which took extensive research to pinpoint—it involves variants altering the binding of a transcriptional repressor within an intronic super-enhancer region of the *FTO* gene, thereby leading to overexpression of two other genes, which ultimately promotes lipid accumulation (Claussnitzer et al., 2015; Laber et al., 2021).

We evaluate Med-Gemini-M 1.5’s ability to parse a large collection of academic papers on the FTO locus and provide a succinct and accessible description of the mechanistic link between FTO and obesity, together with a list of concrete supporting experimental results. As seen in Figure 14, the model provides a concise, informative, and accurate description of how the FTO locus contributes to obesity biology and presents it in a clear and digestible manner. Improvement can be made by the model listing other well-studied variants in high linkage equilibrium with rs1421085, and by providing references of where each piece of information originated from. This example shows how Med-Gemini-M 1.5’s long-context capability has clear potential to reduce cognitive load for genomic researchers and clinicians, enhancing their access to the latest findings regarding gene-disease associations; and the potential has broad relevance in other domains of biomedical and scientific research.

Prompt

...
...

You are a helpful medical video assistant.

**Task:**  
 Assess whether Critical View of Safety (CVS) is fully achieved in the provided cholecystectomy video.  
 The Critical View of Safety (CVS) is fully achieved if the following three criteria are met:

- C1: Clear view of 2 tubular structures connected to the gallbladder.
- C2: A carefully dissected hepatocystic triangle presenting an unimpeded view of only the 2 cystic structures and the cystic plate.
- C3: The lower third of the gallbladder is dissected off the cystic plate.

**Instructions:**

1. Watch the video carefully, and answer which of the Critical View of Safety (CVS) criteria are met.
2. Output Format:
  - CVS Assessment:
    - C1 - achieved / partially achieved / not achieved (Optional explanation)
    - C2 - achieved / partially achieved / not achieved (Optional explanation)
    - C3 - achieved / partially achieved / not achieved (Optional explanation)

**Response**

CVS Assessment:

C1 - partially achieved: Two tubular structures are visible but the view is not entirely clear due to some bleeding obscuring the area.

C2 - not achieved: The hepatocystic triangle is not clearly dissected, and there is not an unimpeded view of only the two cystic structures and the cystic plate.

C3 - not achieved: The lower third of the gallbladder is not dissected off the cystic plate.

**Figure 9 | Example of Med-Gemini-M 1.5’s long-context capabilities on a surgical video.** Med-Gemini-M 1.5 analyzes a video clip from the Cholec80 dataset to assess achievement of the Critical View of Safety (CVS) during a laparoscopic cholecystectomy (a keyhole operation to remove the gallbladder). The model assesses whether the three criteria defining the achievement of the CVS are met, with a detailed explanation for each criterion separately.

**Prompt**



You are a helpful medical video assistant.

**Task:** You are given a video, and a corresponding subtitle with start time and duration, followed by a question. Your task is to extract the precise video timestamps that answer the given question below.

**Instructions:** Provide one single timestamp that spans the entire length of the answer while considering the entire video. It is better to be exhaustive and providing the longest time span for the answer.

**Question:** How to relieve calf strain with foam roller massage?

**Response**

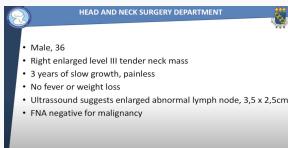
Start = 02:22

End = 02:58

Additional details: Bob demonstrates how to foam roll your calf on a foam roller to relieve calf strain.

**Figure 10 | Example of Med-Gemini-M 1.5’s long-context capabilities on medical instructional videos.** Med-Gemini-M 1.5 analyzes a video from the Medical Video Question Answering (MedVidQA) dataset to answer a specific question about relieving calf strain. The model identifies the relevant video segment (02:22-02:58) where the physical therapist explains and demonstrates the exercise for this condition. The MedVidQA ground truth time span annotation is 02:22-03:00.

**Prompt**



You are a helpful medical video assistant.

**Task:**

Annotate the surgical actions in the provided video, segmenting it into 1-second intervals. Assign each segment to one of the following classes:

- **Cutting:** If the surgeon is using a tool like scissors, scalpel, knife, or electrocautery device to cut or dissect tissues.
- **Tying:** If the surgeon is using their hands or needle holders to create secure knots.
- **Suturing:** If the surgeon is closing an open wound with a needle but not creating secure knots.
- **Background:** If no surgical action is being performed, including actions like using forceps, clamps, retractors, or dilators.

**Instructions:**

1. Watch the video carefully, segmenting it into 1-second intervals.
2. For each segment, identify the primary surgical action being performed and assign the appropriate class.
3. If no surgical action is performed in a segment, label it as "background" and provide a brief description of what is happening (e.g., "adjusting retractors").
4. If the video is too long to annotate every second, annotate each segment with one single surgical action class that best represents the dominant action within that segment.
5. Output Format:
  - Timestamp (Start - End): Class Label
  - Example: 00:00-00:01: background

**Response**

00:00-00:08: Background

00:08-03:09: Cutting

03:09-03:15: Background

03:15-04:02: Suturing

**Figure 11 | Example of Med-Gemini-M 1.5’s long-context capabilities on surgical action video span prediction.** Med-Gemini-M 1.5 analyzes a video from the Surgical Action Recognition (AVOS) dataset to annotate surgical actions in a video. It segments the video and assigns each segment a label based on the surgical action being performed (e.g., cutting, tying, suturing) or classifies it as background if no action is taking place. The ground truth AVOS annotation of this example is: 00:00-00:11: Background, 00:12-03:05: Cutting, 03:05-03:15: Background, 03:15-04:02: Suturing.

## (a) Video dialogue example

Hello, I am a medical student analyzing this surgical procedure video for my studies. Which surgical procedure is shown in this video?



The specific surgical phase shown in this video clip is the dissection and isolation of the cystic duct and artery. This phase comes after the initial access to the abdominal cavity and visualization of the gallbladder, which can be seen retracted superiorly by forceps, and before the clipping and dividing of the cystic duct and artery. The surgeon is using a laparoscopic dissector to carefully separate the cystic duct and cystic artery from the surrounding tissue. This is a critical step as it allows the surgeon to clearly identify and isolate these structures before they are clipped and divided and can help to prevent injury to the common bile duct and hepatic artery.

Which tools and how many are being used by the surgeon?

The surgeon in this video clip is using two tools:  
**Laparoscopic Grasper:** A long, thin instrument with a handle on one end and a grasping mechanism on the other, which can be opened and closed to hold tissues, and manipulate organs and other structures.  
**Hook Electroscautery:** A long thin instrument with a handle on one end and a hook-shaped tip on the other. The hook shape allows the tool to dissect, cut and cauterize tissue. Cauterizing tissue helps prevent bleeding during surgery and seals off blood vessels and lymphatic ducts to minimize the risk of leaks.

Which anatomical structures are visible?

The following anatomical structures are visible:  
**Liver:** The large, reddish-brown organ in the upper right portion of the abdomen.  
**Gallbladder:** The small, pear-shaped organ located below the liver. The gallbladder is grasped and retracted by forceps.  
**Cystic duct:** The thin, tube-like structure that connects the gallbladder to the common bile duct.  
**Common bile duct:** The larger duct visible below the cystic duct that carries bile from the liver and gallbladder to the small intestine.  
**Cystic artery:** The small artery that supplies blood to the gallbladder, it is visible alongside the cystic duct.

Which specific surgical phase is shown this video clip?

Which steps should the surgeon take to avoid injuring those structures?

To prevent such injuries, the surgeon should consider the following:  
**Critical View of Safety (CVS):** This technique involves dissecting and clearing the Calot's triangle to clearly visualize the cystic duct and artery independently before clipping or cutting to ensure that they are not mistaken for the common bile duct or hepatic artery.  
**Intraoperative Cholangiogram (IOC):** Although not frequently performed now, this procedure involves injecting a contrast dye into the cystic duct to visualize the biliary tree on X-ray. This can help the surgeon to confirm the anatomy of the bile ducts and identify any abnormal structures, ensuring the common bile duct is not clipped.  
**Fluoroscopy** By using fluoroscopy during the surgery to accurately identify the cystic duct and common bile duct using contrast.  
**Stay within Calot's triangle** - The area bordered by the cystic duct, the common hepatic duct, and the cystic artery.  
Do you want to learn more about the specific steps?

## (b) Feedback from surgeon

**Pros**

- “Impressive understanding of which surgical procedure is being performed, which specific phase in the procedure this video clip is taken from, which tools are being used, and which anatomical structures are visible.”
- “Correctly identifies the risks in the current video clip, which is injuring the common bile duct, and correctly suggests taking the safety measures, highlighting the Critical View of Safety method first.”

**Cons**

- “The surgical procedure shown in this video is a laparoscopic cholecystectomy – the surgical removal of the gallbladder.” -- this describes “cholecystectomy”, but does not refer to the “laparoscopic” part – that the procedure is performed “keyhole”.
- “Hook Electroscautery: The hook shape allows the tool to dissect, cut and cauterize tissue” -- this section does not explain what electrocautery itself actually is.

**Figure 12 | Example of Med-Gemini-M 1.5’s long-context capabilities on surgical video dialogue.** Med-Gemini-M 1.5 analyzes a video clip from the Cholec80 dataset comprising footage from a laparoscopic cholecystectomy (a keyhole operation to remove the gallbladder). The model demonstrates its ability to analyse the video and conduct a realistic dialogue with a student learning about the procedure.

Based on the medical records provided, provide a summary over all the patient's current and past medical conditions assuming today is 04/02/2024.



...



...



...



#### Active Conditions (as of 04/02/2024):

1. **Coronary Artery Disease (CAD):** The patient had an NSTEMI (a type of heart attack) in 2012 and received a bare-metal stent (BMS) in the LAD (left anterior descending artery). He is currently on aspirin, atorvastatin, and metoprolol for management. His most recent echocardiogram in 2019 showed normal biventricular systolic function. He is due for a follow-up with cardiology.
2. **Diabetes Mellitus Type 2 (DM2):** Initially diet-controlled, the patient's HbA1c levels have been rising over the years. He started Metformin in August 2022 due to worsening glycemic control. He has missed a few appointments with his endocrinologist and is due for retinopathy screening. His last HgA1C in March 2021 was 5.9.
3. **Chronic Kidney Disease (CKD):** The patient developed AKI (acute kidney injury) during a severe pneumonia episode in January 2019 and again during a left lower extremity cellulitis in July 2022. His baseline creatinine has since increased to 1.5, indicating CKD. His diuretic (Azilsartan/Chlorthalidone) was held due to AKI and remains on hold.
4. **Hypertension (HTN):** The patient has a history of hypertension, previously managed with lisinopril, which was switched to an ARB (azilsartan) due to a cough. His blood pressure has been borderline high during office visits, but home logs show reasonable control. He is currently not taking Azilsartan/Chlorthalidone due to CKD.
5. **Hyperlipidemia (HL):** The patient has high cholesterol levels and is currently on atorvastatin 80mg. His lipid panel in March 2021 showed improved levels.

#### History Conditions:

1. **Severe Pneumonia with Sepsis:** In January 2019, the patient was hospitalized for severe pneumonia complicated by sepsis and AKI. He recovered well after treatment with antibiotics and supportive care.
2. **Left Lower Extremity Cellulitis:** In July 2022, the patient was hospitalized for left lower extremity cellulitis, which led to AKI. He was treated with antibiotics and IV fluids and recovered.
3. **Needle Stick Injury:** In 2018, the patient had a needle stick injury and received post-exposure prophylaxis for HIV. Follow-up testing was negative.
4. **Possible HPV Exposure:** The patient was told by a partner in college that she may have given him HPV.

Describe in more detail the patient's history of pneumonia.



**Onset and Symptoms:** In January 2019, the patient started experiencing malaise, a non-productive cough, and shortness of breath. He also had fevers and chills. These symptoms gradually worsened over a week, leading to his visit to the emergency department.

**Diagnosis:** A portable chest X-ray done in the ED revealed multifocal infiltrates in his lungs, which is consistent with multifocal pneumonia. This means that the infection was present in multiple areas of his lungs.

**Treatment and Hospital Course:** Initially, he was started on a fluoroquinolone antibiotic (levofloxacin) and given nebulizer treatments. However, his condition deteriorated, and he became increasingly short of breath and wheezy. He was admitted to the ICU, where he was put on high-flow nasal cannula (HFNC) for oxygen support and switched to broader spectrum antibiotics - vancomycin, zosyn, and azithromycin. He also received IV fluids and other supportive care.

**Complications:** The pneumonia led to several complications:

- **Severe Sepsis:** The infection progressed to severe sepsis, and he even briefly required pressors to maintain his blood pressure. This indicates that the infection was affecting his entire body and not just his lungs.
- **Acute Kidney Injury (AKI):** His kidney function worsened, likely due to a combination of the infection, decreased oral intake, and the effects of some of his medications.
- **Elevated Troponin:** There was a rise in his troponin levels, a marker of heart damage. This was likely due to the stress of the severe infection on his body, in combination with low blood pressure and low oxygen levels.
- **Coagulopathy and Thrombocytopenia:** He developed issues with blood clotting and low platelet count, possibly indicating disseminated intravascular coagulation (DIC), a serious complication of sepsis.
- **Anemia:** He had a sudden drop in hemoglobin, requiring a blood transfusion. The source of bleeding was not clearly identified, but he was started on a PPI to protect against GI bleeding.

**Recovery:** After several days of treatment, his condition began to improve. He was weaned off of oxygen support, his kidney function started to recover, and his blood clotting issues resolved. He was eventually discharged on a course of antibiotics to complete at home.

**Causative Organism:** Blood cultures taken during his hospital stay did not grow any bacteria, but wound cultures from a subsequent hospitalization for cellulitis grew Methicillin-sensitive Staph aureus (MSSA). It's possible that MSSA was also the cause of his pneumonia, or that he had a polymicrobial infection.

**Long-Term Implications:** The patient's pneumonia was a serious event that had a significant impact on his health. He was left with a new baseline of decreased kidney function (CKD) and required adjustments to his medication regimen.

What did the CXR show when the patient had pneumonia?



The chest X-ray taken on January 14, 2019 showed:

- **Bilateral airspace opacities:** These are areas of increased density in the lungs, often indicating inflammation or fluid accumulation, which are common findings in pneumonia.
- **Mild interlobar septal thickening:** The lines between the lobes of the lungs (the interlobar septa) were slightly thicker than normal, possibly due to inflammation or fluid.
- **Mildly decreased lung volumes:** The overall space the lungs occupied was slightly reduced, which can occur in pneumonia due to fluid or inflammation.
- **Small bilateral pleural effusions:** There was a small amount of fluid in the pleural space (the area between the lungs and the chest wall) on both sides. This can be a complication of pneumonia.

**Impression:** The radiologist concluded that the findings were consistent with **multipifocal pneumonia**, meaning the pneumonia was present in multiple areas of the lungs.

**Figure 13 | Example of Med-Gemini-M 1.5's long-context capabilities on EHR dialogue.** Med-Gemini-M 1.5 demonstrates its ability to parse extensive medical records, providing a comprehensive summary of the patient's active and historical conditions. The conversation focuses on the patient's history of severe pneumonia, illustrating how the model's long-context understanding enables a user-friendly interface for interacting with and querying complex medical data. This capability has the potential to reduce cognitive load and improve information accessibility for both clinicians and patients.

**Prompt**

**Task:** Given the provided [Research Articles], what is the mechanism linking locus FTO with disease Obesity?

**Instructions:** Provide the main variants involved, the chain of molecular events linking the variant to the phenotype, and the experimental support for your statement.

📄 [Research Articles]:

BMC Cell Biology  Brain dystrophin-glycoprotein complex: persistent expression of β-dystroglycan, impaired oligomerization of Dp71 and up-regulation of utrophin in animal models of muscular dystrophy  Culligan, K., Glover, L., Dowling, P., & Ohlendieck, K.	International Journal of Biological Sciences  Complex relationship between obesity and the fat mass and obesity locus.  Yang, Q., Xiao, T., Guo, J., & Su, Z.	Frontiers in Endocrinology  Critical enzymatic functions of FTO in obesity and cancer  Deng, X., Su, R., Stanford, S., & Chen, J.
Current Diabetes Reports  FTO and obesity: mechanisms of association.  Zhao, X., Yang, Y., Sun, B. F., Zhao, Y. L., & Yang, Y. G.	The American Journal of Clinical Nutrition  FTO genotype impacts food intake and corticolimbic activation  Melhorn, S. J., Askren, M. K., Chung, W. K., Kratz, M., Bosch, T. A., Tyagi, V., ... & Schur, E. A.	New England Journal of Medicine  FTO obesity variant circuitry and adipocyte browning in humans. New England Journal of Medicine  Clausnitzer, M., Dankel, S. N., Kim, K. H., Quon, G., Meuleman, W., Haugen, C., ... & Kellis, M.
Obesity  Genetics of obesity: what we have learned over decades of research  Bouchard, C.	Science Advances  Linking the FTO obesity rs1421085 variant circuitry to cellular, metabolic, and organismal phenotypes <i>in vivo</i>  Laher, S., Forcisi, S., Bentley, L., Petzold, J., Moritz, F., Smirnov, K. S., ... & Cox, R. D.	Nature Communications  Multiple novel gene-by-environment interactions modify the effect of FTO variants on body mass index  Young, A. I., Wauthier, F., & Donnelly, P.
Cell Metabolism  Obesity and FTO: changing focus at a complex locus  Tung, Y. L., Yeo, G. S., O'Rahilly, S., & Coll, A. P.	Obesity  The fat mass-and obesity-associated (FTO) gene to obesity: lessons from mouse models  Chang, J. Y., Park, J. H., Park, S. E., Shon, J., & Park, Y. J.	Nature Reviews Genetics  The genetics of obesity: from discovery to biology.  Loos, R. J., & Yeo, G. S.

**Response****Mechanism linking the FTO locus with Obesity:**

**Variants Involved:** The primary variant is the single-nucleotide polymorphism (SNP) rs1421085, a T-to-C substitution. This SNP is located in the first intron of the FTO gene and is strongly associated with an increased risk of obesity.

**Molecular Chain of Events:**

- ARID5B Binding Disruption:** In individuals with the non-risk T allele, the transcription repressor ARID5B binds to the region containing rs1421085, effectively repressing enhancer activity.
  - The T-to-C substitution (risk allele) disrupts a critical motif necessary for ARID5B binding. This leads to a loss of repression and a gain of enhancer activity.
- Increased IRX3 and IRX5 Expression:** The increased enhancer activity results in the upregulation of two nearby genes, IRX3 and IRX5, during the early stages of adipocyte differentiation. This is particularly notable in subcutaneous (inguinal) white adipose tissue (iWAT).
  - This upregulation is specific to adipocyte progenitor cells and is not observed in mature adipocytes or other tissues like the brain.
- Adipocyte Fate Switch:** IRX3 and IRX5 promote a shift in the developmental fate of adipocytes. Instead of differentiating into energy-dissipating beige (brite) adipocytes, the progenitors are more likely to become energy-storing white adipocytes.
  - Beige adipocytes are characterized by high mitochondrial content and the expression of uncoupling protein 1 (UCP1), which allows them to burn energy and produce heat.
  - White adipocytes, on the other hand, are primarily involved in lipid storage.
- Reduced Thermogenesis and Increased Lipid Storage:** The decrease in beige adipocytes leads to reduced thermogenesis and increased lipid storage, contributing to weight gain and obesity.

**Experimental Support:**

- Mouse Models:** CRISPR-Cas9 was used to create a mouse model with a deletion of the rs1421085-conserved cis-regulatory module (CRM). Mice homozygous for this deletion exhibited reduced body weight and fat mass when fed a high-fat diet, consistent with the human phenotype. These mice also showed increased mitochondrial activity and thermogenesis marker expression in adipose tissue, as well as decreased lipid storage marker expression and smaller adipocytes.
- Human Adipocytes:** CRISPR-Cas9 editing of rs1421085 in primary human adipocytes from a patient with the risk allele restored IRX3 and IRX5 repression and activated browning expression programs and thermogenesis.
- Gene Expression and Chromatin Conformation:** Studies have shown that the risk allele is associated with increased IRX3 and IRX5 expression in human adipocyte progenitor cells. Additionally, 3D fluorescence *in situ* hybridization experiments in mice demonstrated that the risk allele disrupts the interaction between the enhancer and the Irx3 promoter, providing further evidence for the proposed mechanism.
- Metabolomic Studies:** Untargeted metabolomic analysis of adipose tissue and blood plasma from mice and humans revealed an rs1421085-specific disruption of steroid metabolic homeostasis, with a decrease in acyl steroids in risk allele carriers. This suggests that the FTO obesity variant may also exert its effects through a genetically determined metabotype.

**Figure 14 | Example of Med-Gemini-M 1.5's long-context capabilities on scientific information synthesis.** Med-Gemini-M 1.5 is given the full content of 12 research articles in portable document format (PDF) as input and is instructed to provide a succinct description of the mechanistic link between the FTO locus (a region of BMI- and obesity-associated variants within the gene *FTO*) and obesity, together with a list of concrete supporting experimental results. The model generates a concise, informative, and accurate description of how the FTO locus contributes to obesity biology and presents it in a clear and digestible manner.

## 5. Discussion

Med-Gemini, built upon the Gemini models, demonstrates significant advancements in clinical reasoning, multimodal understanding, and long-context processing within the medical domain. This is evidenced by its strong performance across a diverse range of 25 tasks spanning 14 medical benchmarks, encompassing medical knowledge, clinical reasoning, genomics, waveforms, medical imaging, health records and videos.

**MedQA performance** Notably, Med-Gemini-L 1.0 achieves a new SoTA on MedQA (USMLE), a popular benchmark for medical question answering with the use of self-training based fine-tuning and search integration. Our thorough relabeling of the MedQA test set (performed by attending clinicians) reveals important insights. While MedQA (USMLE) is a useful benchmark for assessing medical knowledge and reasoning, it is essential to acknowledge its limitations. We discover that approximately 4% of the questions contain missing information, and an additional 3% potentially have labeling errors. Establishing definitive ground truth is frequently challenging in medicine, where inter-reader variability and ambiguity are common and medical knowledge is constantly evolving. Our observations suggest that further improvements in SoTA performance on the MedQA (USMLE) benchmark in isolation may not directly correlate to progress in the capabilities of medical LLMs for meaningful real-world tasks and as such it is important to perform more comprehensive benchmarking and evaluation representative of real-world clinical workflows ([Fleming et al., 2023](#)). In general, most benchmarks have limitations around dataset size and quality. While we focus our analysis here on MedQA (USMLE), prior work has suggested similar issues with other popular benchmark datasets ([Xu et al., 2023](#)). Retraining Med-Gemini-M 1.5 with a new split of the PAD-UFES-20 dermatology dataset leads to a drop of 7.1% as compared to our results in Table 2. As such, careful attention needs to be given to the size and quality of datasets when interpreting and contextualizing model performance.

**Web search integration** Med-Gemini’s integration with web search presents exciting possibilities to provide more factually accurate and reliable answers to medical queries with LLMs. In this work, we focus on training Med-Gemini-L 1.0 to issue web search queries when uncertain and integrate the results when producing responses. While the results on MedQA, NEJM CPC, and GeneTuring benchmarks are promising, significant further research is necessary. For example, we haven’t considered restricting the search results to more authoritative medical sources ([Zakka et al., 2024](#)), using multimodal search retrieval or performed analysis on accuracy and relevance of search results and the quality of the citations ([Wu et al., 2024](#)). Further, it remains to be seen if smaller LLMs can also be taught to make use of web search. We leave these explorations to future work.

**Promising multimodal conversational capabilities** The multimodal conversational capabilities of Med-Gemini-M 1.5 are promising given they are attained without any specific medical dialogue fine-tuning. Such capabilities allow for seamless and natural interactions between people, clinicians, and AI systems. As showcased in our qualitative examples, Med-Gemini-M 1.5 has the capability to engage in multi-turn clinical dialogues, request additional information such as images when needed, explain their reasoning in a comprehensible manner, and even help provide information useful for clinical decisions while appropriately deferring the final decision to human experts. This capability has significant potential for helpful real-world applications, including assisting clinicians and patients, but of course also entails highly significant associated risks. While highlighting the potential for future research in this domain, we have not rigorously benchmarked capabilities for clinical conversation in this work as previously explored by others in dedicated research towards conversational diagnostic

AI ([Tu et al., 2024b](#)). In addition, in forthcoming work, we will also rigorously explore the capabilities of Gemini in clinically specific multimodal tasks such as radiology report generation.

**Opportunities with long-context processing** Perhaps the most notable aspect of Med-Gemini is the long-context processing capabilities because they open up new performance frontiers and novel, previously infeasible application possibilities for medical AI systems. In this work, we introduce a novel EHR task focused on identifying and verifying conditions, symptoms and procedures within very long electronic patient records. This “needle-in-a-haystack” retrieval task reflects a real-world challenge faced by clinicians ([Klerings et al., 2015](#)), and Med-Gemini-M 1.5’s performance demonstrates its potential to significantly reduce cognitive load and augment clinicians’ capabilities by efficiently extracting and analyzing crucial information from vast amounts of patient data. The medical video question answering and annotation performance suggests these capabilities can generalize to complex multimodal data. It is worth highlighting that the demonstration of long-context capabilities is in a few-shot fashion without any task-specific fine-tuning. Such capabilities open up the possibilities of fine grained analysis and annotation of genomic and multi-omic sequence data, complex imaging modalities such as pathology or volumetric images and integrative processing with health records to uncover novel insights and assist in clinical workflows.

**Importance of medical specialization and fine-tuning** Gemini models are inherently multimodal and have strong medical knowledge as a result of large-scale multimodal pretraining. This is reflected in impressive out-of-the-box performance on multimodal benchmarks such as NEJM Image Challenge surpassing similar generalist vision-language models such as GPT-4V by a large margin ([Buckley et al., 2023](#)). At the same time, medical knowledge and data (particularly multimodal data) is unique and complex and unlikely to be seen on the public internet commonly used to train LLMs. Gemini is a strong intelligence substrate but further fine-tuning, specialization and alignment of even such powerful models are necessary before use in the medical domain. At the same time, given the general capabilities of Gemini, the amount of data needed for such specialization and alignment is much lower than prior generation of medical AI systems ([Azizi et al., 2023](#)) and it is indeed possible to efficiently adapt such models even to previously unseen but important medical modalities such as ECGs with relative efficiency as demonstrated here.

**Need for rigorous evaluation beyond benchmarks** To the best of our knowledge, this work is the most comprehensive evaluation of medical LLMs and LMMs. The work includes evidence of new capabilities for medical AI and tasks that suggest real-world utility. This is particularly reinforced by strong performance of our models in evaluations of medical summarization and referral note generation. Diagnostic tasks draw considerable attention in research, but carry significant regulatory, clinical and equity-related risks that require addressing before real-world implementation is safe and feasible. The more common real-world use cases of generative AI in healthcare are therefore in non-diagnostic tasks, where errors have a lower risk-profile yet model outputs can significantly improve the efficiency of care providers by alleviating administrative burdens and assisting complex information retrieval or synthesis required in day-to-day work. At the same time, even for such non-diagnostic tasks, assurance of real-world impact requires evaluation grounded in specific use-cases and environments. These evaluations lie beyond the scope of initial benchmarking, and our results should be interpreted with appropriate caution. To assess downstream consequence and generalization of the promise we demonstrate here to real-world clinical workflows, practitioners should adhere to best practices of responsible AI, rigorously measuring multiple endpoints including equity ([Pfohl et al., 2024](#)), fairness and safety in the intended environment while also considering the multiple socio-technical factors that are use-case specific determinants of impact. Finally, it is worth noting

that while we have considered 14 diverse and challenging benchmarks in this study, over 350 medical benchmarks are available in the community ([Meta, 2024](#)).

**Responsible AI** Our work has been primarily focused on capabilities and improvements and the art of the possible with Gemini models. An important focal area for future exploration is the integration of the responsible AI principles throughout the model development process ([Pfohl et al., 2024](#)), including, but not limited to, the principles of fairness, privacy, equity, transparency and accountability. Privacy considerations in particular need to be rooted in existing healthcare policies and regulations governing and safeguarding patient information. Fairness is another area that may require attention, as there is a risk that AI systems in healthcare may unintentionally reflect or amplify historical biases and inequities ([Abràmoff et al., 2023](#); [Char et al., 2018](#); [Cirillo et al., 2020](#); [Gichoya et al., 2022](#); [Obermeyer et al., 2019](#); [Pfohl et al., 2024](#)), potentially leading to disparate model performance and harmful outcomes for marginalised groups. Such health disparities have been identified across gender ([Kent et al., 2012](#)), race ([Obermeyer et al., 2019](#); [Williams and Wyatt, 2015](#)), ethnicity ([Razai et al., 2021](#)), socioeconomic status ([Steptoe and Zaninotto, 2020](#)), sexual orientation ([Medina-Martínez et al., 2021](#)), age ([Jackson et al., 2019](#)), and other sensitive and/or protected personal characteristics. There is an increasing need for a deep intersectional analysis of impact ([Iyer et al., 2008](#); [López and Gadsden, 2017](#)), though this remains a hard technical problem ([Cabrera et al., 2019](#); [Wang et al., 2022a](#); [Yang et al., 2020](#)), and an active area of research.

As we demonstrate new capabilities for LLMs and LMMs, new opportunities arise for potential issues at the confluence of dataset bias ([Ganapathi et al., 2022](#)), model bias ([Liu et al., 2023](#)), and the socio-technical considerations for individual use cases. In the context of the capabilities we have discussed, these issues may potentially occur in in-context learning within the long-context utilization of potentially biased examples and instructions, in search integration, the dynamics of self-training, or multimodal understanding with fine-tuning and customized data encoders. Within each of these capabilities, there could be multiple points at which such biases may need to be considered. When it comes to web search integration, biases could come up at query construction time, get reflected in the returned result set ([Novin and Meyers, 2017](#)), or be embedded within each of the linked external sources, and manifest in various other subtle ways, e.g. how the results are integrated into the generative reasoning process when producing the final answer. With multimodal models, biases may occur in each of the individual modalities separately, or only be apparent jointly, across co-dependent modalities of the data ([Mandal et al., 2023](#); [Srinivasan and Bisk, 2021](#)). A comprehensive analysis of potential issues may need to consider each of these points separately, but also holistically as they are all parts of a complex system. These systems may also need to be thoroughly evaluated not only in isolation, but also with human experts in the loop.

However, these new capabilities also present an opportunity to mitigate prior issues and dramatically improve accessibility across use-cases. For example, new long-context capabilities in medicine may enable a model’s users to solve complex problems at inference time without the need for engaging in model fine-tuning, as the data can be utilized directly within the context of the query, followed by a set of natural language instructions. Previously, users of such systems would have needed to possess engineering expertise and invest additional time and resources in fine-tuning custom models for tackling such complex tasks. Web search integration, on the other hand, may prove to be invaluable when it comes to rapidly integrating newly developed pieces of medical knowledge and external consensus on what is a highly dynamic and non-stationary medical landscape. The COVID-19 pandemic has shown just how quickly the public health understanding and recommendations may need to get updated, and it also highlighted the overall danger posed by medical misinformation ([Kouzy et al., 2020](#)). Models that can reliably consume reputable up-to-date external sources may be far less likely to lead to such misinformation. Similar new opportunities are presented by the other model

capabilities, though further study is needed to develop a robust evaluation framework to assess the associated risk of bias and unfair outputs (whether individually or jointly across complex use-cases), with such assessments sociotechnically grounded in real settings for specific clinical use-cases.

## 6. Conclusion

Large multimodal language models are ushering in a new era of possibilities for health and medicine. The capabilities demonstrated by Gemini and Med-Gemini suggest a significant leap forward in the depth and breadth of opportunities to accelerate biomedical discoveries and assist in healthcare delivery and experiences. However, it is paramount that advancements in model capabilities are accompanied by meticulous attention to the reliability and safety of these systems. By prioritizing both aspects, we can responsibly envision a future where the capabilities of AI systems are meaningful and safe accelerators of both scientific progress and care in medicine.

## 7. Acknowledgements

This project was an extensive collaboration between many teams at Google Research and Google DeepMind. We thank Taylan Cemgil, Jake Sunshine, Daniel Golden, Pete Clardy, Zoubin Ghahramani and Dr. Gary Peltz (Stanford University) for their comprehensive review and detailed feedback on the manuscript. We also thank Sami Lachgar, Lauren Winer, John Guilyard, and Maggie Shiels for contributions to the narratives and visuals. We thank Yun Liu for discussions, design, and preliminary analysis for the MedQA label uncertainty experiments. We are grateful to Noam Velan, Ira Ktena, Eric Aboussouan, Karan Singhal, Shashir Reddy, Aza Tulepbergenov, Priya Gupta, Rory Sayres, Naama Hammel, Jen McKay, Peter Clardy, Chu-ling Ko, Abhinav Das, Haiyang Yu, Chang Liu, Yuchen Liu, Erica Moreira, Jordan Grimstad, Brett Hatfield, Gordon Turner, Jackie Barr, Jim Winkens, Jackie Barr, Brian Cappy, Pinal Bavishi, Tim McConnell, Ines Mezzorag, Annisah Um'rani, Christian Wright, Divya Pandya, Daireen Garcia, Prachant Bradwell, Alyssa Pierce, Sarah-Jane Allen, Erica Harland, Jennifer Ye, Praney Mittal, Donny Cheung, Andy Crowne and Preeti Singh for their valuable technical support during our research. Finally, we are grateful to Shravya Shetty, Sushant Prakash, Susan Thomas, Michael Howell, Karen DeSalvo, and Zoubin Ghahramani for their support of this project.

## 8. Data Availability

Except for the three clinical abstraction tasks, the remaining datasets used for development, benchmarking and evaluation of the AI systems are open source or otherwise accessible publicly with permissions. We will make our re-annotation of the MedQA (USMLE) dataset publicly available.

## 9. Code Availability

We are not open-sourcing model code and weights due to the safety implications of unmonitored use of such a system in medical settings. In the interest of responsible innovation, we will be working with research partners, regulators, and providers to validate and explore safe onward uses of our medical models and expect to make them available via Google Cloud APIs in due course.

## 10. Competing Interests

This study was funded by Alphabet Inc and/or a subsidiary thereof ('Alphabet'). All authors are (or were) employees of Alphabet and may own stock as part of the standard compensation package.

## References

- M. D. Abràmoff, M. E. Tarver, N. Loyo-Berrios, S. Trujillo, D. Char, Z. Obermeyer, M. B. Eydelman, F. P. of Ophthalmic Imaging, D. Algorithmic Interpretation Working Group of the Collaborative Community for Ophthalmic Imaging Foundation, Washington, and W. H. Maisel. Considerations for addressing bias in artificial intelligence for health equity. *NPJ digital medicine*, 6(1):170, 2023.
- J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, et al. PaLM 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- F. Antaki, D. Milad, M. A. Chia, C.-É. Giguère, S. Touma, J. El-Khoury, P. A. Keane, and R. Duval. Capabilities of GPT-4 in ophthalmology: an analysis of model entropy and progress towards human-level medical question answering. *British Journal of Ophthalmology*, 2023.
- S. Azizi, L. Culp, J. Freyberg, B. Mustafa, S. Baur, S. Kornblith, T. Chen, N. Tomasev, J. Mitrović, P. Strachan, et al. Robust and data-efficient generalization of self-supervised machine learning for diagnostic imaging. *Nature Biomedical Engineering*, 7(6):756–779, 2023.
- P. Barham, A. Chowdhery, J. Dean, S. Ghemawat, S. Hand, D. Hurt, M. Isard, H. Lim, R. Pang, S. Roy, et al. Pathways: Asynchronous distributed dataflow for ML. *Proceedings of Machine Learning and Systems*, 4:430–449, 2022.
- M. Besta, N. Blach, A. Kubicek, R. Gerstenberger, M. Podstawski, L. Gianinazzi, J. Gajda, T. Lehmann, H. Niewiadomski, P. Nyczyk, and T. Hoefler. Graph of thoughts: Solving elaborate problems with large language models, 2024.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- T. Buckley, J. A. Diao, A. Rodman, and A. K. Manrai. Accuracy of a vision-language model on challenging medical cases. *arXiv preprint arXiv:2311.05591*, 2023.
- Á. A. Cabrera, W. Epperson, F. Hohman, M. Kahng, J. Morgenstern, and D. H. Chau. Fairvis: Visual analytics for discovering intersectional bias in machine learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pages 46–56. IEEE, 2019.
- T. Cai, X. Wang, T. Ma, X. Chen, and D. Zhou. Large language models as tool makers. *arXiv preprint arXiv:2305.17126*, 2023.
- D. S. Char, N. H. Shah, and D. Magnus. Implementing machine learning in health care—addressing ethical challenges. *The New England journal of medicine*, 378(11):981, 2018.
- W. Chen, J. Feng, J. Lu, and J. Zhou. Endo3d: Online workflow analysis for endoscopic surgeries based on 3d cnn and lstm. In *OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, and Skin Image Analysis: First International Workshop, OR 2.0 2018, 5th International Workshop, CARE 2018, 7th International Workshop, CLIP 2018*,

*Third International Workshop, ISIC 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16 and 20, 2018, Proceedings 5*, pages 97–107. Springer, 2018.

- X. Chen, X. Wang, S. Changpinyo, A. Piergiovanni, P. Padlewski, D. Salz, S. Goodman, A. Grycner, B. Mustafa, L. Beyer, et al. PaLI: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.
- A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. PaLM: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- D. Cirillo, S. Catuara-Solarz, C. Morey, E. Guney, L. Subirats, S. Mellino, A. Gigante, A. Valencia, M. J. Rementeria, A. S. Chadha, et al. Sex and gender differences and biases in artificial intelligence for biomedicine and healthcare. *NPJ digital medicine*, 3(1):1–11, 2020.
- M. Claussnitzer, S. N. Dankel, K.-H. Kim, G. Quon, W. Meuleman, C. Haugen, V. Glunk, I. S. Sousa, J. L. Beaudry, V. Puviindran, et al. Fto obesity variant circuitry and adipocyte browning in humans. *New England Journal of Medicine*, 373(10):895–907, 2015.

Cochrane. Standards for reporting plain language summaries (pls) for cochrane diagnostic test accuracy reviews, 2014. [https://methods.cochrane.org/sites/methods.cochrane.org.sdt/files/uploads/Draft\\_PLS\\_document.pdf](https://methods.cochrane.org/sites/methods.cochrane.org.sdt/files/uploads/Draft_PLS_document.pdf).

- E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020.
- X. Dai, I. Chalkidis, S. Darkner, and D. Elliott. Revisiting transformer-based models for long document classification. *arXiv preprint arXiv:2204.06683*, 2022.
- Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- P. Densen. Challenges and opportunities facing medical education. *Transactions of the American Clinical and Climatological Association*, 122:48, 2011.
- A. Devaraj, I. Marshall, B. Wallace, and J. J. Li. Paragraph-level simplification of medical texts. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 4972–4984. Association for Computational Linguistics, June 2021. URL <https://www.aclweb.org/anthology/2021.naacl-main.395>.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, et al. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- A. V. Eriksen, S. Möller, and J. Ryg. Use of GPT-4 to diagnose complex clinical cases, 2023.
- A. Feder, I. Laish, S. Agarwal, U. Lerner, A. Atias, C. Cheung, P. Clardy, A. Peled-Cohen, R. Fellinger, H. Liu, et al. Building a clinically-focused problem list from medical notes. In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 60–68, 2022.

- W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- S. L. Fleming, A. Lozano, W. J. Haberkorn, J. A. Jindal, E. P. Reis, R. Thapa, L. Blankemeier, J. Z. Genkins, E. Steinberg, A. Nayak, et al. Medalign: A clinician-generated dataset for instruction following with electronic medical records. *arXiv preprint arXiv:2308.14089*, 2023.
- E. Ford, J. A. Carroll, H. E. Smith, D. Scott, and J. A. Cassell. Extracting information from the text of electronic medical records to improve case detection: a systematic review. *Journal of the American Medical Informatics Association*, 23(5):1007–1015, 2016.
- I. R. Galatzer-Levy, D. McDuff, V. Natarajan, A. Karthikesalingam, and M. Malgaroli. The capability of large language models to measure psychiatric functioning. *arXiv preprint arXiv:2308.01834*, 2023.
- S. Ganapathi, J. Palmer, J. E. Alderman, M. Calvert, C. Espinoza, J. Gath, M. Ghassemi, K. Heller, F. Mckay, A. Karthikesalingam, et al. Tackling bias in ai health datasets through the standing together initiative. *Nature Medicine*, 28(11):2232–2233, 2022.
- Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, Q. Guo, M. Wang, and H. Wang. Retrieval-augmented generation for large language models: A survey, 2024.
- Gemini Team, Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Gemini Team, Google. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. 2024. URL [https://storage.googleapis.com/deepmind-media/gemini/gemini\\_v1\\_5\\_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_v1_5_report.pdf).
- J. W. Gichoya, I. Banerjee, A. R. Bhimireddy, J. L. Burns, L. A. Celi, L.-C. Chen, R. Correa, N. Dullerud, M. Ghassemi, S.-C. Huang, et al. Ai recognition of patient race in medical imaging: a modelling study. *The Lancet Digital Health*, 4(6):e406–e414, 2022.
- T. Golany, A. Aides, D. Freedman, N. Rabani, Y. Liu, E. Rivlin, G. S. Corrado, Y. Matias, W. Khoury, H. Kashtan, et al. Artificial intelligence for phase recognition in complex laparoscopic cholecystectomy. *Surgical Endoscopy*, 36(12):9215–9223, 2022.
- E. D. Goodman, K. K. Patel, Y. Zhang, W. Locke, C. J. Kennedy, R. Mehrotra, S. Ren, M. Y. Guan, M. Downing, H. W. Chen, et al. A real-time spatiotemporal ai model analyzes skill in open surgical videos. *arXiv preprint arXiv:2112.07219*, 2021.
- E. D. Goodman, K. K. Patel, Y. Zhang, W. Locke, C. J. Kennedy, R. Mehrotra, S. Ren, M. Guan, O. Zohar, M. Downing, et al. Analyzing surgical technique in diverse open surgical videos with multitask machine learning. *JAMA surgery*, 159(2):185–192, 2024.
- K. K. Grandage, D. C. Slawson, and A. F. Shaughnessy. When less is more: a practical approach to searching for evidence-based answers. *Journal of the Medical Library Association*, 90(3):298, 2002.
- L. D. Gruppen. Clinical reasoning: defining it, teaching it, assessing it, studying it. *Western Journal of Emergency Medicine*, 18(1):4, 2017.
- D. Gupta and D. Demner-Fushman. Overview of the MedVidQA 2022 shared task on medical video question-answering. In D. Demner-Fushman, K. B. Cohen, S. Ananiadou, and J. Tsujii, editors, *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 264–274, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bionlp-1.25. URL <https://aclanthology.org/2022.bionlp-1.25>.

- D. Gupta, K. Attal, and D. Demner-Fushman. A dataset for medical instructional video classification and question answering. *Scientific Data*, 10(1):158, 2023.
- S. Hao, T. Liu, Z. Wang, and Z. Hu. ToolkenGPT: Augmenting frozen language models with massive tools via tool embeddings. *Advances in neural information processing systems*, 36, 2024.
- X. He, Z. Cai, W. Wei, Y. Zhang, L. Mou, E. Xing, and P. Xie. PathVQA: 30000+ questions for medical visual question answering. *arXiv preprint arXiv:2010.12435*, 2020.
- E. Horvitz, D. Heckerman, B. N. Nathwani, and L. M. Fagan. Diagnostic strategies in the hypothesis-directed pathfinder system. pages 630–636, January 1984. URL <https://www.microsoft.com/en-us/research/publication/diagnostic-strategies-hypothesis-directed-pathfinder-system/>.
- W. Hou and Z. Ji. GeneTuring tests GPT models in genomics. *BioRxiv*, 2023.
- J. Huang and K. C.-C. Chang. Towards reasoning in large language models: A survey, 2023.
- J. Huang, L. Neill, M. Wittbrodt, D. Melnick, M. Klug, M. Thompson, J. Bailitz, T. Loftus, S. Malik, A. Phull, et al. Generative artificial intelligence for chest radiograph interpretation in the emergency department. *JAMA Network Open*, 6(10):e2336100–e2336100, 2023.
- S. L. Hyland, S. Bannur, K. Bouzid, D. C. Castro, M. Ranjit, A. Schwaighofer, F. Pérez-García, V. Salvatelli, S. Srivastav, A. Thieme, et al. MAIRA-1: A specialised large multimodal model for radiology report generation. *arXiv preprint arXiv:2311.13668*, 2023.
- J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya, et al. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019.
- A. Iyer, G. Sen, and P. Östlin. The intersections of gender and class in health status and health care. *Global public health*, 3(S1):13–24, 2008.
- S. E. Jackson, R. A. Hackett, and A. Steptoe. Associations between age discrimination and health and wellbeing: cross-sectional and prospective analysis of the english longitudinal study of ageing. *The Lancet Public Health*, 4(4):e200–e208, 2019.
- P. B. Jensen, L. J. Jensen, and S. Brunak. Mining electronic health records: towards better research applications and clinical care. *Nature Reviews Genetics*, 13(6):395–405, 2012.
- D. Jin, E. Pan, N. Oufattolle, W.-H. Weng, H. Fang, and P. Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421, 2021.
- Q. Jin, Y. Yang, Q. Chen, and Z. Lu. GeneGPT: Augmenting large language models with domain tools for improved access to biomedical information. *Bioinformatics*, 40(2):btae075, 2024.
- A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019a.

- A. E. Johnson, T. J. Pollard, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, Y. Peng, Z. Lu, R. G. Mark, S. J. Berkowitz, and S. Horng. MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019b.
- Z. Kanjee, B. Crowe, and A. Rodman. Accuracy of a generative artificial intelligence model in a complex diagnostic challenge. *Jama*, 330(1):78–80, 2023.
- J. A. Kent, V. Patel, and N. A. Varela. Gender disparities in health care. *Mount Sinai Journal of Medicine: A Journal of Translational and Personalized Medicine*, 79(5):555–559, 2012.
- I. Klerings, A. S. Weinhandl, and K. J. Thaler. Information overload in healthcare: too much of a good thing? *Zeitschrift für Evidenz, Fortbildung und Qualität im Gesundheitswesen*, 109(4-5):285–290, 2015.
- R. Kouzy, J. Abi Jaoude, A. Kraitem, M. B. El Alam, B. Karam, E. Adib, J. Zarka, C. Traboulsi, E. W. Akl, and K. Baddour. Coronavirus goes viral: quantifying the covid-19 misinformation epidemic on twitter. *Cureus*, 12(3), 2020.
- S. Laber, S. Forcisi, L. Bentley, J. Petzold, F. Moritz, K. S. Smirnov, L. Al Sadat, I. Williamson, S. Strobel, T. Agnew, et al. Linking the fto obesity rs1421085 variant circuitry to cellular, metabolic, and organismal phenotypes in vivo. *Science advances*, 7(30):eabg0108, 2021.
- T. Le Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2022.
- G. Leifman, A. Aides, T. Golany, D. Freedman, and E. Rivlin. Pixel-accurate segmentation of surgical tools based on bounding box annotations. In *2022 26th International Conference on Pattern Recognition (ICPR)*, pages 5096–5103. IEEE, 2022.
- B. Li, Y. Weng, B. Sun, and S. Li. Towards visual-prompt temporal answering grounding in medical instructional video. *arXiv preprint arXiv:2203.06667*, 2022.
- C. Li, C. Wong, S. Zhang, N. Usuyama, H. Liu, J. Yang, T. Naumann, H. Poon, and J. Gao. LLaVa-Med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36, 2024.
- Y. Li, R. M. Wehbe, F. S. Ahmad, H. Wang, and Y. Luo. A comparative study of pretrained language models for long clinical text. *Journal of the American Medical Informatics Association*, 30(2):340–347, 2023.
- B. Liu, L.-M. Zhan, L. Xu, L. Ma, Y. Yang, and X.-M. Wu. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1650–1654. IEEE, 2021.
- M. Liu, Y. Ning, S. Teixayavong, M. Mertens, J. Xu, D. S. W. Ting, L. T.-E. Cheng, J. C. L. Ong, Z. L. Teo, T. F. Tan, et al. A translational perspective towards clinical ai fairness. *NPJ Digital Medicine*, 6(1):172, 2023.
- N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.
- R. J. Loos and G. S. Yeo. The genetics of obesity: from discovery to biology. *Nature Reviews Genetics*, 23(2):120–133, 2022.

- N. López and V. L. Gadsden. Health inequities, social determinants, and intersectionality. In *Perspectives on health equity and social determinants of health*. National Academies Press (US), 2017.
- M. Y. Lu, B. Chen, D. F. Williamson, R. J. Chen, K. Ikamura, G. Gerber, I. Liang, L. P. Le, T. Ding, A. V. Parwani, et al. A foundational multimodal vision language ai assistant for human pathology. *arXiv preprint arXiv:2312.07814*, 2023.
- R. Luo, L. Sun, Y. Xia, T. Qin, S. Zhang, H. Poon, and T.-Y. Liu. BioGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409, 2022.
- A. Mandal, S. Leavy, and S. Little. Multimodal composite association score: Measuring gender bias in generative multimodal models. *arXiv preprint arXiv:2304.13855*, 2023.
- P. Mascagni, D. Alapatt, A. Garcia, N. Okamoto, A. Vardazaryan, G. Costamagna, B. Dallemande, and N. Padoy. Surgical data science for safe cholecystectomy: a protocol for segmentation of hepatocystic anatomy and assessment of the critical view of safety. *arXiv preprint arXiv:2106.10916*, 2021.
- D. McDuff, M. Schaeckermann, T. Tu, A. Palepu, A. Wang, J. Garrison, K. Singhal, Y. Sharma, S. Azizi, K. Kulkarni, et al. Towards accurate differential diagnosis with large language models. *arXiv preprint arXiv:2312.00164*, 2023.
- J. Medina-Martínez, C. Saus-Ortega, M. M. Sánchez-Lorente, E. M. Sosa-Palanca, P. García-Martínez, and M. I. Mármol-López. Health inequities in lgbt people and nursing interventions to reduce them: A systematic review. *International Journal of Environmental Research and Public Health*, 18(22):11801, 2021.
- Meta. Papers with code - medical, 2024. URL <https://paperswithcode.com/area/medical>. Accessed: 2024-04-26.
- M. Moor, O. Banerjee, Z. S. H. Abad, H. M. Krumholz, J. Leskovec, E. J. Topol, and P. Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023a.
- M. Moor, Q. Huang, S. Wu, M. Yasunaga, Y. Dalmia, J. Leskovec, C. Zakka, E. P. Reis, and P. Rajpurkar. Med-flamingo: a multimodal medical few-shot learner. In *Machine Learning for Health (ML4H)*, pages 353–367. PMLR, 2023b.
- R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, et al. WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- H. Nori, Y. T. Lee, S. Zhang, D. Carignan, R. Edgar, N. Fusi, N. King, J. Larson, Y. Li, W. Liu, et al. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023.
- A. Novin and E. Meyers. Making sense of conflicting science information: Exploring bias in the search engine result page. In *Proceedings of the 2017 conference on conference human information interaction and retrieval*, pages 175–184, 2017.
- C. I. Nwoye, D. Mutter, J. Marescaux, and N. Padoy. Weakly supervised convolutional lstm approach for tool tracking in laparoscopic videos. *International journal of computer assisted radiology and surgery*, 14:1059–1067, 2019.

- Z. Obermeyer, B. Powers, C. Vogeli, and S. Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
- J. Oh, G. Lee, S. Bae, J.-m. Kwon, and E. Choi. Ecg-qa: A comprehensive question answering dataset combined with electrocardiogram. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 66277–66288. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/d0b67349dd16b83b2cf6167fb4e2be50-Paper-Datasets\\_and\\_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/d0b67349dd16b83b2cf6167fb4e2be50-Paper-Datasets_and_Benchmarks.pdf).
- J. A. Omiye, J. C. Lester, S. Spichak, V. Rotemberg, and R. Daneshjou. Large language models propagate race-based medicine. *NPJ Digital Medicine*, 6(1):195, 2023.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- A. G. Pacheco, G. R. Lima, A. S. Salomao, B. Krohling, I. P. Biral, G. G. de Angelo, F. C. Alves Jr, J. G. Esgario, A. C. Simora, P. B. Castro, et al. PAD-UFES-20: A skin lesion dataset composed of patient data and clinical images collected from smartphones. *Data in brief*, 32:106221, 2020.
- M. Parmar, A. Naik, H. Gupta, D. Agrawal, and C. Baral. LongBoX: Evaluating transformers on long-sequence clinical tasks, 2023.
- O. Pelka, S. Koitka, J. Rückert, F. Nensa, and C. M. Friedrich. Radiology objects in context (roco): a multimodal image dataset. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis: 7th Joint International Workshop, CVII-STENT 2018 and Third International Workshop, LABELS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 3*, pages 180–189. Springer, 2018.
- S. R. Pfohl, H. Cole-Lewis, R. Sayres, D. Neal, M. Asiedu, A. Dieng, N. Tomasev, Q. M. Rashid, S. Azizi, N. Rostamzadeh, et al. A toolbox for surfacing health equity harms and biases in large language models. *arXiv preprint arXiv:2403.12025*, 2024.
- M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré. Hyena hierarchy: Towards larger convolutional language models. In *International Conference on Machine Learning*, pages 28043–28078. PMLR, 2023.
- S. Qiao, Y. Ou, N. Zhang, X. Chen, Y. Yao, S. Deng, C. Tan, F. Huang, and H. Chen. Reasoning with language model prompting: A survey, 2023.
- Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian, et al. ToolLLM: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- P. Rajpurkar, E. Chen, O. Banerjee, and E. J. Topol. AI in health and medicine. *Nature medicine*, 28(1):31–38, 2022.

- V. Ramesh, N. A. Chi, and P. Rajpurkar. Improving radiology report generation systems by removing hallucinated references to non-existent priors. In A. Parziale, M. Agrawal, S. Joshi, I. Y. Chen, S. Tang, L. Oala, and A. Subbaswamy, editors, *Proceedings of the 2nd Machine Learning for Health symposium*, volume 193 of *Proceedings of Machine Learning Research*, pages 456–473. PMLR, 28 Nov 2022.
- M. S. Razai, H. K. Kankam, A. Majeed, A. Esmail, and D. R. Williams. Mitigating ethnic disparities in covid-19 and beyond. *bmj*, 372, 2021.
- M. S. Ríos, M. A. Molina-Rodríguez, D. Londoño, C. A. Guillén, S. Sierra, F. Zapata, and L. F. Giraldo. Cholec80-cvs: An open dataset with an evaluation of strasberg’s critical view of safety for ai. *Scientific Data*, 10(1):194, 2023.
- D. E. Sanford and S. M. Strasberg. A simple effective method for generation of a permanent record of the critical view of safety during laparoscopic cholecystectomy by intraoperative “doublet” photography. *Journal of the American College of Surgeons*, 218(2):170–178, 2014.
- L. Sbaffi, J. Walton, J. Blenkinsopp, and G. Walton. Information overload in emergency medicine physicians: a multisite case study exploring the causes, impact, and solutions in four north england national health service trusts. *Journal of medical Internet research*, 22(7):e19126, 2020.
- T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, and T. Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- E. Sieferd, N. Mohanty, and R. J. Holden. After visit summary: Not an afterthought. In *Proceedings of the International Symposium on Human Factors and Ergonomics in Health Care*, volume 8, pages 85–89. SAGE Publications Sage CA: Los Angeles, CA, 2019.
- K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfahl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023a.
- K. Singhal, T. Tu, J. Gottweis, R. Sayres, E. Wulczyn, L. Hou, K. Clark, S. Pfahl, H. Cole-Lewis, D. Neal, et al. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617*, 2023b.
- T. Srinivasan and Y. Bisk. Worst of both worlds: Biases compound in pre-trained vision-and-language models. *arXiv preprint arXiv:2104.08666*, 2021.
- A. Steptoe and P. Zaninotto. Lower socioeconomic status and the acceleration of aging: An outcome-wide analysis. *Proceedings of the National Academy of Sciences*, 117(26):14911–14917, 2020.
- S. M. Strasberg and M. L. Brunt. Rationale and use of the critical view of safety in laparoscopic cholecystectomy. *Journal of the American College of Surgeons*, 211(1):132–138, 2010.
- D. Stutz, A. T. Cemgil, A. G. Roy, T. Matejovicova, M. Barsbey, P. Strachan, M. Schaekermann, J. Freyberg, R. Rikhye, B. Freeman, J. P. Matos, U. Telang, D. R. Webster, Y. Liu, G. S. Corrado, Y. Matias, P. Kohli, Y. Liu, A. Doucet, and A. Karthikesalingam. Evaluating AI systems under uncertain ground truth: a case study in dermatology, 2023.

- R. Tanno, D. Barrett, A. Sellergren, S. Ghaisas, S. Dathathri, A. See, J. Welbl, K. Singhal, S. Azizi, T. Tu, et al. Consensus, dissensus and synergy between clinicians and specialist foundation models in radiology report generation. 2024.
- The New England Journal of Medicine. Image challenge. [https://www.nejm.org/  
image-challenge](https://www.nejm.org/image-challenge), 2024.
- A. Toma, P. R. Lawler, J. Ba, R. G. Krishnan, B. B. Rubin, and B. Wang. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. *arXiv preprint arXiv:2305.12031*, 2023.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- T. Tu, S. Azizi, D. Driess, M. Schaekermann, M. Amin, P.-C. Chang, A. Carroll, C. Lau, R. Tanno, I. Ktena, et al. Towards generalist biomedical AI. *NEJM AI*, 1(3):Aloa2300138, 2024a.
- T. Tu, A. Palepu, M. Schaekermann, K. Saab, J. Freyberg, R. Tanno, A. Wang, B. Li, M. Amin, N. Tomasev, et al. Towards conversational diagnostic AI. *arXiv preprint arXiv:2401.05654*, 2024b.
- A. P. Twinanda, S. Shehata, D. Mutter, J. Marescaux, M. De Mathelin, and N. Padoy. Endonet: a deep architecture for recognition tasks on laparoscopic videos. *IEEE transactions on medical imaging*, 36(1):86–97, 2016.
- L. K. Umapathi, A. Pal, and M. Sankarasubbu. Med-halt: Medical domain hallucination test for large language models. *arXiv preprint arXiv:2307.15343*, 2023.
- N. Varshney, W. Yao, H. Zhang, J. Chen, and D. Yu. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *arXiv preprint arXiv:2307.03987*, 2023.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- P. Wagner, N. Strodthoff, R.-D. Bousseljot, D. Kreiseler, F. I. Lunze, W. Samek, and T. Schaeffter. PTB-XL, a large publicly available electrocardiography dataset. *Scientific data*, 7(1):1–15, 2020.
- Z. Wan, C. Liu, X. Wang, C. Tao, H. Shen, Z. Peng, J. Fu, R. Arcucci, H. Yao, and M. Zhang. Electrocardiogram instruction tuning for report generation, 2024.
- A. Wang, V. V. Ramaswamy, and O. Russakovsky. Towards intersectionality in machine learning: Including more identities, handling underrepresentation, and performing evaluation. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 336–349, 2022a.
- X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022b.
- Y. Wang, W. Chen, X. Han, X. Lin, H. Zhao, Y. Liu, B. Zhai, J. Yuan, Q. You, and H. Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning, 2024.
- A. Ward, J. Li, J. Wang, S. Lakshminarasimhan, A. Carrick, B. Campana, J. Hartford, T. Tiyasirichokchai, S. Virmani, R. Wong, et al. Crowdsourcing dermatology images with google search ads: Creating a real-world skin condition dataset. *arXiv preprint arXiv:2402.18545*, 2024.

- L. W. Way, L. Stewart, W. Gantert, K. Liu, C. M. Lee, K. Whang, and J. G. Hunter. Causes and prevention of laparoscopic bile duct injuries: analysis of 252 cases from a human factors and cognitive psychology perspective. *Annals of surgery*, 237(4):460–469, 2003.
- J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Y. Weng and B. Li. Visual answer localization with cross-modal mutual knowledge transfer. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- D. R. Williams and R. Wyatt. Racial bias in health care and health: challenges and opportunities. *Jama*, 314(6):555–556, 2015.
- K. Wu, E. Wu, A. Cassasola, A. Zhang, K. Wei, T. Nguyen, S. Riantawan, P. S. Riantawan, D. E. Ho, and J. Zou. How well do llms cite relevant medical references? an evaluation framework and analyses. *arXiv preprint arXiv:2402.02008*, 2024.
- S. Xu, L. Yang, C. Kelly, M. Sieniek, T. Kohlberger, M. Ma, W.-H. Weng, A. Kiraly, S. Kazemzadeh, Z. Melamed, et al. ELIXR: Towards a general purpose x-ray artificial intelligence system through alignment of large language models and radiology vision encoders. *arXiv preprint arXiv:2308.01317*, 2023.
- F. Yang, M. Cisse, and S. Koyejo. Fairness with overlapping groups; a probabilistic perspective. *Advances in neural information processing systems*, 33:4067–4078, 2020.
- S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.
- J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun, et al. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- C. Zakka, R. Shad, A. Chaurasia, A. R. Dalal, J. L. Kim, M. Moor, R. Fong, C. Phillips, K. Alexander, E. Ashley, et al. Almanac—retrieval-augmented language models for clinical medicine. *NEJM AI*, 1 (2):A1oa2300068, 2024.
- J. M. Zambrano Chaves, S.-C. Huang, Y. Xu, H. Xu, N. Usuyama, S. Zhang, F. Wang, Y. Xie, M. Khademi, Z. Yang, et al. Training small multimodal models to bridge biomedical competency gap: A case study in radiology imaging. *arXiv preprint arXiv:2403.08002*, 2024.
- E. Zelikman, J. Mu, N. D. Goodman, and Y. T. Wu. Star: Self-taught reasoner bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- T. Zhang, S. G. Patil, N. Jain, S. Shen, M. Zaharia, I. Stoica, and J. E. Gonzalez. Raft: Adapting language model to domain specific rag. *arXiv preprint arXiv:2403.10131*, 2024.
- D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le, and E. Chi. Least-to-most prompting enables complex reasoning in large language models, 2023.

## Appendix

### A. Supplementary Table for Figure 1

Capability	Task	Metric	Med-Gemini	Previous SoTA	Best GPT-4 Method	SoTA Reference
Advanced Text Reasoning	NEJM CPC	Top-10 Accuracy	72.3	59.1	50.0	McDuff et al. (2023)
	GeneTuring	Averaged accuracy	53.3	48.6	48.6	Hou and Ji (2023)
	MedQA	Accuracy	91.1	90.2	90.2	Nori et al. (2023)
Multimodal Understanding	NEJM Image	Accuracy	69.7	61.0	61.0	Buckley et al. (2023)
	USMLE-MM	Accuracy	93.5	80.4	80.4	Reproduced
	ECG-QA	Accuracy	57.7	51.6	51.6	Oh et al. (2023)
	MMMU-HM	Accuracy	67.3	64.7	64.7	Yue et al. (2023)
	Path-VQA	Token F1	64.7	62.7	36.0 (Reproduced)	Tu et al. (2024a)
	PAD-UFES-20	Accuracy	85.9	88.0	50.0 (Reproduced)	Tu et al. (2024a)
Long-context Processing	Slake-VQA	Token F1	87.5	89.3	41.0 (Reproduced)	Tu et al. (2024a)
	MedVidQA	mIoU	43.4	27.5	N/A	Li et al. (2022)
	MedVidQA w/ subtitles	mIoU	65.8	58.3	N/A	Weng and Li (2023)
	Long EHR	F1	0.77	0.78	N/A	Feder et al. (2022)
Surgery Video CVS Assessment		Accuracy	50.0	67.0	0.290	Reproduced

Table A1 | **Performance results of bar plot in Figure 1.** We display our aggregated results comparing Med-Gemini to the previous state-of-the-art (SoTA) and the best GPT-4 methods across text-based, multimodal, and long-context tasks. For benchmarks where we could not find GPT-4 (or GPT-4V) reported numbers in literature, we run evaluations on the same test sets using public APIs for a head-to-head comparison, using same few-shot prompts as the corresponding Med-Gemini model including instructions to ensure outputs are correctly formatted. Note that GPT-4 results are not available (N/A) for three long context tasks due to limitations of the context window of the public GPT-4 / GPT-4V APIs

## B. Related Works

**Overview of large language model in medicine** Large language models (LLMs) have revolutionized machine learning and artificial intelligence. Researchers have employed novel network architectures, such as transformers (Vaswani et al., 2017) and pathways (Barham et al., 2022), to train these models on massive datasets. This self-supervised training across diverse domains includes models like BERT (Devlin et al., 2018), GPT (Radford et al., 2018), T5 (Raffel et al., 2020), FLAN (Wei et al., 2021), BLOOM (Le Scao et al., 2022), Flamingo (Alayrac et al., 2022), PaLM and PaLM2 (Anil et al., 2023; Chowdhery et al., 2023), LLaMA (Touvron et al., 2023), PaLI (Chen et al., 2022), PaLM-E (Driess et al., 2023), and the recent Gemini models (Gemini Team, Google, 2023, 2024). By processing text or multimodal information, these pretrained models develop a robust understanding of language, patterns, and relationships with remarkable adaptability.

Minimal fine-tuning allows these models to adapt to diverse downstream tasks. In the medical domain, Med-PaLM (Singhal et al., 2023a) and Med-PaLM 2 (Singhal et al., 2023b) represent pioneering medical LLMs fine-tuned on EHRs, exam questions, and research literature. To achieve the goal of generalist medical AI (GMAI) Moor et al. (2023a), researchers use general LLMs with prompting strategies [e.g., GPT-4 with Medprompt (Nori et al., 2023)], or refine them with multimodal data for enhanced medical understanding [e.g., Med-PaLM-M (Tu et al., 2024a)]. These models show promise in diagnosis assistance (McDuff et al., 2023), risk prediction, drug discovery, diagnostic dialogue (Tu et al., 2024b) and assessing psychiatric functioning (Galatzer-Levy et al., 2023). Our work leverages the latest Gemini models, using either direct instruction prompting or further fine-tuning for specialized medical tasks. Below, we discuss related works across the areas of language, multimodal learning, and long-context modeling.

**Model reasoning and tool-use for language-based tasks** Reasoning is a process of logical thinking that leads to a conclusion, which can be significantly enhanced by recent advances in LLMs and large multimodal models (LMMs). These improvements stem from a combination of better models

and methods that directly imitate human reasoning. Language model based reasoning techniques have been surveyed in prior works (Huang and Chang, 2023; Qiao et al., 2023), with such surveys extended into multimodal reasoning (Wang et al., 2024). Strategies to enhance language reasoning include prompt engineering, improved processes, and enhancing reasoning with access to external elements such as tools or knowledge. Prompt engineering is exemplified by approaches such as Chain-of-Thought (CoT) prompting (Wei et al., 2022), which involves generating a series of intermediate reasoning steps, Least-to-Most prompting, which involves breaking down a problem into smaller subproblems and then sequentially solving them (Zhou et al., 2023), and other methods that explore different reasoning paths to arrive at a conclusion (Besta et al., 2024; Yao et al., 2023). Improved processes arise from methods such as model updates via self-improvement (Zelikman et al., 2022) or ensemble-based approaches (Wang et al., 2022b).

Access to external elements such as tools (Hao et al., 2024; Schick et al., 2024) or external knowledge bases through the use of retrieval augmented generation (RAG) (Gao et al., 2024; Zhang et al., 2024) has also demonstrated improvements in language model reasoning. Recently LLMs have also evolved to interact with information and web tools. For tool-use, LLMs can learn to execute external tools or application programming interface (APIs), enabling them to perform actions in the real world like searching, calendar use, or using translation service via APIs (Qin et al., 2023; Schick et al., 2024). For web search specifically, LLMs incorporate traditional search engines by understanding complex queries and providing summaries that synthesize information from multiple sources (Nakano et al., 2021; Varshney et al., 2023). Furthermore, LLMs are able to not only retrieve information but also utilize tools and create ones based on user-defined needs (Cai et al., 2023). Zakka et al. (2024) have demonstrated that search tool-use can be particularly useful in medical guideline and treatment recommendations. In this work, we integrate a strategy of self-training with search to improve Med-Gemini’s capabilities for model reasoning.

**Large multimodal models in medicine** Medical practice often requires integration of multiple modalities to deliver effective care, for example, integrating data sources from patient history, medical imaging, genetic testing and lab results. Models that can integrate such modalities may provide a more comprehensive picture of a patient’s condition. Existing approaches fall into two broad categories: specialist and generalist. Specialist models excel at specific tasks within a medical discipline. Examples include models optimized for radiology report generation (Tanno et al., 2024; Zambrano Chaves et al., 2024), pathology question answering or histopathology captioning (Lu et al., 2023), radiology-related tasks (Xu et al., 2023), and cardiology electrocardiogram captioning (Wan et al., 2024).

Conversely, “generalist medical AI” (GMAI) systems (Moor et al., 2023a), such as Med-PaLM M (Tu et al., 2024a) and LLaVA-Med (Li et al., 2024), tackle a wider range of tasks across multiple specialties, aiming for broader applicability in clinical settings. The diversity of tasks performed by systems such as Med-PaLM M performance remains noteworthy as one of the earliest examples of generalist multimodal models in medicine, capable of addressing radiology, pathology, dermatology, and genomics tasks with competitive performance or exceeding SoTA across different specialties using a strong pretrained LLM with appropriate fine-tuning strategies. In this report, we further advance the evidence that AI systems can deliver strong generalist multimodal capabilities in medicine with Med-Gemini but the primary focus is on developing a model family considering application specific trade-offs.

**Long-context capability of large language models** Prior works addressing tasks with long-context windows have been limited by the capabilities of LLMs to effectively utilize large spans of text due to the memory and computation limitation of the Transformer-based models (Liu et al., 2024; Vaswani

et al., 2017). Initial efforts used hierarchical approaches to derive representations of clinical text that could not fit into a model’s limited context window (Dai et al., 2022). Subsequent work such as Clinical-Longformer and Clinical-BigBird (Li et al., 2023) focused on extending context lengths from 512 to 4096 tokens, enabling improvements in performance in question answering, document classification and information retrieval tasks. Subsequent approaches explored the use of such models in combination with imaging encoders to tackle multimodal tasks such as medical visual question answering (Gupta and Demner-Fushman, 2022). With the advancements of hardware and efficient algorithms, researchers have developed LLMs with larger context windows toward 100K tokens (Dai et al., 2019; Poli et al., 2023). Recently, Gemini further advanced the boundary of long-context capability to one million tokens (Gemini Team, Google, 2024).

However, in the domain of medicine the majority of LLMs continue to be evaluated on relatively short texts (Parmar et al., 2023) and single images. Despite their importance to medicine and clinical practice, long-context capabilities in medicine, especially in multimodal settings, are underexplored. We address this unmet need and investigate the potential of Med-Gemini on different long-context use cases, including video and long EHR-related tasks.

## C. Additional details on advanced reasoning text-based tasks

### C.1. Text-based fine-tuning & evaluation datasets

Task type	Datasets	Sample size	Description	Reference
Multiple-choice question answering	MedQA	10177	Multiple-choice questions from MedQA	Jin et al. (2021)
Multiple-choice question answering with reasoning	MedQA-R, MedQA-RS	20354	Multiple-choice questions from MedQA with synthetically generated reasoning examples	Novel
Long-form question answering	HealthSearchQA, LiveQA, MedicationQA	260	Clinician-written long-form responses	Singhal et al. (2023a)
Summarization	MIMIC-summaries	65	Clinician-written summaries of medical notes	Tu et al. (2024b)

Table C1 | Overview of datasets used for text-based instruction fine-tuning. The dataset mixture and synthetic data are curated to improve Med-Gemini-L 1.0’s reasoning and ability to make use of web search.

Task Type	Modality	Dataset	Test sample size	Description	Reference
Close-ended QA	Text	MedQA	1273	US medical licensing exam-style, multiple-choice	Jin et al. (2021)
Open-ended QA	Text	NEJM CPC	303	Complex diagnostic challenging in NEJM	McDuff et al. (2023)
Open/Close-ended QA	Text	GeneTuring	600	Commonly seen tasks in genomics research	Hou and Ji (2023)
Long-form generation	Text	Clinical Abstraction	81	Meaningful summarization in clinical practice and research	Appendix C.4

Table C2 | Overview of the evaluation benchmarks used for text-based reasoning tasks.

### C.2. MedQA (USMLE) Relabeling

The main objectives of this rater study are to identify (a) unanswerable questions due to missing information, (b) potential label errors, and (c) potentially ambiguous questions (Stutz et al., 2023). To this end, we carefully design a two-step study as follows:

- **Step 1:** Given the MedQA (USMLE) question and all four answer options:
  - (Q1) We ask “Are any of the options appropriate to answer this question?”
  - (Q2) If yes, “Select one or more options to answer the question.” (Multi-select)
  - (Q3) We ask “Is there any additional information (such as figures, plots, lab results, or similar) referenced in the question that is missing?”
  - (Q4) If yes, we ask “Do you think having access to the missing information would change your answer?”

- **Step 2:** After the rater completes step 1, they are presented with the ground truth answer from MedQA:
  - (Q1) We ask “Having revealed the question bank’s answer key, does your answer from before change?”
  - (Q2) If yes, we repeat the first two questions from above.

A key consideration that leads to this **two-step approach** is to reveal the MedQA (USMLE) ground truth at the right time to avoid biasing the rater with the ground truth when answering questions about potentially missing information in the question (Q3 and Q4). For properly identifying label errors, however, we present the raters with the ground truth so they can decide to disagree (Q1 and Q2 in step 2). To identify potentially ambiguous questions (allowing multiple “good” or true answers), we further allow raters to select multiple options as answers<sup>1</sup>. When asking about potentially missing information, we aim to identify whether this missing information is critical to answer the question.

Rater agreement in % against vote				
		All	Med-Gemini incorrect	Med-Gemini correct
info missing	majority	95.6	95.6	95.6
info missing	unanimous	94.0	94.4	94.0
label errors	majority	89.6	80.8	90.5
label errors	unanimous	87.6	74.6	88.8
ambiguous	majority	94.9	92.9	95.0
ambiguous	unanimous	94.6	92.3	94.9

Rater agreement of answers as average overlap in %				
		All	Med-Gemini incorrect	Med-Gemini correct
before	majority	49.9	36.4	51.2
before	unanimous	49.1	36.0	50.4
after	majority	75.9	54.8	77.9
after	unanimous	74.6	53.9	76.6

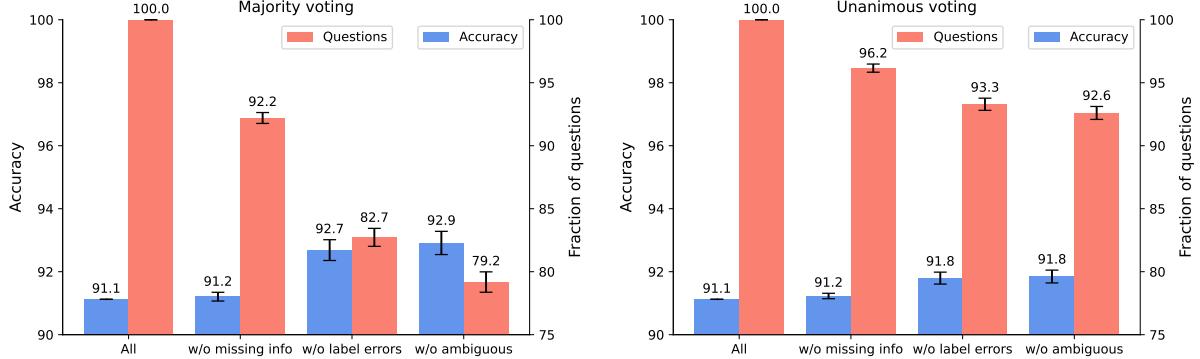
**Table C3 | Annotation agreement.** **Top:** Agreement of individual ratings against the majority or unanimous vote for various rating tasks of interest. **Bottom:** Agreement of raters’ answers in terms of average overlap before and after having revealed the MedQA ground truth.

We recruit a total of 18 primary care physicians (PCPs) from the US to participate in the study. We select PCPs located in the US because MedQA comprises USMLE-style questions across many specialties. For each MedQA (USMLE) question, we collect at least three ratings from independent raters. While the original MedQA work in (Jin et al., 2021) evaluated expert performance with access to additional text material, our raters are not instructed to use any material. However, we do not explicitly control for this. PCPs take an average of 255 seconds to complete one question; 98% take less than 10 minutes.

For each question, we aggregate the ratings in order to identify, e.g., label errors with high certainty. First, we evaluate the agreement of each rater against the majority or unanimous vote in Table C3. Specifically, we consider the rater agreement for four rating tasks of interest: whether information is missing, whether there is a label error (i.e., the rater’s answer after revealing the MedQA (USMLE) ground truth does *not* include the ground truth answer from MedQA (USMLE)), whether a question is ambiguous (i.e., the rater’s answer includes more than one option even after revealing the MedQA (USMLE) ground truth) and agreement between raw answer options selected in terms of average overlap (each rater can select none or multiple options). For the former three, agreement is generally high (>87%), although it is usually lower on questions where Med-Gemini-L

<sup>1</sup>The additional question asking whether any option is appropriate (before revealing the multi-select) is due to technical constraints.

1.0 makes mistakes. For the third, in contrast, agreement in terms of average overlap between all pairs of answers is significantly lower: typically around 75% when raters have seen the MedQA (USMLE) ground truth, but agreement drops to around 50% if the ground truth is not revealed to raters.



**Figure C1 | MedQA (USMLE) results after re-annotation.** Results complementary to Figure 3b showing accuracy (blue) and remaining MedQA (USMLE) questions (red) after filtering questions with missing information, label errors or questions deemed ambiguous when aggregating ratings using majority voting (left) or unanimous voting (right).

To measure the impact of filtering MedQA (USMLE) questions with missing information or label errors on evaluation while taking into account annotation uncertainty, we perform a bootstrapping experiment. Specifically, we repeatedly sample a committee of three raters per question (with replacement). For each committee of raters, we perform majority or unanimous voting to identify questions with missing information or label errors to be filtered for evaluation. This can be seen as an instance of the evaluation framework proposed in (Stutz et al., 2023). The advantage of bootstrapping over simple voting is that we get reliable uncertainty estimates that make sure we can identify performance changes as statistically significant. We repeat this experiment 1000 times and report average and standard deviation of accuracy and fraction of remaining questions in Figure C1.

While questions with missing information or label errors can be identified with high confidence due to high agreement, it is more difficult to judge whether a question is *ambiguous*. Here, we define a question as being ambiguous if it allows for multiple answer options to be correct. Most questions in the MedQA (USMLE) test set specifically ask for the “best”, “most likely” or “most appropriate” option. However, it is largely unclear whether answers do indeed only allow for one option to be e.g. the “best next step in management” of a case. After excluding questions with missing information and label errors using majority voting, raters selected on average 1.065 options, indicating that some questions might indeed be ambiguous. This increases to 1.119 after revealing the ground truth. To take this into account during evaluation, we define a rating as ambiguous if the rater selects more than one option after revealing the ground truth. We then follow the same analysis as above and show results in Figure C1. Overall, we find that filtering label errors has the biggest impact on Med-Gemini-L 1.0’s performance, while filtering for missing information or ambiguous questions can reduce the number of questions but does not change accuracy significantly.

### C.3. Additional results on NEJM clinico-pathological conference dataset

We present Top-1 and Top-10 performance on the NEJM clinical pathology case studies as broken down by the primary speciality of the cases (as identified by NEJM) in Table C4, for all specialties with at least 10 cases. In most specialties, Internal Medicine, Pediatrics, and Psychiatry, the best Top-1 and Top-10 performance is achieved by either Med-Gemini-L 1.0 without search or the model with

search.

Metrics	Previous SoTA AMIE		Clinician Without Search      with Search				Med-Gemini			
	Top-1↑	Top-10↑	Top-1↑	Top-10↑	Top-1↑	Top-10↑	Top-1↑	Top-10↑	Top-1↑	Top-10↑
Internal Med (159 Cases)	27.7%	61.6%	15.5%	34.6%	24.5%	47.8%	24.5%	64.8%	31.4%	74.8%
Neurology (42 Cases)	26.8%	56.1%	17.1%	31.7%	22.0%	36.6%	31.0%	66.7%	26.2%	61.9%
Pediatrics (33 Cases)	30.3%	45.5%	6.1%	22.7%	12.1%	33.3%	21.2%	45.5%	12.1%	48.5%
Psychiatry (10 Cases)	50.0%	70.0%	20.0%	50.0%	20.0%	60.0%	50.0%	100.0%	60.0%	90.0%

Table C4 | Breakdown of performance on NEJM case studies by speciality. Performance is reported for the specialities with at least 10 cases.

#### C.4. Real-world use cases for advanced reasoning on text-based tasks

We instruction fine-tune and evaluate Med-Gemini-M 1.0 on three challenging real-world tasks requiring long-form text generation. Summary results are shown in Figure 5. Detailed results of additional evaluation axes are shown in Table C5. Datasets and evaluation procedures for each task are described in more detail below.

Task	Num. Examples	Evaluation-Axis	Expert Preferred	Tied	Med-Gemini	p-value
After-Visit Summary	31	Accuracy	19%	55%	26%	<b>p &lt; 0.001</b>
		Coverage	48%	16%	35%	<b>p = 0.861</b>
		Succinctness	29%	10%	61%	<b>p = 0.017</b>
		Coherence	29%	13%	58%	<b>p = 0.017</b>
		Overall	32%	3%	65%	<b>p = 0.046</b>
Referral Letter	25	Overall	0%	8%	92%	<b>p &lt; 0.001</b>
Cochrane Simplification	25	Accuracy	52%	8%	40%	<b>p = 0.846</b>
		Coverage	12%	12%	76%	<b>p &lt; 0.001</b>
		Succinctness	4%	4%	92%	<b>p &lt; 0.001</b>
		Reading Level	0%	0%	100%	<b>p &lt; 0.001</b>
		Overall	12%	0%	88%	<b>p &lt; 0.001</b>

Table C5 | Evaluation of Med-Gemini on long-form text-based tasks via side-by-side comparison with experts. Tasks include generation of after-visit summaries, referral letters, and simplified summaries of systematic biomedical reviews. Evaluation is performed by clinician raters. P-values are for whether the rate at which Med-Gemini-M 1.0 is preferred or tied with experts is 0.5 (two-sided t-test).

**Medical summarization evaluation** This task involves generating an after-visit summary (AVS) from a de-identified history and physical (H&P) note. An H&P note is a detailed document in which a healthcare provider records the findings of a patient visit, including a patient's health background, their current symptoms, and the findings of a physical examination. It is largely written for other **healthcare providers** to ensure coordinated care. An AVS, on the other hand, is a structured report that **patients** receive at the end of a medical visit summarizing the most important aspect of the visit and their health status.

A set of 31 de-identified H&P notes is sampled from a dataset of de-identified medical notes from outpatient visits to family medicine or internal medicine providers. The expert after-visit summaries are written by U.S based clinicians following guidelines based on (Sieferd et al., 2019), and further refined by a second round of clinicians to further increase quality.

Med-Gemini-M 1.0 is prompted to generate an after-visit summary given the de-identified H&P note as follows:

Please read through the provided medical note describing an outpatient visit and extract the relevant information for each of the following 12 fields:

- Patient name/age/gender: This should summarize the patient's name, age and gender. It should use the format: "[Patient name], [age] year old [gender]". If the name is not mentioned in the note, please answer "Not available".
- Today I was seen by: This field should provide the name of the provider. If the provider seen for the note being summarized is not mentioned, please answer "Not available".
- I came in today for: This field should indicate the chief complaint or complaints that caused the visit.
- New health issues identified today are: This field should indicate any new diagnoses or other issues identified as a result of the visit being summarized. If the issue is a pre-existing condition identified in the past, please answer "No new diagnosis".
- Other health issues I have are: This field should indicate any pre-existing health issues identified in notes.
- Today we accomplished: This field should summarize the main topics of discussion and results of any procedures performed during the current visit. The summary could be a short list of procedures, or could be a text description of the patient's experience. Please be as brief as possible when providing details, such as test results or medication names. Describing the experience from the patient's point of view, using phrases like "my visit", "my condition".
- My important numbers: This field should provide the results of any measurements relevant to the visit, including vitals. Provide the results of any numeric measurements relevant to the visit, including vitals, laboratory studies, or pain scores. Please include the numbers that should be monitored. Do not fabricate numbers that are not presented in the note.
- Changes to my medications are: This field should specify any medications that were added, for which the doses were updated, or which are no longer needed after the visit. Please specify both newly added and stopped medications when possible. If no changes are apparent from the note, please answer "no changes".
- Other medications I have are: If the note indicates any existing medications for the patient that the patient should continue taking without changes, list them here. If no medications are indicated in the note, please "Not specified".
- My next steps are: This field should document the patient's next steps, including any actions they should take, test results they should expect, and follow-up visits they should schedule, along with the appropriate time frames for each.
- I should seek immediate medical attention if: If the note specifies any conditions for which the patient should immediately seek care, specify it here. Be sure to only include conditions that are mentioned in the note. If no conditions are mentioned, write "Not specified".
- Other comments from my provider: This is an optional extra field that captures any additional relevant information the provider indicated in the notes that it would be useful for the patient to know. Do not include information that is already listed in the previous field.

For each field, write at a sixth-grade reading level and avoid using abbreviations or jargon.

Note: {MEDICAL\_NOTE}  
After Visit Summary:

Physician raters are presented with the H&P note, the clinician generated AVS and our model's generated AVS. Each example is evaluated once by one of three different U.S.-based physicians across the following axes:

**Accuracy:**

Which summary is more accurate? (Are all statements in the summary correct?)

- A - B - Tie

**Coverage:**

Which summary has better coverage? (Does it include all relevant aspects of the note?)

- A - B - Tie

**Coherence:**

Which summary is easier to read? (Is the summary comprehensible to a consumer with no specific medical knowledge at a 6th-grade reading level?)

- A - B - Tie

**Succinctness:**

Which summary is more succinct? (Is the summary longer than it needs to be?)

- A - B - Tie

**Overall:**

Which summary feels higher quality to you? (Beyond these metrics, is there a gut feeling about the quality of the summary?)

- A - B - Tie

**Referral letter generation evaluation** This task involves generating a referral letter to another healthcare provider given a de-identified outpatient medical note that contains a recommendation for a referral. A medical referral letter is a formal document written by a healthcare professional that requests another healthcare professional to evaluate or treat a patient. It serves as a communication tool between healthcare providers, ensuring continuity of care and facilitating appropriate treatment for the patient.

A set of de-identified medical notes requiring inter-specialty evaluation are manually selected by clinicians from a de-identified electronic healthcare record dataset. They then generate referral letters, which are further reviewed for quality by a U.S. board-certified clinician.

Med-Gemini-M 1.0 is prompted to generate a referral letter given the medical note as follows:

You will be provided with a medical note describing a patient visit. The medical note will contain a recommendation that the patient be referred to another healthcare provider. Your task is to generate the medical referral letter for this healthcare provider.

A medical referral letter is a formal document written by a healthcare professional that requests another healthcare professional to evaluate or treat a patient. It serves as a communication tool between healthcare providers, ensuring continuity of care and facilitating appropriate treatment for the patient.

Medical Note: {MEDICAL\_NOTE}

Referral Letter:

Physician raters are presented with the outpatient note, the clinician generated referral letter and our model's generated referral letter. They are blinded to the source of each referral letter and asked to perform the following comparison:

**Instructions:**

You are given a medical note that mentions a referral to another healthcare provider. Imagine you need to write a referral letter based on the information in the note. You are provided with draft referral letters written by two different assistants. Which draft do you prefer as a starting point for editing into a final version? Please also provide a brief justification for your preference in the 'Notes' column.

**WARNING:** Unfortunately, it is not guaranteed that the draft letters accurately reflect the referral reason or patient history. This will need to be ascertained based on the provided medical note and should heavily factor into your preference.

**Options:**

- A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

Three different U.S. board certified physicians are recruited and each of them evaluates all 25 examples. Ratings are aggregated by mapping the Likert scales to a numerical range ( $[-2,2]$ ) and taking the sign of the median value.

**Medical simplification evaluation** This task involves generating a plain language summary (PLS) from a technical abstract from a biomedical systematic review. A PLS is a version of the technical abstract that is written in plain English and meant to be understood by most readers without a university education (Cochrane, 2014).

A set of 25 technical abstracts and plain language summaries from systematic reviews conducted by Cochrane is sampled from the test split of the dataset introduced by Devaraj et al. (2021). The expert plain language summaries are written by the original authors of the Cochrane systematic reviews.

Med-Gemini-M 1.0 is prompted to generate a PLS given the technical abstract as follows:

Please read through the provided technical summary of a body of medical research and provide a simplified summary that is accessible to a lay audience without medical expertise.

Technical Summary: {TECHNICAL\_ABSTRACT}  
Simplified Summary:

Clinicians are presented with the technical abstract, the original PLS and our model's generated PLS. They are blinded to the source of each PLS and asked to perform the following comparisons:

**Grounding:**

Is all the information in the simple summary factually supported by the technical summary?  
 - A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

**Coverage:**

Are the most important takeaways for a lay audience included in the simple summary?  
 - A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

**Succinctness:**

Does the simple summary only contain the most important takeaways for a lay audience?  
 - A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

**Reading Level:**

Is the reading-level of the simple summary appropriate for a lay audience?  
 - A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

**Overall:**

What is the overall quality of the simple summary for a lay audience?  
 - A Strongly Preferred - A Somewhat Preferred - Tied - B Somewhat Preferred - B Strongly Preferred

Three different U.S. board certified physicians each evaluates all 25 examples. Ratings are aggregated similar to the referral letter task.

## D. Additional details on multimodal understanding tasks

### D.1. Multimodal fine-tuning datasets

For Med-Gemini-M 1.5’s multimodal fine-tuning, we use four image-text datasets from MultiMed-Bench (Tanno et al., 2024; Tu et al., 2024a) including Slake-VQA (Liu et al., 2021), Path-VQA (He et al., 2020), MIMIC-CXR (Johnson et al., 2019a,b), PAD-UFES-20 (Pacheco et al., 2020), in addition to the Radiology Objects in COntext (ROCO) dataset (Pelka et al., 2018). We further use a subset of ECG-QA (Oh et al., 2023) to develop the health signal encoder for encoding sensor input in Med-Gemini-S 1.0. We describe the datasets in details below:

Task type	Dataset	Fine-tuning sample size	Dataset Description	Reference
VQA	Slake-VQA	9849	Close/open-ended English/Chinese VQA of radiology images	Liu et al. (2021)
VQA	Path-VQA	19755	Close/open-ended VQA of pathology images	He et al. (2020)
VQA	ROCO	29907	Close/open-ended VQA of radiology and non-radiology images	Pelka et al. (2018)
Classification	PAD-UFES-20	1838	Close-ended, multiple-choice, 6-class dermatology condition	Pacheco et al. (2020)
Classification	MIMIC-CXR	164512	Close-ended, multiple-choice, 13-class CXR condition	Johnson et al. (2019a,b)
Classification	MIMIC-CXR	237962	Close-ended, binary-choice, normal vs. abnormal classification	Johnson et al. (2019a,b)
Text Report Generation	MIMIC-CXR	90968	Open-ended, predicting CXR finding given image and indication	Johnson et al. (2019a,b)
Signal QA	ECG-QA	159306	Close-ended signal QA of electrocardiograms.	Oh et al. (2023)

Table D1 | Overview of the datasets used for multimodal instruction fine-tuning.

- **MIMIC-CXR** is a CXR dataset with free-text reports (Johnson et al., 2019a,b), consisting of 377110 chest X-ray images along with the corresponding protected health information (PHI)-removed text reports from 65379 patients (227835 image studies, with one or more image view positions). Each report is annotated with 13 common radiological conditions using the CheXpert labelling software (Irvin et al., 2019). We use the official train/test

split as described in the MIMIC-CXR for all tasks. We consider four fine-tuning tasks using MIMIC-CXR: (1) normal vs. abnormal binary classification, (2) CXR abnormality condition VQA, (3) synthetic CXR VQA, and (4) text report generation. For the normal vs. abnormal binary classification task, we classify each image into either normal or abnormal category based on the CheXpert “no finding” label using all frontal view images [anterior-posterior (AP) and posterior-anterior (PA) views] with the task prompt listed in Figure D1. For CXR abnormality condition VQA, we exclude all images with normal findings, and group positive and uncertain labels as positive class for 13 abnormal conditions: atelectasis, cardiomegaly, consolidation, edema, enlarged cardiomedastinum, fracture, lung lesion, lung opacity, pleural effusion, pleural other, pneumonia, pneumothorax, and support devices. Then we frame the abnormality detection problem into a close-ended multi-class multiple-choice question setup as shown in Figure D1. To further enrich these VQA tasks, we generate a collection of synthetic question-and-answer pairs from radiology reports by querying Gemini base models. We specifically prompt the LLM to extract pairs of yes-or-no question and the corresponding answer from each report such that they are independent of the presence of the above 13 conditions. We ensure that for each question, the number of “yes” and “no” are the same to avoid introducing spurious correlation. All VQA tasks are added as the auxiliary tasks for the report generation task which combines the image with the contextual information from the INDICATION section (reason for the study) as the model input to generate the FINDINGS and IMPRESSION sections of the report as the target, similar to prior works (Hyland et al., 2023; Tu et al., 2024a). Furthermore, following the procedure proposed in Tanno et al. (2024), we filter out the training examples whose reports reference prior studies and only keep examples where the report only refers to findings present in the input image. This aims to mitigate hallucination of references to non-existing prior reports, a common issue raised by multiple lines of research Ramesh et al. (2022) and Hyland et al. (2023). The evaluation of MIMIC-CXR will be reported in a subsequent paper.

- **PAD-UFES-20** includes 2298 clinical skin lesion images collected from various smartphone devices with different resolutions, sizes, and lighting conditions through the Dermatological and Surgical Assistance Program at the Federal University of Espírito Santo (UFES-Brazil) (Pacheco et al., 2020). Six types of skin lesions are included in the dataset: basal cell carcinoma, melanoma, squamous cell carcinoma, actinic keratosis, melanocytic nevus, and seborrheic keratosis. Each image is correlated with up to 21 clinical features (e.g., patient demographics, family cancer history, lesion location, lesion size). Given no published official splits, we adopt two PAD-UFES-20 split setup. We use Med-PaLM M split (the image-level split) for a direct, fair evaluation and comparison against the previous SoTA method. We also evaluate on a new split, which is a split at the patient level (Table 2). We set up three classification tasks for fine-tuning: (1) 6-class classification using the original label distribution and 14 clinical features (age, gender, smoke, drink, skin cancer history, cancer history, region, Fitzpatrick, horizontal and vertical diameters, itch, grew, bleed, and elevation); (2) 6-class classification using images and clinical features as the previous task, but with image augmentation on the training set using 8 RandAugment (Cubuk et al., 2020) operations: autoContrast, equalize, invert, rotate, posterize, solarize, color, and contrast; (3) 6-class classification the same as previous task, but using an upsampled subset for four minor skin conditions (melanoma, squamous cell carcinoma, seborrheic keratosis, and nevus) with image augmentation during training to mitigate the class imbalance problem. The latter two auxiliary tasks are included in the training mixture to help the model to distinguish among different types of clinical observations. We also formulate the skin condition classification problem as a close-ended multiple-choice question setup as shown in Figure D1, and report the prediction accuracy for this task.
- **Path-VQA** is a pathology VQA dataset, which consists of 998 pathology images with 32799 QA

pairs ([He et al., 2020](#)). All images are extracted from medical textbooks and online digital libraries. Each image is associated one or more questions regarding different aspects of the pathology imaging including color, location, appearance, shape, etc. 50.2% of the QA pairs are open-ended questions (divided into 7 categories: what, where, when, whose, how, and how much/how many). 49.8% of the QA pairs are close-ended questions with simple "yes/no" answer. We adopt the official splits where the training/validation/testing splits contain 19755, 6279, and 6761 QA pairs, respectively.

- **Slake-VQA** is a bilingual (English and Chinese) radiology image VQA dataset ([Liu et al., 2021](#)), containing 642 annotated images with 14028 question-answer pairs covering three imaging modalities (CT, MRI, and chest X-Rays), 39 organ systems, and 12 diseases. Questions are either open-ended or closed-ended related to various aspects of the radiology images, including plane, quality, position, organ, abnormality, size, color, shape, knowledge graph, etc. The training/validation/testing splits contain 9849, 2109, and 2070 QA pairs, respectively.
- **ROCO** (Radiology Objects in Context) dataset is a large-scale medical and multimodal imaging dataset ([Pelka et al., 2018](#)). The ROCO images are from publications available on the PubMed Central Open Access FTP mirror, which are automatically labeled as either radiology or non-radiology. Each image has its caption, keywords, the corresponding UMLS Semantic Types (SemTypes), and UMLS Concept Unique Identifiers (CUIs). We use the official training set across radiology and non-radiology, which contain 29907 image-caption pairs, and set up a captioning task for fine-tuning. We only include the images under CC BY, CC BY ND, CC BY SA and CC0 licenses in ROCO.
- **ECG-QA** is a sensor-text multimodal benchmark for assessing cardiac health ([Oh et al., 2023](#)). It is the first QA dataset specifically designed for electrocardiogram analysis based on PTB-XL ([Wagner et al., 2020](#)), containing diverse question templates, each validated by an ECG expert to ensure clinical utility. Strong performance on ECG-QA indicates the ability to grasp complex medical concepts and their connections to raw waveform signals. ECG-QA contains two types of questions involving (1) single ECG and (2) comparing two ECGs; each question type consists of (1) yes/no questions, (2) multiple-choice questions, and (3) open-ended questions to provide ECG-related attributes. We focus on single ECG questions in this work, which contain 159306, 31137, and 41093 samples for train, validation and test sets, respectively.

## D.2. Multimodal evaluation datasets

In addition to the in-distribution datasets (details in the above section), we include three out-of-distribution datasets to evaluate the multimodal capability of Gemini:

Task Type	Modality	Dataset	Test sample size	Description	Reference
VQA	Radiology	Slake-VQA	2070	English-Chinese bilingual VQA on radiology images	<a href="#">Liu et al. (2021)</a>
	Pathology	Path-VQA	6761	Close/open-ended VQA on pathology images	<a href="#">He et al. (2020)</a>
	Dermatology	PAD-UFES-20	460	6-class skin condition multiple-choice	<a href="#">Pacheco et al. (2020)</a>
	Cross-specialty	NEJM Image Challenge	934	OOD, Close-ended VQA on open domain medical images	<a href="#">The New England Journal of Medicine (2024)</a>
	Cross-specialty	USMLE-MM	46	OOD, Close-ended VQA on open domain medical images	<a href="#">Novel</a>
Signal QA	Cross-domain	MMHU-HM	150	OOD, Close/open-ended VQA on health and medical images	<a href="#">Yue et al. (2023)</a>
	Cardiology	ECG-QA	41093	Close-ended signal QA of electrocardiograms	<a href="#">Oh et al. (2023)</a>

Table D2 | Overview of the datasets used for multimodal understanding evaluation. OOD: out-of-distribution dataset.

- **New England Journal of Medicine (NEJM) Image Challenge** is a renowned clinical case challenge series that tests the diagnostic acumen and visual observation skills of medical professionals worldwide ([The New England Journal of Medicine, 2024](#)). Every week, the NEJM presents a clinical image accompanied by a brief case description. The images include radiographic images, natural and dermatoscopic skin images, electrocardiograms, histopathology

images, endoscopy images, and ophthalmoscopy images. Readers are invited to carefully analyze the photograph, consider the patient’s history, and select the final diagnosis from five possible diagnosis candidates. We collect 942 NEJM Image Challenge cases from 2005 to 2023. Each case consists of a medical image and an associated question (e.g., “What is the most likely diagnosis?”), five multiple-choice options, and a correct answer. Some cases additionally provide text captions with relevant clinical context or other background information in the question. We have collected 942 cases in total, yet 934 cases are evaluated in the end for the fair comparison (until October 12, and two cases, 20160519 and 20111103 were not evaluated due to GPT-4V filters preventing images that are assumed to be sexually explicit (Buckley et al., 2023)).

- **USMLE-MM (Multimodal)** is a multimodal multiple-choice question dataset with 46 questions identified in the sample exams provided by [www.usmle.org](http://www.usmle.org), which includes images in the question. The sample exams are used for USMLE preparation.
- **MMMU-HM (health and medicine)** is a subset of the publicly available benchmark, MMMU (Massive Multi-discipline Multimodal Understanding) validation set (Yue et al., 2023). MMMU-HM includes 150 questions related to basic medical science, clinical medicine, diagnostics and laboratory medicine, pharmacy, and public health domains.

### D.3. Additional results for multimodal tasks

**ECG-QA** To expand Med-Gemini’s capability to process raw biomedical signals for ECG-QA tasks, we augment Gemini 1.0 Nano with an ECG-specific encoder and fine-tune using two approaches: keeping Gemini model unchanged (*frozen*) and fine-tuning Gemini model (*unfrozen*). We compare our Med-Gemini-S 1.0 to their baseline counterparts: our model with *frozen* Gemini model to GPT-4 with SE-WRN ECG features in input prompts (Oh et al., 2023) and our model with *unfrozen* Gemini model to an ECG foundation model based on M<sup>3</sup>AE (Oh et al., 2023). Med-Gemini-S 1.0 with *frozen* and *unfrozen* Gemini yield accuracies of 57.7% and 58.4% on single ECG questions, respectively, outperforming GPT-4 (51.6%) by 6.1% and M<sup>3</sup>AE (57.6%) by 0.8%.

## E. Additional details on long-context understanding tasks

### E.1. Long-context evaluation datasets

- **MIMIC-III Needle-in-a-Haystack** is a specially curated dataset from MIMIC-III (Johnson et al., 2016) for subtle medical condition search-retrieval task over long EHRs. It is designed to mimic a clinically-relevant “needle-in-a-haystack” challenging problem (Gemini Team, Google, 2023). MIMIC-III is a large publicly-available medical database that contains medical records of patients admitted to intensive care units. We randomly select unstructured medical notes from 44 unique patients with more than 100 “high-value”<sup>2</sup> clinical notes. To construct “needle-in-a-haystack” examples, we use our prior work (Feder et al., 2022), which aims at identifying the problem list (conditions/symptoms/procedures) from patients’ collection of EHR documents through (1) labeling all mentions (text spans) of problems on the medical records with machine learning based annotators; (2) rule-based selection and aggregation of mentions to decide whether a problem is actually existent or not. We select the examples where there is only 1 evidence snippet identified in the aggregation step, and then randomly sample 100 negative and 100 positive examples determined by the rule-based method. 200 selected examples are then sent to 3 human medical raters to decide whether the problem actually exists or not. Specifically, the

<sup>2</sup>Consult notes, Progress notes, History and Physical notes and Discharge Summary notes authored by physicians/PAs/NPs/APRNs and Operative notes by physicians/PAs.

<p><b>Open-ended Visual QA (Path-VQA)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful medical assistant. The following are questions about medical knowledge. Solve them in a step-by-step fashion, referring to authoritative sources as needed. Question: What does the wall of the artery show with protein deposition and inflammation?</p> <p><b>Response</b> a circumferential bright pink area of necrosis</p>	<p><b>Image Classification (PAD-UFES-20 6-condition classification)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful dermatology assistant. The following are questions about skin lesions. Categorize the skin lesions into the most likely class given the patient history. Output a single option letter from the provided options as the final answer. Patient History: Age: 51, Gender: female, Smoke: false, Drink: false, Family skin cancer history: true, Family any cancer history: false, Lesion region: back, Lesion Itch: false, Lesion grew: false, Lesion bled: false, Lesion elevation: false, Fitzpatrick scale: 10, Diameters (mm): [12.0, 8.0]. Question: Which of the following is the most likely diagnosis of the patient's skin lesion? (A) Nevus (B) Basal Cell Carcinoma (C) Squamous Cell Carcinoma (D) Actinic Keratosis (E) Seborrheic Keratosis (F) Melanoma.</p> <p><b>Response</b> (A)</p>
<p><b>Open-ended Visual QA in Chinese (Slake-VQA)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful medical assistant. The following are questions about medical knowledge. Solve them in a step-by-step fashion, referring to authoritative sources as needed. Question: 图像里包含的区域属于身体哪个部分?</p> <p><b>Response</b> 腹部</p>	<p><b>Image Classification (MIMIC-CXR 13-condition classification)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful radiology assistant. The following are questions about findings in chest X-ray in the frontal view. Identify if a specific type of abnormality is shown in the X-ray. Given the -VIEW-X-ray image. Question: Which of the following abnormalities are indicated by the image? (A) Atelectasis (B) Cardiomegaly (C) Consolidation (D) Edema (E) Enlarged Cardiomediastinum (F) Fracture (G) Lung Lesion (H) Lung Opacity (I) Pleural Effusion (J) Pneumothorax (K) Pneumonia (L) Pneumothorax (M) Support Devices</p> <p><b>Response</b> (A)</p>
<p><b>Close-ended Visual QA (NEJM Image Challenge, USMLE-MM)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a medical expert. Only output the final (diagnosis, answer). Do not output the reasoning or explanation. Output the final diagnosis in the format "Final (Diagnosis, Answer); X" where X is the most (possible medical diagnosis, correct letter choice). Question: Infection? with which one of the following organisms is the most likely cause of this rash? (A) Coxsackie virus type A16 (B) Echovirus type 16 (C) Group A streptococcus (D) Herpes simplex virus type 1 (E) Norwalk virus</p> <p><b>Response</b> Final Answer: (A)</p>	<p><b>Image Classification (MIMIC-CXR normal vs abnormal classification)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful radiology assistant. The following are questions about findings in chest X-ray in the frontal view. Identify if a specific type of abnormality is shown in the X-ray. Given the -VIEW-X-ray image. Question: are there any abnormalities indicated by the image? (A) Yes (B) No.</p> <p><b>Response</b> (A)</p>
<p><b>Waveform Signal Visual QA (ECG-QA)</b></p> <p><b>Raw sensor input*</b></p> <p><b>Instruction</b> Given this ECG sequence, please answer the following question. From the provided options, select all that apply. List your selections alphabetically, separated by commas. Question: What signs of a rhythm-related disorder can be found in this ECG recording? Options: atrial fibrillation, atrial flutter, bigeminal pattern, normal functioning artificial pacemaker, sinus arrhythmia, sinus bradycardia, sinus rhythm, sinus tachycardia, supraventricular tachycardia</p> <p><b>Response</b> atrial fibrillation, atrial flutter</p>	<p><b>Text Report Generation (MIMIC-CXR)</b></p> <p><b>Visual input</b></p> <p><b>Instruction</b> You are a helpful radiology assistant. Question: Describe the findings and impression in the chest X-ray following the instructions.</p> <p><b>Response</b> <b>FINDINGS:</b> Low lung volumes are present. This accentuates the size of the cardiac silhouette which is likely mildly enlarged. Mediastinal and hilar contours are likely within normal limits. A right brachiocephalic venous stent is re-demonstrated. There is crowding of the bronchovascular structures with probable mild pulmonary vascular congestion. No pleural effusion or pneumothorax is identified. <b>IMPRESSION:</b> Low lung volumes with mild pulmonary vascular congestion.</p>

**Figure D1 | Representative examples and prompts of multimodal understanding tasks.** Med-Gemini is evaluated on various tasks, including image classification and visual question answering (VQA), demonstrating its ability to analyze and interpret diverse biomedical data formats. Note that the input for ECG-QA is the raw ECG sensor sequence, visualized here as a 12-lead ECG image from PTB-XL (Wagner et al., 2020). Also note that three MIMIC-CXR tasks are only used for instruction fine-tuning.

raters are presented with the condition name and retrieved supporting evidence snippets. The raters are then asked to answer the question: “Select ALL problems the patient HAS HAD based on the evidence in the provided note excerpts.” As a result, we have 121 positive examples and 79 negative examples based on the majority voting<sup>3</sup> (Krippendorff’s alpha at 0.77, see Table E3). The majority-vote labels are then used as ground truth labels for the subsequent evaluation. For each example, it consists of a set of medical records, a test question regarding to whether or not a condition of interest exists, and a binary ground truth label. The length of the medical records varies from 200K to 700K words.

- **Medical Instructional Video QA (MedVidQA)** is a video-language cross-modal dataset for the Medical Visual Answer Localization (MVAL) task (Gupta et al., 2023). Three medical informatics experts created 3010 health-related instructional questions for 899 videos extracted

<sup>3</sup>Example of a negative example: a patient’s records with one mention of *Sepsis* in a text segment “Received IV Ceftriaxone for concern of UTI/sepsis.”. Here the patient should not be labeled as having history of sepsis as there is no definitive diagnosis of the condition without other context.

from YouTube, and localized the visual answer to those questions by annotating their timestamps in the video, i.e., identifying the timestamp span given text question query. The mean duration time of these videos is 383.29 seconds. We follow the official data split, where 2710, 145, and 155 questions and visual answers are used for training, validation, and testing respectively. However, 7 questions are excluded due to the YouTube video access restriction (private videos, removed videos).

- **Cholec80 and Cholec80-CVS.** Cholec80 is a dataset containing 80 high-quality videos of laparoscopic cholecystectomy performed by 13 surgeons ([Twinanda et al., 2016](#)). Cholec80 is one of the most popular benchmarks for research in laparoscopic cholecystectomy video analysis with deep learning, and it has been widely used in recent research, on different video understanding tasks, including temporal segmentation of surgical phases ([Chen et al., 2018](#); [Golany et al., 2022](#)), and surgical tool detection ([Leifman et al., 2022](#); [Nwoye et al., 2019](#)). Cholec80-CVS ([Ríos et al., 2023](#)) contains Critical View of Safety (CVS) criteria annotations, provided by skilled surgeons, for each video in the Cholec80 dataset. The CVS ([Strasberg and Brunt, 2010](#)) is a mandatory method, defined by three visual criteria, used for secure identification of the cystic duct and cystic artery to minimize the risk of Bile Duct Injury (BDI). For each video in Cholec80, skilled surgeons selected different video segments where at least one CVS criteria was satisfied, and then for each selected video segment, the surgeons assigned a score of 0, 1, or 2 for each of the three CVS criteria, following an extension of the original scoring system proposed by ([Sanford and Strasberg, 2014](#)) and ([Mascagni et al., 2021](#)). In total, Cholec80-CVS provides CVS criteria annotations for 572 video segments within Cholec80 videos. We assess the performance of Med-Gemini-M 1.5 in comparison to GPT-4V and Resnet3D. It is important to note that GPT-4V does not officially support video data as input. Therefore, we sample frames from each video clip at a rate of 1 frame per second and combine a sequence of frames as input to the model. During our experimentation, we observe that GPT-4V's vision context length is limited, and we are able to insert up to 300 low-resolution images. Consequently, we filter out all video clips longer than 5 minutes. For fair comparison, we evaluate Med-Gemini-M 1.5 on the same filtered subset of video clips. To conduct the evaluation on Resnet3D, we randomly split the dataset into 5 consecutive folds and assess the performance on each validation fold separately. The average accuracy across all five folds is reported.

Task Type	Modality	Dataset	Test sample size	Description	Reference
EHR	Text	MIMIC-III-Needle-in-a-Haystack	200	Problem identification from EHR records	<a href="#">Curated based on Feder et al. (2022)</a>
MVAL	Video ± Text	MedVidQA	148	Video span prediction	<a href="#">Gupta et al. (2023)</a>
Classification	Video	Cholec80/Cholec80-CVS	572	Critical view of safety assessment	<a href="#">Ríos et al. (2023)</a>

Table E1 | Overview of the datasets used for the long-context capability evaluation. MVAL: medical visual answer localization.

Methods	Precision	Recall	F1
Heuristic-based baseline	<b>0.85 (0.78, 0.92)</b>	0.73 (0.64, 0.80)	<b>0.78 (0.72, 0.84)</b>
Med-Gemini-M 1.5 (one-shot)	0.77 (0.66, 0.86)	<b>0.76 (0.67, 0.86)</b>	0.77 (0.68, 0.84)

Table E2 | Performance comparison of Med-Gemini-M 1.5 versus the heuristic-based annotation-aggregation baseline.

## E.2. Rater agreement metrics for the long EHR understanding task

To ensure the reliability of the EHR benchmark, we collect ratings from three independent raters for each of the 200 example questions. The following metrics demonstrate strong consistency among raters:

- Jaccard Similarity Index: Measures the overlap between sets of rater selections. Let  $A$ ,  $B$ , and  $C$  represent the sets of selections made by each rater. The Jaccard similarity index for unanimous selections is defined as  $J_{=3} = \frac{|A \cap B \cap C|}{|A \cup B \cup C|}$ . The Jaccard similarity index for at least 2 raters being in agreement is defined as  $J_{\geq 2} = \frac{|(A \cap B) \cup (A \cap C) \cup (B \cap C)|}{|A \cup B \cup C|}$ .
- Krippendorff's Alpha: A reliability coefficient designed for multiple raters.

	Num. tasks	$J_{=3}$	$J_{\geq 2}$	Krippendorff's alpha
Existence of condition	200	0.83	0.915	0.77

Table E3 | Rater agreement metrics on long EHR understanding task.

A Jaccard similarity index of 0.83 for unanimous selections indicates substantial agreement when all three raters select identical choices. An even higher Jaccard index of 0.915 reflects strong consistency when at least two out of three raters make the same selections. A Krippendorff's alpha of 0.77 indicates good agreement on the existence of medical conditions within the EHR data.