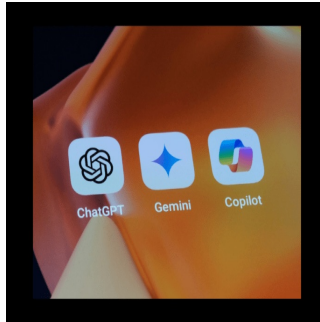

Google Med-Gemini Outperforms GPT-4



Google Research and DeepMind have introduced Med-Gemini, an advanced family of AI models tailored for medical applications. These models represent a significant leap in clinical diagnostics, offering immense potential in real-world healthcare scenarios. They address the multifaceted challenges doctors face daily, such as understanding patient records, staying updated on medical advancements, and nurturing doctor-patient relationships. Unlike previous AI models, Med-Gemini can process various types of information, including text, images, videos, and audio, enabling it to perform tasks like language processing, understanding diverse data, and long-context reasoning. Researchers have recently fine-tuned these models specifically for medical tasks, resulting in impressive capabilities showcased in [their paper posted on ArXiv](#).

Advancing Medical AI Through Self-Training and Web Search Capabilities

Med-Gemini, an AI model developed by Google, aims to assist doctors in diagnosing and treating patients by amalgamating various sources of medical information. Doctors often navigate through patient symptoms, medical history, lab results, and treatment responses to formulate a diagnosis and treatment plan. Recognising the dynamic nature of medical treatments, which constantly evolve with updates and new findings, Google incorporated web-based searching into Med-Gemini to enhance its clinical reasoning capabilities.

Med-Gemini underwent training on the MedQA dataset, which comprises multiple-choice questions resembling those found in the US Medical License Exam (USMLE). This training equips Med-Gemini with a solid foundation in medical knowledge and reasoning across diverse scenarios. Additionally, Google introduced two novel datasets to further refine Med-Gemini's performance. The first, MedQA-R (Reasoning), supplements MedQA with synthetically generated reasoning explanations known as 'Chain-of-Thoughts' (CoTs). The second, MedQA-RS (Reasoning and Search), instructs the model to leverage web search results to enhance answer accuracy, particularly when faced with uncertain queries.

In evaluating Med-Gemini's performance, researchers conducted tests on 14 medical benchmarks. Impressively, Med-Gemini achieved state-of-the-art performance on 10 of these benchmarks, surpassing the performance of the GPT-4 model family in every comparable scenario. Notably, on the MedQA (USMLE) benchmark, Med-Gemini attained an accuracy rate of 91.1% using its uncertainty-guided search strategy, outperforming Google's previous medical LLM, Med-PaLM 2, by 4.5%.

Moreover, Med-Gemini excelled on seven multimodal benchmarks, including the challenging New England Journal of Medicine (NEJM) image challenge. In these assessments, Med-Gemini outperformed GPT-4 by an average relative margin of 44.5%, showcasing its superiority in processing and analysing diverse types of medical data.

Despite these promising results, the researchers acknowledge the need for further investigation. Areas for future exploration include refining search results to prioritise authoritative medical sources, exploring multimodal search retrieval techniques, and assessing the adaptability of smaller language models to utilise web search effectively.

Retrieving Information from Electronic Health Records with Precision and Efficiency

Electronic health records (EHRs) are often lengthy and contain various complexities such as textual similarities, misspellings, acronyms, and synonyms, posing challenges for AI systems like Med-Gemini. To evaluate Med-Gemini's ability to comprehend and reason from extensive medical data, researchers conducted a 'needle-in-a-haystack' task using the MIMIC-III database, which houses de-identified health data from intensive care patients.

In this task, Med-Gemini was tasked with identifying mentions of rare medical conditions within vast EHR notes, simulating a real-world challenge

© For personal and private use only. Reproduction must be permitted by the copyright holder. Email to copyright@mindbyte.eu.

clinicians face. The dataset comprised 200 examples, each containing de-identified notes from 44 ICU patients with long medical histories. Each example met specific criteria, including having more than 100 medical notes, containing a single mention of the condition of interest, and ranging from 200,000 to 700,000 words in length.

The task involved two steps: first, Med-Gemini had to retrieve all mentions of the specified medical problem from the extensive records, and second, it had to evaluate the relevance of each mention, categorise them, and determine whether the patient had a history of the problem, providing clear reasoning for its decision.

Med-Gemini performed well on the task, achieving a precision of 0.77 and surpassing the state-of-the-art method in recall, scoring 0.76 compared to 0.73. The model's ability to process long-context information demonstrated its potential to significantly reduce cognitive load for clinicians by efficiently extracting and analyzing relevant information from vast amounts of patient data, thus augmenting their decision-making capabilities in clinical settings.

Facilitating Clinical Conversations with Precision and Caution

In a practical test, Med-Gemini demonstrated its real-world utility by engaging in a conversation with a patient about an itchy skin lump. After requesting an image, the model asked relevant follow-up questions and accurately diagnosed the rare lesion, providing appropriate recommendations. Additionally, Med-Gemini interpreted a chest X-ray for a physician while they awaited a formal radiologist's report and crafted a simplified version of the report in plain English for the patient.

The researchers highlighted the promising multimodal conversation capabilities of Med-Gemini-M 1.5, achieved without specific medical dialogue fine-tuning. These capabilities enable seamless interactions between individuals, clinicians, and AI systems. However, they acknowledge the need for further research in this area.

While acknowledging the potential for helpful real-world applications, the researchers also caution about significant risks associated with such capabilities. They note the absence of rigorous benchmarking of Med-Gemini's clinical conversation capabilities in this study, which has been explored in dedicated research on conversational diagnostic AI by others.

Ensuring Responsible AI: Prioritising Privacy and Fairness in the Development

Researchers acknowledge that there is still extensive work ahead, despite the promising initial capabilities of the Med-Gemini model. They emphasise the importance of integrating responsible AI principles, such as privacy and fairness, into the model development process.

Privacy considerations are particularly crucial, aligning with existing healthcare policies and regulations governing patient information. The researchers also recognise the potential for AI systems in healthcare to inadvertently perpetuate historical biases and inequities, emphasising the need for fairness and the prevention of disparate model performance and harmful outcomes for marginalised groups.

However, authors ultimately view Med-Gemini as a tool for positive impact. They assert that large multimodal language models are opening up new possibilities for health and medicine, offering significant potential to accelerate biomedical discoveries and improve healthcare delivery and experiences. Yet, they stress the importance of balancing advancements in model capabilities with a meticulous focus on reliability and safety. By prioritising both aspects, they envision a future where AI systems can responsibly contribute to scientific progress and healthcare with meaningful and safe acceleration.

Source: [arXiv](#)

Image Credit: [iStock](#)

Published on : Wed, 8 May 2024