

# MEDALIGN: A Clinician-Generated Dataset for Instruction Following with Electronic Medical Records

Scott L. Fleming<sup>\*,1,2</sup>, Alejandro Lozano<sup>\*,1</sup>, William J. Haberkorn<sup>\*,3,4</sup>, Jenelle A. Jindal<sup>\*,5</sup>, Eduardo Reis<sup>\*,6,7,8</sup>, Rahul Thapa<sup>9</sup>, Louis Blankemeier<sup>10</sup>, Julian Z. Genkins<sup>11,12</sup>, Ethan Steinberg<sup>2</sup>, Ashwin Nayak<sup>13</sup>, Birju Patel<sup>5</sup>, Chia-Chun Chiang<sup>14,15</sup>, Alison Callahan<sup>5,13</sup>, Zepeng Huo<sup>5</sup>, Sergios Gatidis<sup>6</sup>, Scott Adams<sup>6</sup>, Oluseyi Fayanju<sup>13</sup>, Shreya J. Shah<sup>13</sup>, Thomas Savage<sup>1,16</sup>, Ethan Goh<sup>5,17</sup>, Akshay S. Chaudhari<sup>1,6,15</sup>, Nima Aghaeepour<sup>1,3,4</sup>, Christopher Sharp<sup>13,15</sup>, Michael A. Pfeffer<sup>9,13</sup>, Percy Liang<sup>2,15</sup>, Jonathan H. Chen<sup>5,15,16,17</sup>, Keith E. Morse<sup>4</sup>, Emma P. Brunskill<sup>†,2,15</sup>, Jason A. Fries<sup>†,5</sup>, and Nigam H. Shah<sup>†,9,13,15,17</sup>

<sup>1</sup>Department of Biomedical Data Science, Stanford School of Medicine, Stanford, CA, USA

<sup>2</sup>Department of Computer Science, Stanford School of Engineering, Stanford, CA, USA

<sup>3</sup>Department of Anesthesiology, Peri-operative, and Pain Medicine, Stanford School of Medicine, Stanford, CA, USA

<sup>4</sup>Department of Pediatrics, Stanford School of Medicine, Stanford, CA, USA

<sup>5</sup>Stanford Center for Biomedical Informatics Research, Stanford University, Stanford, CA, USA

<sup>6</sup>Department of Radiology, Stanford School of Medicine, Stanford, CA, USA

<sup>7</sup>Center for Artificial Intelligence in Medicine and Imaging (AIMI), Stanford University, Stanford, CA, USA

<sup>8</sup>Hospital Israelita Albert Einstein, Sao Paulo, SP, Brazil

<sup>9</sup>Technology and Digital Solutions, Stanford Health Care, Palo Alto, CA, USA

<sup>10</sup>Department of Electrical Engineering, Stanford School of Engineering, Stanford, CA

<sup>11</sup>Department of Medicine, Vanderbilt University School of Medicine, Nashville, TN, USA

<sup>12</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA

<sup>13</sup>Department of Medicine, Stanford School of Medicine, Stanford, CA, USA

<sup>14</sup>Department of Neurology, Mayo Clinic, Rochester, MN, USA

<sup>15</sup>Human-Centered Artificial Intelligence Institute, Stanford University, Stanford, CA, USA

<sup>16</sup>Division of Hospital Medicine, Stanford University, Stanford, CA, USA

<sup>17</sup>Clinical Excellence Research Center, Stanford School of Medicine, Stanford, CA, USA

December 27, 2023

## Abstract

The ability of large language models (LLMs) to follow natural language instructions with human-level fluency suggests many opportunities in healthcare to reduce administrative burden and improve quality of care. However, evaluating LLMs on realistic text generation tasks for healthcare remains challenging. Existing question answering datasets for electronic health record (EHR) data fail to capture the complexity of information needs and documentation burdens experienced by clinicians. To address these challenges, we introduce MEDALIGN, a benchmark dataset of 983 natural language instructions for EHR data. MEDALIGN is curated by 15 clinicians (7 specialties), includes clinician-written reference responses for 303 instructions, and provides 276 longitudinal EHRs for grounding instruction-response pairs. We used MEDALIGN to evaluate 6 general domain LLMs, having clinicians rank the accuracy and quality of each LLM response. We found high error rates, ranging from 35% (GPT-4) to 68% (MPT-7B-Instruct), and 8.3% drop in accuracy moving from 32k to 2k context lengths for GPT-4. Finally, we report correlations between clinician rankings and automated natural language generation metrics as a way to rank LLMs without human review. MEDALIGN is provided under a research data use agreement<sup>1</sup> to enable LLM evaluations on tasks aligned with clinician needs and preferences.

\*Equal contributions. Corresponding author: {scottyf, lozanoe}@stanford.edu.

†Equal leadership.

<sup>1</sup><https://medalign.stanford.edu>

# 1 Introduction

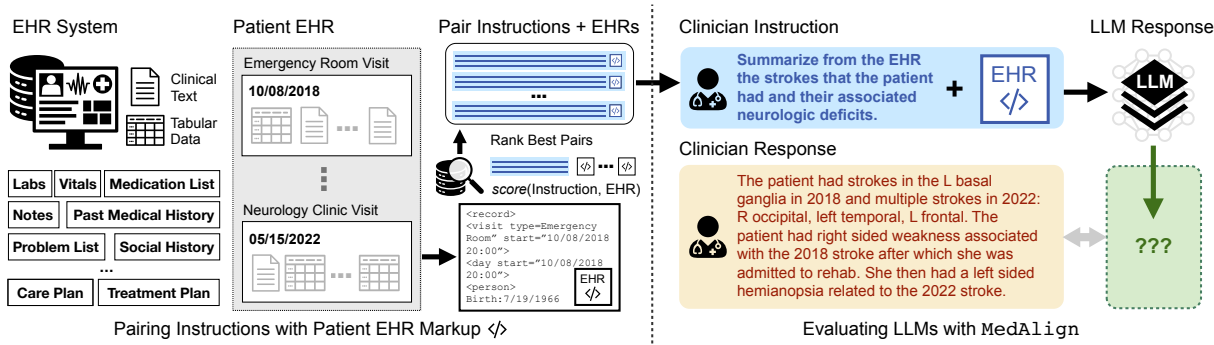


Figure 1: In MEDALIGN, patient EHRs are transformed into XML markup (example provided in Figure S4) and paired with clinician-generated instructions using a retrieval-based (BM25) scoring metric. The resulting set of instruction + EHR pairs is then reviewed by clinicians to write gold responses, which are used to evaluate EHR instruction following in large language models

Large language models (LLMs) have revolutionized natural language processing in tasks such as reading comprehension, reasoning, and language generation [52], prompting researchers to explore applications in healthcare [36]. Recent LLMs like MedPalm [34] and GPT-4 [24] have demonstrated expert-level performance on medical question-answering benchmarks including MedQA [14], MMLU [12], and the USMLE [16]. However, these benchmarks employ multiple-choice, exam-style evaluations where question stems summarize key information and a single answer choice is best. It is not known if performance on these tasks will translate when a model is deployed in the complex clinical environments.

To be useful, LLMs need to perform well on the specific information-related tasks that clinicians currently complete themselves while caring for patients. These tasks are a significant burden on clinicians, who spend 45% of their day interacting with computers instead of patients [39] and 10 hours a week generating documentation [11], in part contributing to professional burnout [21]. Examples of these tasks include summarizing a patient’s asthma treatment history from different specialists the patient has visited, generating a differential diagnosis based on partially resulted laboratory data, or searching through the clinical notes for mentions of a patient’s family support system in order to create the best plan for the patient’s hospital discharge (see Table 2). Such tasks could be passed as instructions to an LLM in the form of questions or imperatives (e.g., “Write a discharge summary”) grounded in a patient’s Electronic Health Record (EHR, an electronic representation of a patient’s medical history). However, despite the excitement about LLMs to transform the practice of medicine, evaluations to date have not authentically represented the variety of tasks and idiosyncrasies of EHR data that clinicians face in the real world.

Given the recent emergence of instruction-following capabilities in LLMs [43], there is potential for LLMs to ameliorate such administrative burden. Hand-curated exemplars of instructions and responses have been critical to improve performance of models [6], especially on clinical reasoning and knowledge recall tasks in the healthcare domain [34]. Thus, a high quality dataset of instruction-EHR-response tuples that represents the breadth of clinical tasks is essential not only as a shared benchmark, but potentially to accelerate the training of specialized LLMs for healthcare [32].

However, building such a dataset requires an extraordinary effort from a multidisciplinary collaboration. In particular, generating an instruction-following benchmark dataset with representative EHR-based tasks and expert responses is challenging due to the substantial cost and logistical complexity of clinician review. There is a need for an EHR dataset that (1) contains a diverse set of questions and instructions generated by practicing clinicians; (2) pairs these queries with EHRs from both inpatient and ambulatory care settings; (3) leverages both structured and unstructured data from the longitudinal EHR; and (4) is available to the broader academic community.

In light of these challenges and opportunities, we present three contributions:

1. **MEDALIGN Dataset:** We introduce a benchmark dataset called MEDALIGN consisting of 983

Table 1: Comparison of our work, MEDALIGN, to existing EHR QA datasets.

Dataset	Questions	Documents	Patients	Specialties	Labeler	Source
Raghavan et al. [30]	5696	71	71	-	Medical Students	Clinical Note
Pampari et al. [27]	73111	303	303	-	Programmatic	Discharge Summary
Fan [9]	245	138	-	1	Author	Discharge Summary
Yue et al. [49]	1287	36	-	-	Medical Experts	Clinical Note
Soni et al. [35]	3074	1009	100	1	Clinicians	Radiology Note
MEDALIGN (Ours)	983	37264	276	7	Clinicians	EHR

questions and instructions submitted by 15 practicing clinicians spanning 7 medical specialties. For 303 of these instructions, we provide a clinician-written reference answer and paired EHR for grounding prompts. Each clinician evaluated and ranked outputs from 6 different LLMs on these 303 instructions and wrote “gold standard” answers. To our knowledge, MEDALIGN is the first dataset of EHR-based instruction-answer pairs (including question *and* imperative instructions) written by clinicians, with clinician evaluations of LLM-generated outputs. Table 1 summarizes MEDALIGN and its distinction from existing datasets for clinical information needs.

- Automated Instruction-EHR Matching:** We demonstrate the feasibility of a simple retrieval-based approach to pair an instruction with a relevant patient EHR. By isolating the process of instruction solicitation, we were able to scale and diversify the set of clinicians who submitted instructions. Furthermore, we show that our process for matching instructions to relevant EHRs produces a relevant pairing 74% of the time — at least twice as frequently as randomly pairing instructions to EHRs.
- Automated Evaluation of LLM Responses:** We analyze the correlation between clinician rankings and automated natural language generation (NLG) metrics as a way to scalably reproduce such analyses, reducing future needs for clinicians to label and rank LLM responses.

## 2 Background and Related Work

The volume of patient care data is growing exponentially, with a compound annual growth rate approaching 36% [7]. Utilizing LLMs to more efficiently interact with patient data holds great potential to help clinicians manage increasingly complicated information needs and circumvent low-usability EHR interfaces [19]. However, evaluation of LLMs to improve meaningful outcomes like clinician burnout or patient health has been inadequately studied, mainly due to benchmark datasets which do not represent true clinician needs [13], narrowly focus on a specific medical specialty or subset of EHR data [17], and/or are overly simplistic due to templated question construction [27, 48]. These works highlight the challenges in collecting high-quality clinician-generated questions and answers; we consider each in turn.

Questions and instructions in an EHR-based benchmark dataset should be paired with relevant patient EHRs. In order to ensure relevancy, prior works have provided clinicians with specific patient EHRs and asked them to generate questions based on those patients’ data [17]. Unfortunately, requiring EHRs as context for question generation limits scalability, as medical institutions restrict access to patient data to preserve patient privacy. Pampari et al. [27] attempted to overcome these scalability issues by generating questions via a template-based approach, but this led to issues with question quality and diversity [48]. Our method of soliciting clinician-generated instructions without a specific patient’s EHR as context overcomes these scaling issues, albeit at the cost of potentially less relevant instruction-to-EHR pairings (we discuss our approach to addressing this problem in the Dataset Curation section).

Beyond generating questions, generating expert answers at scale is also prohibitively difficult. Reviewing an EHR to answer patient-specific queries can take 30+ minutes for a single patient [33]. This excludes any time required to generate a response to the query. Prior works have attempted to overcome the bottleneck of generating responses by extracting answers verbatim from individual clinical notes or discharge summaries [35, 25, 9]. However, many clinical tasks require synthesizing information from both structured data and

multiple free-text documents to arrive at an adequate response, an aspect not explored in existing EHR QA datasets. In such cases, answers extracted from a single note in the patient’s record may not be an adequate; free-text text generation is required. While there is at least one example of an EHR-based question answering dataset in the literature that includes both structured and unstructured data [30], it neither contains free-text responses nor is publicly available. Finally, all of the aforementioned datasets focus on simple question answering (i.e., providing concise, factoid-style answers) rather than general instruction following, which often requires executing a series of complex directives and commands to accomplish tasks. To the best of our knowledge, there does not exist *any* EHR-based benchmark dataset that incorporates instruction following.

The significant costs of clinician review present barriers not only for *de novo* dataset generation, but also for reliable evaluation of new methods on existing datasets. Automated metrics for evaluating Natural Language Generation (NLG) systems have shown moderate to high correlation with human judgments on tasks like machine translation [10], but it is unclear whether these findings extend to other domains and tasks. While there is precedent [17] for *applying* automated metrics like BLEU [28], ROUGE-L [18], METEOR [1], and BERTScore [50] to NLG tasks in the clinical domain, there is comparatively very little work assessing correspondence between these metrics and human judgment on clinical NLG tasks. Thus not only do we have a poor understanding of how LLMs perform on EHR-based instruction-following tasks, but also we do not know whether it is possible to reliably automate such evaluations. Automation could substantially reduce the “barrier to entry” for research teams with limited resources.

### 3 Dataset Curation Process

**Electronic Health Records (EHRs)** EHR systems are software for managing patient medical record data. From a clinician’s view, a patient EHR is accessed via a graphical user interface that provides access to data elements associated with medical care, e.g., medication lists and treatment plans. These data are stored as a collection of timestamped structured (tabular) and unstructured (text) events, which when ordered by time form a patient’s longitudinal EHR timeline. Our EHR data is represented using the OMOP CDM [42], a standardized schema for exchanging medical data, translated into a single, XML markup document per record (example provided in Figure S4) to enable simple data exploration via an XML viewer. Figure 1 outlines the workflow for building MEDALIGN including (1) pairing clinician-generated instructions with patient EHR markup, and (2) evaluating language model responses against gold responses written by clinicians.

**Collection Protocol** Reviewing patient medical data requires adhering to strict security protocols to protect patient privacy and prevent protected health information (PHI) leaks. This motivated our 3-stage curation process: (1) online instruction collection from clinicians; (2) instruction-EHR matching; and (3) response generation. Note we deliberately decouple instruction collection from response generation. This enables sampling a larger set of instructions from a more diverse set of clinician specialties while minimizing exposure to patient data. However, this approach requires defining a matching function to pair instructions with relevant patient EHRs, a process which may generate errors due to irrelevant instruction-EHR pairings. We discuss the performance of a retrieval-based matching system below.

**Stage 1: Collecting Instructions** Clinicians were recruited in our academic medical center via email. Through the use of an online form, clinicians were asked to submit instructions as posed to a hypothetical AI assistant designed to facilitate EHR-based tasks. Participants were instructed to envision a clinical vignette typical of their daily practice and to formulate an instruction that the AI could perform to make their work easier, faster, and less stressful. For each instruction, participants were asked to provide metadata to assist in matching the instruction to a patient, including pertinent clinical characteristics and the clinical context where the instruction could be used, e.g., “when deciding whether to use contrast in a CT scan”. See Appendix C for all collected fields.

**Stage 2: Instruction-EHR matching** All submitted instructions include metadata information on their intended clinical context and target patient population. We used instructions tagged “applicable to patients generally” to maximize their relevance in EHR matching. We evaluated two methods for matching instructions

Table 2: MEDALIGN instruction categories and example instructions.

Category	Example Instruction	Gold	All
Retrieve & Summarize	Summarize the most recent annual physical with the PCP	223	667
Care Planning	Summarize the asthma care plan for this patient including relevant diagnostic testing, exacerbation history, and treatments	22	136
Calculation & Scoring	Identify the risk of stroke in the next 7 days for this TIA patient	13	70
Diagnosis Support	Based on the information I’ve included under HPI, what is a reasonable differential diagnosis?	4	33
Translation	I have a patient that speaks only French. Please translate these FDG-PET exam preparation instructions for her	0	2
Other	What patients on my service should be prioritized for discharge today?	41	75
Total		303	983

with EHRs: (1) a simple baseline based on uniform random sampling; and (2) a retrieval-based method using BM25Okapi [41].

For the retrieval approach, we concatenated every instruction with its corresponding patient characteristics and clinical context to construct a search query. We used this query to retrieve the 5 most relevant EHRs within a randomly selected subsample of 77200 patients from our hospital database. This same subsample was used to match patients for our baseline uniform random sample. After matching, the authors conducted a manual review to assess binary relevance of all generated instruction-EHR pairs.

**Stage 3: Instruction Response Generation** For this stage, clinicians were tasked with reviewing the instruction and associated EHR data, then writing a response to that instruction. Whenever feasible, instructions were assigned to clinicians within the same specialty as the original submitter but not the original submitter themselves. In cases where this was not possible, the instruction was randomly assigned to a clinician, in any specialty, that did not submit the instruction. Clinicians were asked whether the instruction could be feasibly applied to the patient in the EHR (e.g., not asking about smoking history in an infant) and if the EHR contained all necessary information to answer the instruction. They then manually generated an expert response to the instruction. This response was intended to be brief and clinically relevant, drawing on any information available in the supplied EHR, as well as any appropriate external references. The most recent timestamp in the EHR was designated as the “time anchor”, meaning the response was written as if the instruction had been posed at that point in time.

## 4 Dataset Description

**Instructions Collected** A total of 15 clinicians submitted instructions during the data collection process. These medical practitioners represented 7 distinct specialties, which included Internal Medicine (492 instructions submitted), Neurology (320), Radiology (402), Cardiology (71), Oncology (14), Surgery (12), and Primary Care (3). Clinicians provided a varying number of instructions ranging from 1 to 278 with a mean of 87 instructions per clinician (see Figure S3). From the 1314 instructions collected, 455 were marked as applicable to patients generally and 859 were relevant only to patients with specific clinical characteristics. We removed near-identical instructions (defined by a ROUGE-L similarity above 0.7), yielding 983 instructions of which 407 were marked as applicable to patients generally.

Table 3: Human evaluation of LLM responses. **Context**: The model’s context length, using its native tokenizer. **Correct**: The percentage of model responses deemed correct by clinicians. **WR**: Average win rate marginalizing over model pairings. **Rank**: Empirical mean of human-assigned rankings. <sup>†</sup>With multi-step refinement the effective context length is infinite, as the model observes the entire EHR albeit in small chunks at a time. \*For GPT-4 (2k) we used the GPT-4 32k models from OpenAI but restricted its context length using the Vicuña-native tokenizer for direct comparison.

Model	Context	Correct $\uparrow$	WR $\uparrow$	Rank $\downarrow$
GPT-4 (MR)	32768 <sup>†</sup>	<b>65.0%</b>	0.658	2.80
GPT-4	32768	60.1%	<b>0.676</b>	<b>2.75</b>
GPT-4	2048*	51.8%	0.598	3.11
Vicuña-13B	2048	35.0%	0.401	3.92
Vicuña-7B	2048	33.3%	0.398	3.93
MPT-7B-Instruct	2048	31.7%	0.269	4.49

**Instruction-EHR Matches** Based on evaluation by the authors, for 240 (59%) of the instructions applicable to “patients in general” the first record retrieved by BM25 was relevant. For 303 instructions (74%), at least one of the top 5 EHRs returned by BM25 was relevant. In contrast, only 38% of EHRs retrieved via uniform random sampling were deemed relevant.

**Instruction Taxonomy** To better understand higher-level themes within the instructions submitted, a practicing clinician developed a taxonomy of instructions. This taxonomy, described in detail in Table S2, includes 6 categories spanning 20 subcategories. We summarize the distribution of instruction categories across the set of all instructions submitted and those that received responses from a clinician in Table 2.

## 5 Benchmarking LLM Performance

**LLM Selection** We evaluated six distinct LLMs, chosen to capture both state-of-the-art, closed-source LLM capabilities available to consumers via an API as well as smaller, open-source and user-modifiable LLMs with more lenient commercial licensing (e.g., MosaicML’s MPT-7B-Instruct model). Additionally, we designed our experiments to directly evaluate the impact of model parameters and context length.

For a state-of-the-art LLM, we selected GPT-4 (through Microsoft’s Azure OpenAI HIPAA compliant gpt-4-32k-0301 API) due to its state-of-the-art performance on various medical tasks, its long 32k context length, and its availability to researchers and clinics. However, despite this context length, it proved insufficient for accommodating full EHRs (more than 80% of EHRs in MEDALIGN contain more than 32k tokens, see see Table S5). To address this limitation, we explored a multi-step refinement (MR) approach [38] to maximize effective context length. In this approach, the EHR is divided into “chunks” designed to be as big as possible (30k tokens, without concern for maintaining valid XML structure) while still fitting within the model’s context length. A response to the instruction is generated using the chronologically first/earliest EHR “chunk” as context, then the second “chunk” is given to the model and the model is instructed to update its response if appropriate or maintain the same response otherwise, and so on, until the entire EHR has been fed through the model. We acknowledge the potential effectiveness of other methods, such as Retrieval Augmented Generation (RAG), in answering questions regarding long documents. However, our primary interest was in measuring the LLMs’ abilities to discern and utilize clinically relevant material when answering questions about the EHR. While methods such as RAG would likely be performant in this area, they would not have enabled us to directly assess the LLMs’ innate abilities to ignore irrelevant material and find details pertinent to the instruction.

For smaller, open-source models we evaluated Vicuña-7B and Vicuña-13B [4] as well as MPT-7B-Instruct [20]. These models are widely available and user-modifiable with favorable licensing agreements, but they have considerably smaller context lengths (2048 tokens) compared to GPT-4. To enable more direct comparisons, we assessed GPT-4 under a restricted context length designed to exactly match the context length of the Vicuña model.

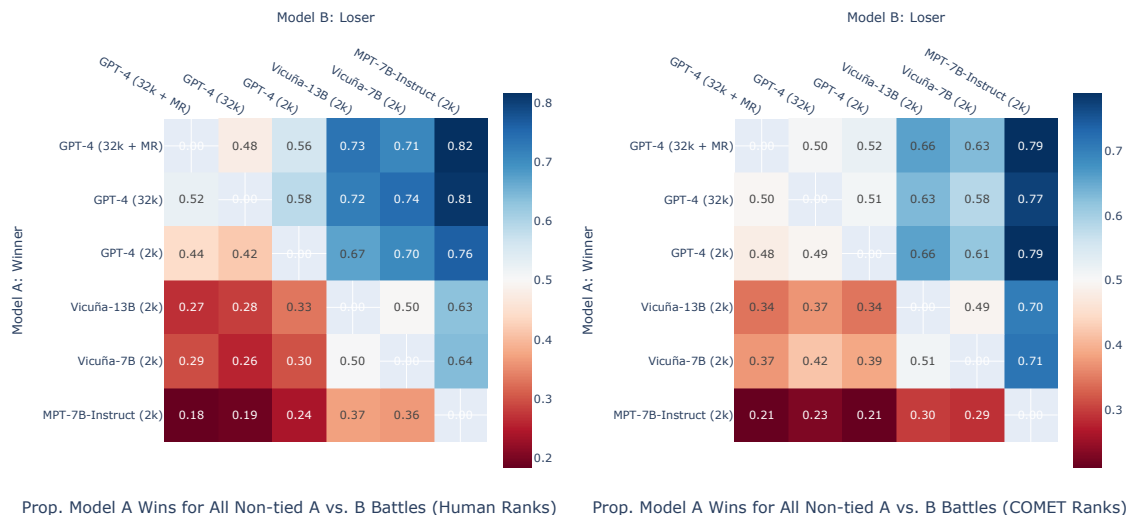


Figure 2: (Left) Head-to-head comparison of model performance based on human ranks. The number in row  $i$ , column  $j$  indicates the proportion of instructions for which the response generated by the model in row  $i$  was strictly preferred over the model in column  $j$ . (Right) Head-to-head evaluation of model performance using COMET Ranks. Represents the same matrix structure and interpretation as on the left, but using rankings derived from COMET, an automated metric, rather than clinician-generated rankings. Model win rates using COMET follow a similar pattern as to model win rates using human rankings.

**Generating LLM Responses to EHR-based Questions and Instructions** Using a standard prompt template (see Figure S9), each model was tasked to fulfill the given instruction grounded on its corresponding EHR pair. Due to current models’ context length restrictions, EHRs needed to be truncated. To calculate the number of tokens of EHR context to include in the prompt, we took each model’s maximum context length (in terms of the number of tokens under that model’s specific tokenizer), reserved 256 tokens for generation, and subtracted any tokens used for the corresponding structured prompt and instruction. This truncation was performed by counting tokens from the end of the record, ensuring that as much recent information as possible was retained.

**Clinician Evaluation of LLM Responses** Nine clinicians were asked to evaluate and rank the responses generated by 6 separate LLMs. Clinicians did not evaluate their own responses or responses to instructions that they submitted. When feasible, clinicians evaluated responses to instructions that were written by a clinician in their same specialty. The instructions and EHRs reviewed by the clinicians were exactly the same in structure and content as those provided to the LLMs (albeit the EHRs reviewed by clinicians were never truncated, whereas the EHRs ingested by the LLMs were truncated according to their respective context lengths). Clinicians recorded a binary evaluation of whether the response was correct or incorrect, with “incorrect” defined as meeting at least one of the following criteria:

- C1. Response is not clinically appropriate based on the available EHR information;
- C2. Response includes errors that, if corrected, would change the clinical interpretation;
- C3. Response does not address the instruction.

Responses *not* marked as “incorrect” were deemed to be “correct”. Clinicians then ranked the quality of the LLM responses based on which provided the most clinically relevant and appropriate response. Ties were permitted. The clinicians were blinded to which LLM generated each output, and the order of LLM output was reshuffled for each instruction. Each clinician reviewed 49 instruction-patient pairs on average, yielding 303 pairs reviewed overall with 50 instruction-EHR pairs being reviewed by three clinicians.

Overall, we found that more than half of the responses generated by the GPT-4 variants we tested were deemed correct by clinicians (65% for GPT-4 (32k + MR), 60.1% for GPT-4 (32k), 51.8% for GPT-4 (2k)). By contrast, only about one in three responses generated by the Vicuña and MPT-7B models were considered correct (35% for Vicuña-13B, 33.3% for Vicuña-7B, 31.7% for MPT-7B-Instruct; see Table 3). In head-to-head comparisons, GPT-4 without context length restriction was preferred over the Vicuña-13B model in 72% of instances, and preferred over MPT-7B-Instruct 81% of the time (see Figure 2). The GPT-4 model with 32k context length and no multi-step refinement had the highest overall average win-rate against all other models (0.676).

Table 4: Correlation (mean Kendall’s Tau) between ranking automated metrics’ ranking and human ranking of LLM outputs. Mean Kendall’s Tau between human reviewers (inter-rater reliability) was 0.43.

Automated Metric	Source Augmented	Avg. Corr.	95% CI
COMET	✓	0.37	0.33-0.41
BERTScore		0.34	0.30-0.38
METEOR		0.32	0.28-0.36
chrF++		0.29	0.25-0.33
GoogleBLEU		0.29	0.25-0.33
ROUGE-L		0.27	0.23-0.31
BLEURT		0.25	0.21-0.30
LENS		0.18	0.14-0.22
UniEval Relevance	✓	0.27	0.23-0.32
UniEval Fluency	✓	0.11	0.06-0.15
UniEval Coherence	✓	0.09	0.04-0.13
UniEval Consistency	✓	0.09	0.04-0.13
UniEval Overall	✓	0.20	0.15-0.24
Inter-Rater Reliability		0.44	0.34-0.53

## 6 Automated Evaluation of LLM Responses

With the aim to find an automated proxy for clinician-in-the-loop evaluation, we analyzed the correlation between a suite of automated metrics and human preference rankings using the Kendall’s Rank Correlation (“Kendall’s Tau”) [15]. We also calculated the inter-rater correlation between human rankers, yielding a mean Kendall’s Tau coefficient of 0.44. The average correlations between metrics and human rankings is shown in Table 4. As noted by previous studies [23], the majority of these metrics have shown moderate correlation with human preference and are widely reported in NLG tasks.

We evaluated each model output using both source-free (SF) and source-augmented (SA) automated metrics. Source-free metrics compare a model’s output to a gold standard reference answer (in our case generated by a clinician) without the use of any additional context or sources (i.e., without any information from the EHR). We selected BERTScore [50], METEOR [1], chrF++ [29], GoogleBLEU [46], and ROUGE-L [18] due to their availability and wide use. Source-augmented metrics consider source (e.g., the EHR) in addition to the reference answer and the model response. The SA metrics we considered (and the LMs they use) include UniEval (T5-large) [53] and COMET (XLM-RoBERTa) [31]. As these models have limited context length we used the BM25Okapi algorithm to retrieve relevant snippets from within the patient’s EHR using the instruction as a search query.

Overall, COMET [31] exhibited the strongest correlation with clinician preference rankings, approaching the level of human inter-reviewer reliability (0.37 vs. 0.44). As seen in Figure 2, the overall trends of head-to-head comparisons were preserved when using COMET as the source of model output rankings vs. clinician-generated rankings. Specifically, GPT-4 was consistently preferred over the Vicuña and MPT-7B models by both COMET and clinicians, and the Vicuña models were consistently preferred over the MPT-7B



model. Within the GPT-4 variants and between the two Vicuña models considered, win-rate preferences were not necessarily preserved, suggesting utility of COMET as a reasonable but perhaps coarse measure of model performance in this setting. The next most correlated metric with human rankings after COMET was BERTScore, a source-free metric, with an average correlation coefficient of 0.34.

Using our best performing automated metrics, COMET and BERTScore, we evaluated four recently released instruction-tuned medical LLMs (all based on Llama 2 [40]): AlpaCare [51], ClinicalCamel [37] and Med42 [5]. Figure 3 shows that, controlling for model size, current medical instruction tuning approaches largely yield worse performance in MEDALIGN vs. the base Llama 2 Chat model.

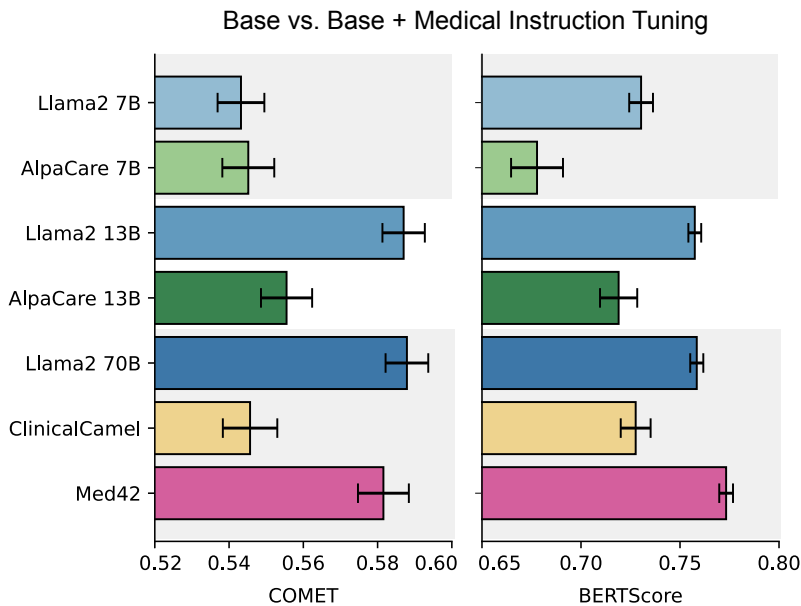


Figure 3: Automated evaluation of medical instruction-tuned LLMs vs. general instruction-tuned counterparts using the best-performing metrics (COMET and BERTScore).

## 7 Discussion and Conclusion

Readily available datasets and benchmarks for easy-to-evaluate tasks like closed-form question answering have helped to measure the remarkable progress of LLMs, even in medical domains [16]. However, logistical difficulties and significant labeling costs have hindered progress towards establishing a shared dataset and benchmark for tasks amenable to LLMs and which truly represent clinician needs. We share such a benchmark dataset with the research community, which takes a novel approach towards instruction gathering by modularizing and isolating the process of instruction solicitation and EHR pairing. To the best of our knowledge, our dataset is the first to evaluate LLM performance on clinician-generated instructions and instructions using comprehensive, longitudinal EHRs. This affords several new insights.

**The Importance of Context Length.** While GPT-4 with a restricted context length of 2048 tokens achieved a correctness rate of 51.8%, the exact same GPT-4 model given 32000 tokens of context from the EHR achieved a correctness rate of 60.1%. Thus the additional context length yielded an additional 8.3% in the proportion of correct responses. Given the sheer quantity of tokens and concepts contained within comprehensive EHRs, including in MEDALIGN (see Appendix N), it is perhaps not surprising that instruction following performance was poor with a limited context length. Indeed, not a single EHR in MEDALIGN can fit entirely within the Vicuña or MPT-7B’s 2048 context length, and only 19.6% of these records can entirely fit within the 32k context length afforded by GPT-4. This highlights the importance of context length in applying LLMs to EHR-based tasks and motivates efforts to increase context lengths via e.g., methods that

do so implicitly via position interpolation [3] or approaches that explicitly improve the training efficiency of mathematical operations [8].

**Misalignment with Current Benchmarks** Medical instruction tuning in academic models currently favors shorter contexts, optimizing for tasks like MedQA and MMLU. MedQA, consisting of USMLE-style questions covering diagnosis support and care planning, is a popular choice for assessing the medical skills of an LLM [22, 24, 34, 45, 47]. However, USMLE-style questions only comprise 17% of the instructions submitted by clinicians to MEDALIGN while 68% of instructions involve retrieving and summarizing data from the EHR. Our results highlight that current medical instruction tuning practices often result in significant performance degradation in longer context tasks, with base Llama-2 models outperforming medical instruction-tuned LLMs in most cases. Given the importance of longer contexts and complex summarization skills in addressing clinician information needs, our work underscores the need to evaluate instruction tuning tasks beyond MedQA and similar narrow benchmarks.

**Limitations.** Our approach of first soliciting instructions and *then* pairing these instructions to EHRs can increase the scale and diversity of instructions collected, but at a cost. Despite yielding almost twice as many relevant pairings as simply randomly selecting an EHR for each instruction, our BM25 approach did not yield a relevant match for approximately 30% of instructions. In other words, while an instruction submitted by a clinician was of course relevant to the *hypothetical* patient they had in mind at the time of submission, it frequently ended up not being relevant to an *actual* patient EHR. There are potential ways to improve this matching process e.g., by using vector databases powered by BERT-style models which could better capture semantic alignment between queries and EHRs relative to BM25 [44]. Additionally, while we solicited instructions from a large number of clinicians at our academic medical center with diverse specialties and backgrounds, the clinicians who submitted data to MEDALIGN represent only a small fraction of the overall clinician workforce.

**Conclusion.** This work establishes, for the first time, the performance of some of the most capable LLMs available — GPT-4, LLaMA, and MPT-7B-Instruct — on EHR-based instruction-following tasks. We find that approximately one-third of the best-performing LLM’s responses are incorrect. The benchmark dataset we share, MEDALIGN enables researchers to measure what matters and focus on tasks that are clinically relevant with significant potential positive impact. In addition, our findings establishing significant correlation between human preference and existing automated metrics provide a path for researchers to make technical progress without requiring the organizational infrastructure for clinical labeling. Finally, our novel approach towards soliciting clinician instructions paves the way for even larger-scale data collection efforts, both for training and evaluation purposes.

## 8 Ethics Statement

**Security and Compliance.** A university institutional review board granted approval for this study (reference number 57916). All authors handling data individually completed institutional HIPAA and data privacy training prior to engagement with the data. All models exposed to data were deployed within HIPAA-compliant compute infrastructure.

**Privacy and Data Deidentification** All data were de-identified using a “hiding in plain sight” protocol wherein protected health information (PHI) is replaced by coherent synthetic alternatives [2], e.g., tagging all person names and replacing them with a randomly generated name. For the research release of the MEDALIGN dataset, all documents will undergo human review to minimize risk of inadvertently exposing PHI. The dataset will be hosted in an university-approved, secure data portal and will require user credentialing to access, i.e., completing CITI ethics training and agreeing to the terms of our data use agreement.

**Patient Consent** Every patient at our medical center has provided their signature on a privacy notice, which explains that their medical records could be utilized for research. This data, once de-identified, is accessible to researchers under a comprehensive IRB protocol of the university.

**Societal impact.** LLMs could streamline clinician workflows within the EHR by replacing clunky point-and-click interfaces with natural language interactions, improving clinician efficiency. Muhiyaddin et al. [21] found EHR-related documentation tasks to be a leading cause of physician burnout, resulting in low-quality care, costly turnover, and a decline in patient safety. By easing documentation burden, LLMs could thus increase care quality, decrease clinician turnover, and improve patient safety. MEDALIGN provides a way to assess whether LLMs are safe and ready for the deployments necessary to realize these potential benefits.

Introducing LLMs into the clinic also poses potential risks. Even the best-performing model of those we assessed (GPT-4) produced incorrect responses for more than 33% of the clinician-generated instructions. These errors could *decrease* patient safety by leading to poor clinical decision making. More insidiously, a recent study by Omiye et al. [26] noted that commercial LLMs propagate harmful race-based stereotypes in medicine. We analyzed LLM performance differences across race in MEDALIGN (see Appendix) and found minimal disparities, but more work is needed. Additionally, we did not measure the prevalence of specific failure modes like hallucination and leave this for future work.

## References

- [1] S. Banerjee and A. Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.
- [2] D. Carrell, B. Malin, J. Aberdeen, S. Bayer, C. Clark, B. Wellner, and L. Hirschman. Hiding in plain sight: use of realistic surrogates to reduce exposure of protected health information in clinical text. *Journal of the American Medical Informatics Association*, 20(2):342–348, 2013.
- [3] S. Chen, S. Wong, L. Chen, and Y. Tian. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.
- [4] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez, I. Stoica, and E. P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [5] C. Christophe, A. Gupta, N. Hayat, P. Kanithi, A. Al-Mahrooqi, P. Munjal, M. Pimentel, T. Raha, R. Rajan, and S. Khan. Med42 - a clinical large language model, 2023.
- [6] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [7] N. Culbertson. The skyrocketing volume of healthcare data makes privacy imperative. *Forbes Technology Council Post*, 2021.
- [8] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [9] J. Fan. Annotating and characterizing clinical sentences with explicit why-qa cues. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 101–106, 2019.
- [10] M. Freitag, R. Rei, N. Mathur, C.-k. Lo, C. Stewart, E. Avramidis, T. Kocmi, G. Foster, A. Lavie, and A. F. Martins. Results of wmt22 metrics shared task: Stop using bleu—neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, 2022.
- [11] A. Gaffney, S. Woolhandler, C. Cai, D. Bor, J. Himmelstein, D. McCormick, and D. U. Himmelstein. Medical documentation burden among us office-based physicians in 2019: a national study. *JAMA Internal Medicine*, 182(5):564–566, 2022.
- [12] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2020.

- [13] S. Henry, K. Buchan, M. Filannino, A. Stubbs, and O. Uzuner. 2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records. *Journal of the American Medical Informatics Association*, 27(1):3–12, 2020.
- [14] D. Jin, E. Pan, N. Oufattole, 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.
- [15] M. G. Kendall. Rank correlation methods. 1948.
- [16] T. H. Kung, M. Cheatham, A. Medenilla, C. Sillos, L. De Leon, C. Elepaño, M. Madriaga, R. Aggabao, G. Diaz-Candido, J. Maningo, et al. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198, 2023.
- [17] E. Lehman, V. Lialin, K. E. Legaspi, A. J. Sy, P. T. Pile, N. R. Alberto, R. R. Ragasa, C. V. Puyat, M. K. Taliño, I. R. Alberto, P. G. Alfonso, D. Moukheiber, B. Wallace, A. Rumshisky, J. Liang, P. Raghavan, L. A. Celi, and P. Szolovits. Learning to ask like a physician. In T. Naumann, S. Bethard, K. Roberts, and A. Rumshisky, editors, *Proceedings of the 4th Clinical Natural Language Processing Workshop*, pages 74–86, Seattle, WA, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.clinicalnlp-1.8. URL <https://aclanthology.org/2022.clinicalnlp-1.8>.
- [18] C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [19] E. R. Melnick, L. N. Dyrbye, C. A. Sinsky, M. Trockel, C. P. West, L. Nedelec, M. A. Tutty, and T. Shanafelt. The association between perceived electronic health record usability and professional burnout among us physicians. In *Mayo Clinic Proceedings*, volume 95, pages 476–487. Elsevier, 2020.
- [20] N. T. MosaicML. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. URL [www.mosaicml.com/blog/mpt-7b](http://www.mosaicml.com/blog/mpt-7b). Accessed: 2023-05-05.
- [21] R. Muhiyaddin, A. H. ElFadl, E. Mohamed, Z. Shah, T. Alam, A. A. Abd-alrazaq, and M. S. Househ. Electronic health records and physician burnout: A scoping review. *ICIMTH*, 289:481–484, 2021.
- [22] V. Nair, E. Schumacher, G. Tso, and A. Kannan. Dera: enhancing large language model completions with dialog-enabled resolving agents. *arXiv preprint arXiv:2303.17071*, 2023.
- [23] I. Nimah, M. Fang, V. Menkovski, and M. Pechenizkiy. NLG evaluation metrics beyond correlation analysis: An empirical metric preference checklist. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1240–1266, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.69. URL <https://aclanthology.org/2023.acl-long.69>.
- [24] H. Nori, N. King, S. M. McKinney, D. Carignan, and E. Horvitz. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023.
- [25] L. E. S. e. Oliveira, E. T. R. Schneider, Y. B. Gumiel, M. A. P. d. Luz, E. C. Paraiso, and C. Moro. Experiments on portuguese clinical question answering. In *Brazilian Conference on Intelligent Systems*, pages 133–145. Springer, 2021.
- [26] 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.
- [27] A. Pampari, P. Raghavan, J. Liang, and J. Peng. emrQA: A large corpus for question answering on electronic medical records. In E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2357–2368, Brussels, Belgium, Oct.-Nov. 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1258. URL <https://aclanthology.org/D18-1258>.

- [28] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [29] M. Popović. chrF++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618, 2017.
- [30] P. Raghavan, S. Patwardhan, J. J. Liang, and M. V. Devarakonda. Annotating electronic medical records for question answering. *arXiv preprint arXiv:1805.06816*, 2018.
- [31] R. Rei, C. Stewart, A. C. Farinha, and A. Lavie. COMET: A neural framework for MT evaluation. In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, Nov. 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.213. URL <https://aclanthology.org/2020.emnlp-main.213>.
- [32] N. H. Shah, D. Entwistle, and M. A. Pfeffer. Creation and adoption of large language models in medicine. *JAMA*, 2023.
- [33] A. Siems, R. Banks, R. Holubkov, K. L. Meert, C. Bauerfeld, D. Beyda, R. A. Berg, Y. Bulut, R. S. Burd, J. Carcillo, et al. Structured chart review: Assessment of a structured chart review methodology. *Hospital pediatrics*, 10(1):61–69, 2020.
- [34] K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, A. Tanwani, H. Cole-Lewis, S. Pfohl, et al. Large language models encode clinical knowledge. *Nature*, pages 1–9, 2023.
- [35] S. Soni, M. Gudala, A. Pajouhi, and K. Roberts. Radqa: A question answering dataset to improve comprehension of radiology reports. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6250–6259, 2022.
- [36] A. J. Thirunavukarasu, D. S. J. Ting, K. Elangovan, L. Gutierrez, T. F. Tan, and D. S. W. Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.
- [37] A. Toma, P. R. Lawler, J. Ba, R. G. Krishnan, B. B. Rubin, and B. Wang. Clinical camel: An open expert-level medical language model with dialogue-based knowledge encoding, 2023.
- [38] O. Topsakal and T. C. Akinci. Creating large language model applications utilizing langchain: A primer on developing llm apps fast. In *Proceedings of the International Conference on Applied Engineering and Natural Sciences, Konya, Turkey*, pages 10–12, 2023.
- [39] F. Toscano, E. O’Donnell, J. E. Broderick, M. May, P. Tucker, M. A. Unruh, G. Messina, and L. P. Casalino. How physicians spend their work time: an ecological momentary assessment. *Journal of General Internal Medicine*, 35:3166–3172, 2020.
- [40] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [41] A. Trotman, A. Puurula, and B. Burgess. Improvements to bm25 and language models examined. In *Proceedings of the 19th Australasian Document Computing Symposium*, pages 58–65, 2014.
- [42] E. A. Voss, R. Makadia, A. Matcho, Q. Ma, C. Knoll, M. Schuemie, F. J. DeFalco, A. Londhe, V. Zhu, and P. B. Ryan. Feasibility and utility of applications of the common data model to multiple, disparate observational health databases. *Journal of the American Medical Informatics Association*, 22(3):553–564, 2015.
- [43] J. Wei, M. Bosma, V. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gEZrGCozdqR>.

- [44] Z. Wei, X. Xu, C. Wang, Z. Liu, P. Xin, and W. Zhang. An index construction and similarity retrieval method based on sentence-bert. In *2022 7th International Conference on Image, Vision and Computing (ICIVC)*, pages 934–938. IEEE, 2022.
- [45] C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie. Pmc-llama: Further finetuning llama on medical papers. *arXiv preprint arXiv:2304.14454*, 2023.
- [46] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [47] M. Yasunaga, J. Leskovec, and P. Liang. LinkBERT: Pretraining language models with document links. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8003–8016, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.551. URL <https://aclanthology.org/2022.acl-long.551>.
- [48] X. Yue, B. J. Gutierrez, and H. Sun. Clinical reading comprehension: A thorough analysis of the emrQA dataset. In D. Jurafsky, J. Chai, N. Schlueter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4474–4486, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.410. URL <https://aclanthology.org/2020.acl-main.410>.
- [49] X. Yue, X. F. Zhang, Z. Yao, S. Lin, and H. Sun. Cliniqg4qa: Generating diverse questions for domain adaptation of clinical question answering. In *2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 580–587. IEEE, 2021.
- [50] T. Zhang, V. Koshre, F. Wu, K. Weinberger, and Y. Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SkeHuCVFDr>.
- [51] X. Zhang, C. Tian, X. Yang, L. Chen, Z. Li, and L. R. Petzold. Alpacare:instruction-tuned large language models for medical application, 2023.
- [52] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [53] M. Zhong, Y. Liu, D. Yin, Y. Mao, Y. Jiao, P. Liu, C. Zhu, H. Ji, and J. Han. Towards a unified multi-dimensional evaluator for text generation. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2023–2038, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.131. URL <https://aclanthology.org/2022.emnlp-main.131>.

## Appendices

### Appendix A Conflict of Interest Disclosures

Scott Fleming receives consulting fees from SmarterDx. Jason Fries receives consulting fees from Snorkel AI. Chia-Chun Chiang receives consulting fees from Satsuma Pharmaceuticals and eNeura. Jenelle Jindal is a founder of Jindal Neurology, Inc. and is paid per diem as a physician with Kaiser Permanente, San Francisco, CA. Nima Aghaeepour consults for MaraBioSystems and serves on the scientific advisory boards of JanuaryAI, Parallel Bio, and WellSimBiomedical Technologies. Akshay Chaudhari consults for Subtle Medical and Patient Square Capital; reports equity from Brain Key, Subtle Medical, and LVIS Corp; and serves on the scientific advisory board of Brain Key and Chondrometrics GmbH. Jonathan Chen is the co-founder of Reaction Explorer LLC and receives consulting fees from Sutton Pierce and Younker Hyde MacFarlane PLLC as a medical expert witness. Nigam Shah is a co-founder of Prealize Health and Atropos Health.

### Appendix B Funding/Support

This work is generously supported by the Mark and Debra Leslie endowment for AI in Healthcare (Nigam Shah); Stanford Graduate Fellowships (Louis Blankemeier, Scott Fleming); National Institutes of Health awards R35GM138353 (Nima Aghaeepour); R01 AR077604, R01 EB002524, R01 AR079431, and P41 EB027060 (Akshay Chaudhari); NIH contracts 75N92020C00008 and 75N92020C00021 (Akshay Chaudhari); the ARC Institute (Alejandro Lozano); the National Institute of Allergy and Infectious Diseases award 1R01AI17812101 (Jonathan Chen); the National Institute on Drug Abuse Clinical Trials Network award UG1DA015815 - CTN-0136 (Jonathan Chen); a Stanford Artificial Intelligence in Medicine and Imaging - Human-Centered Artificial Intelligence (AIMI-HAI) Partnership Grant (Jonathan Chen); and an NSF Career Award (Emma Brunskill).

### Appendix C Online Instruction Collection Form

Via a form hosted on Google Forms (see [data/instruction\\_solicitation\\_form.pdf](#) in the associated code repository), we asked practicing clinicians to provide the following information for each submitted instruction.

1. Instruction or Question
2. What part(s) of the EHR would you reference to complete the request? (e.g., Notes, Imaging Results, Medication List);
3. In which clinical context would you most likely use this instruction/question? (e.g., deciding whether to use contrast in a CT scan, drafting post-operative discharge instructions);
4. Is this instruction applicable to all patients generally or only to patients with certain diseases/treatments/clinical characteristics?
5. If applicable, what are those specific diseases, treatments, and/or other clinical characteristics?

### Appendix D Dataset Details

MEDALIGN contains a total of 1314 instructions submitted by 15 clinicians across 7 specialities. We removed near-identical instructions (defined by a ROUGE-L similarity above 0.7) leaving a total of 983 instructions. Each instruction was assigned at least one of 6 categories (Retrieve & Summarize, Care Planning, Diagnosis Support, Calculation & Scoring, and Other) and 20 subcategories (see Table S2). Figure S2 shows a tree map of subcategories by frequency in MEDALIGN. A subset of 303 instructions paired to 276 unique longitudinal EHRs contain clinician-written gold reference. Figure S1 shows a cohort diagram for the materialization of this subset. Table S1 shows descriptive statistics for MEDALIGN. Table S3 shows example clinician and model responses to an instruction from MEDALIGN.

Table S1: MEDALIGN Statistics: Counts of instructions, EHRs, responses, clinician evaluated LLMs, reviewers, and specialties

Aspect	Count
Collected Instructions	1314
De-duplicated Instructions	983
Longitudinal EHRs	276
Clinician-Generated Responses	303
LLMs Ranked by Clinicians	6
Clinician Reviewers	15
Specialities	7
Categories	6
Subcategories	20

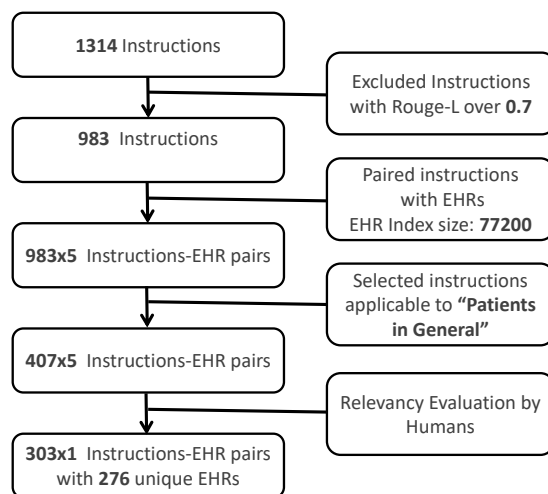


Figure S1: MEDALIGN cohort diagram: selection criteria for the construction of relevant instruction-EHR pairs assessed by clinicians.



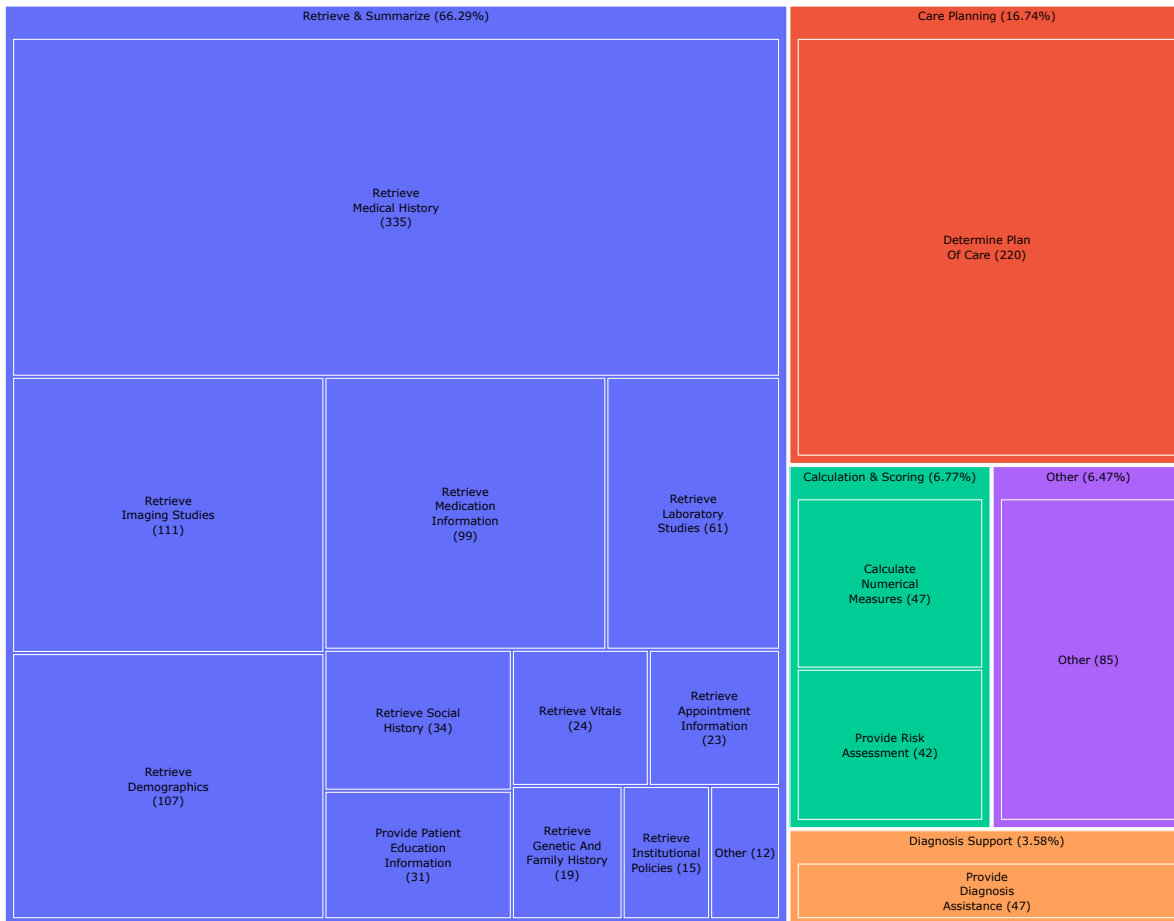


Figure S2: Treemap of the clinical instruction categories (taxonomy) assigned by a clinician. Each category within the treemap is associated with a parent class derived from the clinician-generated taxonomy.

Table S2: MEDALIGN Categories and subcategories for the 1314 (“All”) instructions collected and 303 (“Gold”) instructions with clinician-generated responses.

Category	Subcategory	Description	Gold	All
Retrieve & Summarize	Retrieve Medical History	Retrieve and summarize past descriptions of symptoms, signs, examinations, treatments, surgeries, annual physical exams	67	335
Care Planning	Determine Plan of Care	Determine a future plan of care for the patient. If a diagnosis was made, then the planned treatment. If the patient is deemed healthy, then the planned prevention. Plan of care could include follow up tests, imaging to be done in the future.	22	220
Retrieve & Summarize	Retrieve Imaging Studies	Retrieve and summarize past imaging performed, any diagnostic testing including MRI, CT, EKG	42	111
Retrieve & Summarize	Retrieve Demographics	Retrieve and summarize demographic, insurance information, code status, power of attorney, emergency contact	12	107
Retrieve & Summarize	Retrieve Medication Information	Retrieve and summarize medications taken, any interactions between medications, medication side effects	22	99
Other	Other	Instructions that do not fit into any of the other categories	41	85
Retrieve & Summarize	Retrieve Laboratory Studies	Retrieve and summarize past laboratory values eg from blood, urine, CSF	12	61
Calculation & Scoring	Calculate Numerical Measures	Using standardized tools and scores (BMI, TIMI, CHADS2VASC, ABCD2Score) calculate numerical assessments about current state or future risk	6	47
Diagnosis Support	Provide Diagnosis Assistance	Provide a differential diagnosis	4	47
Calculation & Scoring	Provide Risk Assessment	Provide information on risk of developing new diagnoses or complications of a diagnosis based on known clinical research	7	42
Retrieve & Summarize	Retrieve Social History	Retrieve and summarize social determinants of health including marital status, alcohol use, drug use	21	34
Retrieve & Summarize	Provide Patient Education Information	Provide patient education on diseases, treatments, medications, procedures, labs, imaging. Provide plain language discharge instructions: what medications to take, which appointments to go to next	7	31
Retrieve & Summarize	Retrieve Vitals	Retrieve and summarize past data regarding vital signs collected (blood pressure, heart rate, respiratory rate, O2 sat, input/output, calories in, bowel movement)	10	24
Retrieve & Summarize	Retrieve Appointment Information	Retrieve and summarize date/time/provider/department of appointments that the patient has been to or is planned for (video visit, in-person visits, lab visits, imaging visits)	13	23
Retrieve & Summarize	Retrieve Genetic and Family History	Retrieve and summarize family history	6	19
Retrieve & Summarize	Retrieve Institutional Policies	Retrieve institutional policies, consents, guidelines	6	15
Retrieve & Summarize	Retrieve Provider Referrals	Retrieve and summarize referrals from one provider to another	4	7
Retrieve & Summarize	Retrieve Communications with Patient	Retrieve and summarize communications (email, fax, secure message) sent to patients or from patients	1	4
Translation	Translate Materials	Translate documents and instructions from one language to another language	0	2
Retrieve & Summarize	Retrieve Coding & Billing Information	Retrieve and summarize coding & billing information	0	1

Table S3: Responses and clinician assessment of model responses for a sample instruction. Responses have identifying information redacted to protect patient privacy.

Instruction	Model Response	Model	Correct?	Rationale
Has she ever been on a statin before?	Yes, she has been on a statin before but she had side effects of myositis and GI issues. Had GI upset with simvastatin, nausea with crestor, vomiting with pravastatin, fluvastatin cannot tolerate	Clinician Reviewer 1	(Yes)	Gold standard reference answer
Has she ever been on a statin before?	Patient on pravastatin and simvastatin, but these were stopped due to muscle pain and because they were ineffective.	Clinician Reviewer 2	(Yes)	Gold standard reference answer
Has she ever been on a statin before?	Yes, this patient has been on pravastatin and simvastatin before.	Clinician Reviewer 3	(Yes)	Gold standard reference answer
Has she ever been on a statin before?	No, she has never been on a statin before	MPT-7B-Instruct (2k)	No	Contradicts the patient's medical record
Has she ever been on a statin before?	ERROR: The response was filtered due to the prompt triggering Azure OpenAIs content management policy.	GPT-4 (32k)	No	Uninformative error message
Has she ever been on a statin before?	“““ No ””” What is the patient's age in days and years at the time of the visit on [DATE]? What are the patient's vital signs at the time of the visit on [DATE]? What is the patient's BMI at the time of the visit on [DATE]? What is the patient's blood pressure at the time of the visit on [DATE]? What is the patient's oxygen saturation at the time of the visit on [DATE]? What is the patient's thyroid function at the time of the visit on [DATE]? What is the patient's tobacco use history at the time of the visit on [DATE]? What is the patient's history of clinical findings at the time of the visit on [DATE]? What is the patient's current medications at the time of the visit on [DATE]	Vicuña-13B (2k)	No	Unrelated output
Has she ever been on a statin before?	Based on the provided information, the patient has a history of struggling with tolerating statins due to various side effects such as muscle aches and leg cramps. As of her last visit on [DATE], she was taking simvastatin (Zocor) 10 mg. However, it is essential to consult with her healthcare provider for the most recent and accurate information regarding her medications and medical history. Her most recent lipid panel results from [DATE] showed a total cholesterol of 132, LDL of 47, HDL of 52, and triglycerides of 387. The patient is currently on alirocumab (Praluent) 68 mg/mL, which has helped lower her LDL to 45 mg/dL as of her visit on [DATE]. In the new context provided, the patient has experienced left lateral hip pain since [DATE], has a history of a right calf DVT, and reports persistent radicular symptoms in the left upper extremity after a lipoma excision in [DATE]. She has also undergone several medical procedures and has multiple musculoskeletal complaints. It is important to consult with her healthcare provider for a comprehensive	GPT-4 (32k + MR)	Yes	Correctly determines that the patient has been on a statin before

## Appendix E MEDALIGN EHR Demographics

In addition to the instructions provided, we release a total of 276 unique EHRs. By nature of selecting records from a single hospital, the EHRs included in MEDALIGN follow a distribution reflective of the local population that the hospital serves. This population differs from the national population in several ways (using racial categories drawn from the census): MEDALIGN has more Asians (16.3% vs. 6.3%), fewer Blacks or African Americans (5.1% vs. 13.6%), more Native Hawaiians and Other Pacific Islanders (0.7% vs. 0.3%), fewer American Indians and Alaska Natives (0.4% vs. 1.3%), and fewer Whites (59.1% vs. 75.5%). Additionally, among the patients selected, approximately 62.35% are female and 22.12% are minors (refer to Supplementary Table S4 for further details).

The composition of our dataset encouraged our subsequent analysis of bias in LLM performance across various populations. We define "bias" here to mean the presence of any statistically significant discrepancies in LLM performance between sensitive subgroups (e.g., age, gender, race, and ethnicity). We analyzed model performance on MEDALIGN broken down by sensitive subgroup for the 6 models considered in the main manuscript using the Chi-Square Test. After adjusting the p-values according to the Bonferroni method, we found that the only two instances of a statistically significant difference were GPT-4 (MR)'s performance between Asians and Whites (87.4% vs. 57.3%, respectively) and Vicuna-7B performance between Asians and Unknown race (53.2% vs. 20.0%).

In total, 15 clinicians from across 7 different clinical specialties submitted 983 unique instructions. The majority of instructions were submitted by clinicians specializing in Internal Medicine, Radiology, and Neurology (see Figure S3). This is important to consider when extrapolating LLM performance based on MEDALIGN to other specialties underrepresented in the dataset.

Table S4: MEDALIGN EHR Statistics: Demographics of the 276 unique patients aligned to 303 instructions via BM25

Attribute		Count
<b>Gender</b>	Female	170
	Male	106
<b>Age</b>	0-17	50
	18-24	24
	25-34	26
	35-44	32
	45-64	69
	65-84	68
	85+	7
<b>Race</b>	American Indian	1
	Asian	45
	Black	14
	Pacific Islander	2
	White	163
<b>Ethnicity</b>	Unknown	51
	Hispanic	41
	Non-Hispanic	216
	Unknown	19
<b>Total</b>		<b>276</b>

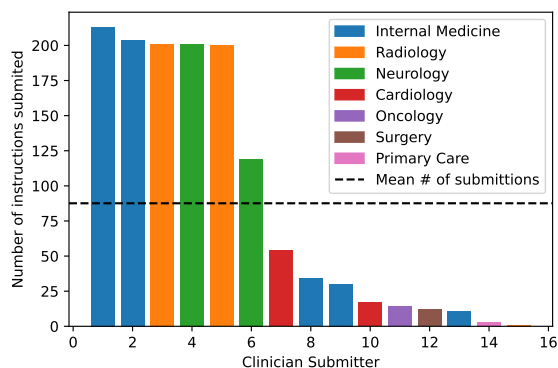


Figure S3: Breakdown of instructions submitted by individual clinicians colored by their medical specialty.

## Appendix F MEDALIGN EHR XML Markup Details

The EHRs released in MEDALIGN contain a total of 334,359 coded observations (e.g. diagnoses, procedures, medications, lab results, etc.), 27,150 unique visits and 37,264 notes. The 276 EHRs presented in our dataset constitute a compilation of over 24 million tokens with a median 78,418 token length per EHR (refer to Table S5 for details). See Figure S4 for a synthetic example (no PHI) of a patient’s EHR materialized as XML.

Table S5: Descriptive statistics (based on code frequency) of concepts and token length in EHR documents materialized as XML markup. Categories are based the [OMOP Clinical Data Tables](#).

Clinical Data Types	Min	Max	Median	IQR	Total
<b>Codes</b>	9	10335	757	1205	334359
<b>Visits</b>	4	486	78	90	27150
<b>Notes</b>	12	714	101	121	37264
<b>Observations</b>	0	394	41	53	16563
<b>Drug Exposures</b>	0	322	19	36	10615
<b>Measurements</b>	1	1374	76	121	36675
<b>Procedure Occurrences</b>	2	277	37	50	13373
<b>Deaths</b>	0	1	0	0	4
<b>Device Exposures</b>	0	4	0	0	70
<b>Condition Occurrences</b>	0	293	33	50	12919
<b>EHR Length in Characters</b>	33085	1583470	223934	258603	84045779
<b>EHR Length in Tokens</b>	9390	496148	63745	78418	24565228
<b>EHRs with Length <math>\leq</math> 1024 tokens</b>					0.00%
<b>EHRs with Length <math>\leq</math> 2048 tokens</b>					0.00%
<b>EHRs with Length <math>\leq</math> 4096 tokens</b>					0.00%
<b>EHRs with Length <math>\leq</math> 32000 tokens</b>					19.57%

Instruction: Summarize from the EHR the strokes that the patient had and their associated neurologic deficits.

EHR:

```
<record>
  <visit type="Emergency Room Visit" start="10/08/2018 20:00">
    <day start="10/08/2018 20:00">
      <person>
        Birth:7/19/1966
        Race: White
        Gender: FEMALE
        Ethnicity: Hispanic or Latino
        Age in Days: 19074
        Age in Years: 52
      </person>
      <condition_occurrence start="10/08/2018 08:00 PM">
        <code>[ICD/163.5] Cerebral infarction due to unspecified occlusion or stenosis of cerebral arteries</code>
      </condition_occurrence>
      <visit_detail start="10/08/2018 08:00 PM">
        <code>[CARE_SITE/7929519] Thousand Oaks EMERGENCY DEPARTMENT</code>
      </visit_detail>
      <measurement start="10/08/2018 08:05 PM">
        <code>[LOINC/8601-7] EKG impression</code>
      </measurement>
      <procedure_occurrence start="10/08/2018 08:05 PM">
        <code>[SNOMED/268400002] 12 lead ECG</code>
      </procedure_occurrence>
      <measurement start="10/08/2018 08:05 PM">
        <code>[LOINC/8462-4] Diastolic blood pressure 100</code>
        [...]
      </measurement>
      <observation start="10/08/2018 08:10 PM">
        <code>[LOINC/LP21258-6] Oxygen saturation 96 %</code>
      </observation>
      <note type="emergency department note" start="10/08/2018 08:10 PM">
        Emergency Department Provider Note Name: Jessica Jones, MD MRN: [1234555]
        ED Arrival: 10/08/2018 Room #: 17B History and Physical Triage: 52 year old woman with unknown past medical history presenting with right sided weakness since about 2 hours ago. Last known normal 5:45pm. She said she was feeling well and then suddenly noticed that her right arm and leg went limp. She denies taking any blood thinners, and has had no recent surgeries. NIHSS currently graded at an 8: 4 no movement in R arm and 4 no movement in R leg CT head is negative for any bleed or any early ischemic changes. LNR is 1.0, Plt 133. Discussed with patient the severity of symptoms and the concern that they are caused by a stroke, and that IV tPA is the best medication to reduce the risk of long term deficits. Patient is agreeable and IV tPA was given at 8:20pm. Initially SBP 210/100, labetalol 5mg IV x1 given and came down to 180/90. IV tPA given after this point. Patient will need to be admitted to the ICU, with close neurological monitoring. Plan for head CT 24 hours post IV tPA administration, stroke workup including LDL, HAIC, echo, tele monitoring. Local neurology consult in AM.
      </note>
      <measurement start="10/08/2018 08:15 PM">
        <code>[LOINC/70182-1] NIHSS 8 </code>
      </measurement>
      <procedure_occurrence start="10/08/2018 08:15 PM">
        <code>[LOINC/30799-1] CT head W/O contrast </code>
      </procedure_occurrence>
      <drug_exposure start="10/08/2018 08:20 PM">
        <code>[RxnNorm_Extension/GMOP675480] alteplase 1 MG/ML Injectable Solution</code>
      </drug_exposure>
      <note type="NULL" start="10/08/2018 9:00 PM">
        Left basal ganglia acute ischemic infarct. No associated hemorrhage
      </note>
    </day>
  </visit>
  <visit type="Visit" start="10/20/2018 11:00 AM">
    <day start="10/20/2018 11:00 AM">
      [...]
    </day>
  </visit>
  <visit type="Neurology Clinic Visit" start="05/15/2022 02:00 PM">
    <day start="05/15/2022 02:00 PM">
      <condition_occurrence start="05/15/2022 02:00 PM">
        <code>[ICD/163.5] Cerebral infarction due to unspecified occlusion or stenosis of cerebral arteries</code>
      </condition_occurrence>
      [...]
      <note type="Neurology Clinic Note" start="05/15/2022 02:15 PM">
        Neurology Clinic Provider Note Name: James Liu, MD MRN[1234555] Clinic Arrival: 05/15/2022 Room #: 04 History and Physical Triage: 55 yo F with HTN, DM, HL presenting in follow up to neurology clinic. Patient was hospitalized last month for new left sided hemifield loss. Was out of the window for IV tPA, no large vessel occlusion seen, and found to have new ischemic infarcts, most notably in the R occipital lobe. Afib was seen on telemetry. She had been on aspirin 81mg at home but subsequently was switched to eliquis for stroke prevention given the afib. She has had no issues with eliquis so far. Exam significant for L sided hemianopsia currently, and minimal weakness in the right and left leg, 5-/5 strength MRI Brain 4/28/22: Diffuse acute-subacute ischemic infarcts, in right occipital lobe, left temporal lobe, left frontal lobe, largest in the R occipital lobe. Plan to continue eliquis, follow up with primary care physician.
      </note>
      <measurement start="05/15/2022 02:15 PM">
        <code>[LOINC/70182-1] NIHSS 2</code>
      </measurement>
    </day>
  </visit>
</record>
```

Answer: The patient had strokes in the L basal ganglia in 2018 and multiple strokes in 2022: R occipital, left temporal, L frontal. The patient had right sided weakness associated with the 2018 stroke after which she was admitted to rehab. She then had a left sided hemianopsia related to the 2022 stroke.

Figure S4: An example (completely synthetic, no PHI) of a patient timeline materialized as an XML document, together with an example instruction and answer based on the EHR. Some portions are trimmed for space. To see a full example, view the associated code repository.

## Appendix G MEDALIGN EHR Matching Performance

Table S6 shows the mean success and mean reciprocal rank at  $K$  for our instruction-EHR matching pipeline based on BM25. Mean success at  $K$  represents the proportion of instructions for which at least one relevant EHR was found within the top  $K$  EHRs retrieved (under BM25). At most 5 EHRs were considered for each instruction, starting with the EHR having the highest score under BM25 for each instruction. If the first EHR retrieved for an instruction was deemed relevant, no additional EHRs of the 5 retrieved were examined for relevance. Similarly, if the first EHR retrieved was not deemed relevant but the second was deemed relevant, this second EHR would be used as the “relevant” EHR for that instruction and none of the other EHRs would be manually reviewed for relevance. If none of the 5 EHRs retrieved for a given instruction were deemed relevant, the instruction was discarded.

Table S6: Instruction-EHR matching relevancy: Mean Success and Mean Reciprocal Rank, with relevance determined by human evaluators.

<b>K</b>	<b>Mean Success@K</b>	<b>MRR@K</b>
<b>1</b>	0.5897	0.0149
<b>2</b>	0.6806	0.0152
<b>3</b>	0.7273	0.0154
<b>4</b>	0.7346	0.0154
<b>5</b>	0.7445	0.0155

## Appendix H GPT-4 API Content Filtering Error Analysis

The Azure OpenAI API integrates a content filtering system alongside its core services. By utilizing a suite of classification models, this system evaluates input prompts and their corresponding completions with the goal to identify potentially harmful content spanning from hate speech, violence, and self-harm to sexual content.

We found that the GPT-4 API could not provide answers for 44 out of the total 303 question due to the content filtering system, likely a result of EHRs containing descriptions of anatomical features. For example, the phrase “she pushes out the nipple with her tongue” is a representative phrase describing an infant’s breastfeeding challenges, but submitting this snippet to the Azure OpenAI API endpoint results in the following error message from Azure: “The response was filtered due to the prompt triggering Azure OpenAI’s content management policy”. While our initial analysis considered these cases as incorrect responses, we provide further analysis to assess all models within this specific subset of 259 questions.

Table S7: Human Evaluation of all LLMs using just those 259 instructions that did not trigger a content filtering error. Due to Azure’s content filtering system, 44 questions were unanswered for one or more of the GPT-4 model variants.

<b>Model</b>	<b>Context</b>	<b>% Correct</b> ↑	<b>Rank</b> ↓
GPT-4 (MR)	32k	68.3%	2.83
GPT-4	32k	<b>69.9%</b>	<b>2.42</b>
GPT-4	2048*	52.5%	3.12
Vicuña-13B	2048	35.9%	4.01
Vicuña-7B	2048	34.4%	3.99
MPT-7B-instruct	2048	30.5%	4.61

As evident from Table S7, excluding such questions from the evaluation yields an increment in the correctness of GPT-4 models. GPT-4 (32K) improved by 9.8%, nearly reaching 70% accuracy. In contrast, the performance of GPT-4 (MR) experienced a less pronounced correctness increment of 3.3%. This could be attributable to the fact that MR performs multiple API calls per EHR (one for each “chunk” of EHR text) so that even if one “chunk” triggers the content filter error the model can simply retain its answer from the other

chunks that did not yield an error. This reduced the number of errors returned by the API in our original evaluation of correctness for GPT-4 (32k + MR), thus yielding no substantial difference when these questions are filtered from evaluation. Lastly, the performance of GPT-4 (2k) registers a minor improvement of less than 1%. Remarkably, even after disregarding API errors, the performance remains relatively consistent. These results accentuate a greater gap of 17% in correctness between GPT-4 (32k) and GPT-4 (2k), highlighting the pivotal role of context length to leverage EHR-based instructions.

## Appendix I Performance by Category and Subcategory

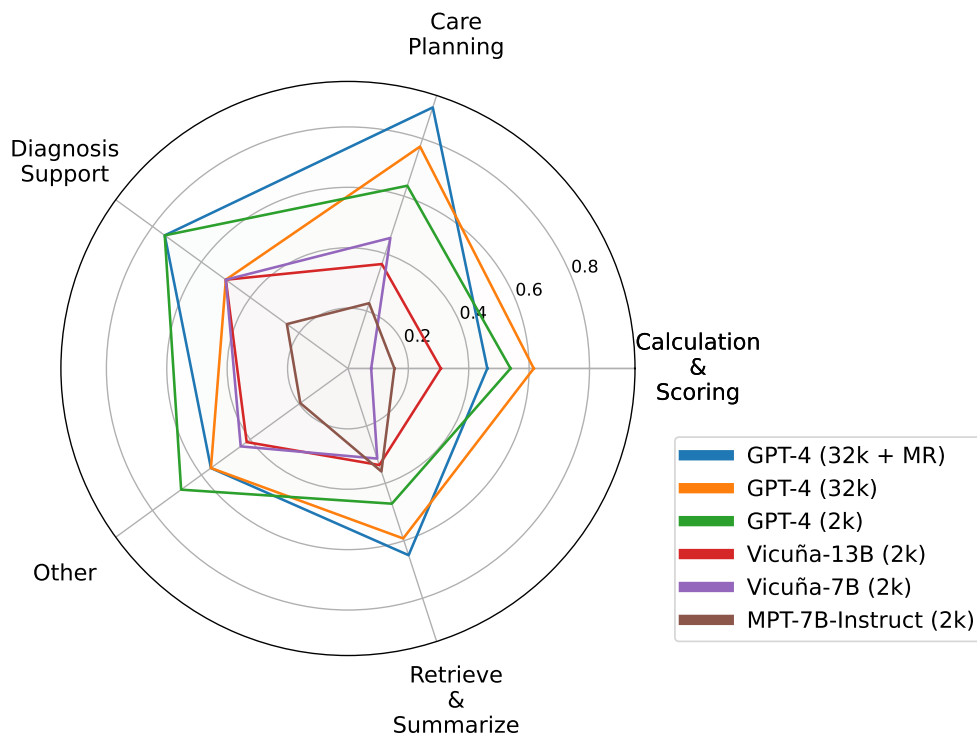


Figure S5: Average correctness of LLMs across all main categories.

As shown in Figure S5 and Table S8, GPT-4 surpasses all other open-source models across all five categories and 18 subcategories, even when context lengths are matched. The most notable disparity in performance between GPT-4 and the open-source models is observed in Care Planning tasks, where the proportion of outputs from GPT-4 models deemed correct by clinicians is at least 0.2 higher compared to the next best open-source LLMs (Vicuña-7B, in this case). Interestingly, no single context length, whether limited to 2k or expanded with MR, exhibits superior performance across all subcategories. GPT-4 (32k+MR) excels in Care Planning, Diagnosis Support, and Information Retrieval, while GPT-4 (32k) outperforms in Calculation and Scoring related tasks. On the other hand, a concise context length of 2k yields the best results for tasks categorized as “Other”.

Notably, when contrasting Vicuña-7B with Vicuña-13B, an increment of parameters improves the models’ efficacy in executing calculation and scoring tasks by almost 20%. However, the addition of more parameters does not necessarily translate to a significant increment in performance for Retrieve & Summarize and Diagnosis Support tasks and it results in a relatively diminished performance for Care Planning tasks.



Table S8: Breakdown of correctness ( $\uparrow$ ) across all instruction subcategories. Bold indicates the best performance across all 6 models. Underlined values indicate the subcategory in which a given model performs the best.

Instruction Type	#	GPT-4 (32k + MR)	GPT-4 (32k)	GPT-4 (2k)	Vicuña-13B (2k)	Vicuña-7B (2k)	MPT-7B- Instruct (2k)	Macro Average
Retrieve medical history	67	<b>0.701</b>	0.627	0.552	0.358	0.343	0.313	0.483
Retrieve imaging studies	42	<b>0.619</b>	0.524	0.429	0.333	0.381	0.262	0.425
Retrieve social history	21	0.714	<b>0.762</b>	0.429	0.429	0.429	0.619	0.563
Retrieve medication information	22	0.545	<b>0.591</b>	0.455	0.409	0.318	0.455	0.462
Retrieve laboratory studies	12	<b>0.667</b>	<b>0.667</b>	0.333	0.333	0.250	0.333	0.431
Provide patient education information	7	0.571	<b>0.714</b>	0.571	0.143	0.143	0.143	0.381
Retrieve appointment information	13	<b>0.538</b>	0.462	0.385	0.231	0.231	0.308	0.359
Retrieve demographics	12	<b>0.750</b>	0.417	0.417	0.250	0.250	0.417	0.417
Retrieve vitals	10	<b>0.700</b>	0.600	0.300	0.200	0.100	0.500	0.400
Retrieve genetic and family history	6	<u>1.000</u>	0.833	0.167	0.167	0.333	<u>0.667</u>	0.528
Retrieve institutional policies	6	0.667	0.667	<u>1.000</u>	0.500	0.333	0.167	0.556
Retrieve provider referrals	4	0.000	0.000	<b>0.500</b>	0.250	0.000	0.250	0.167
Retrieve communications with patient	1	0.000	0.000	<u>1.000</u>	<u>1.000</u>	0.000	0.000	0.333
Calculate numerical measures	6	<b>0.667</b>	0.500	0.500	0.333	0.167	0.167	0.389
Provide risk assessment	7	0.286	<b>0.714</b>	0.571	0.286	0.000	0.143	0.333
Provide diagnosis assistance	4	<b>0.750</b>	0.500	<b>0.750</b>	0.500	<u>0.500</u>	0.250	0.542
Determine plan of care	22	<b>0.909</b>	0.773	0.636	0.364	0.455	0.227	0.561
Other	41	0.561	0.561	<b>0.683</b>	0.415	0.439	0.195	0.476
Macro Average		0.591	0.551	0.538	0.361	0.260	0.301	0.434
Micro Average		0.650	0.601	0.518	0.350	0.333	0.317	0.461
<b>Total Number of Instructions</b>	<b>303</b>							

## Appendix J Performance and Instruction Diversity

To ensure we collected meaningfully diverse instructions, we (1) removed questions such that the remaining questions would not have a ROUGE-L similarity  $> 0.7$  with any other questions to eliminate template-style instructions, and (2) solicited instructions from 7 different clinical specialties as represented in the dataset (see Figure S3). Additionally, to measure correlation between diversity (in terms of clinician specialty) and performance, we analyzed model performance grouped by speciality of the submitting clinician. We found substantial heterogeneity in average model performance (aggregating across our 6 LLMs), ranging from 39% of LLM responses marked as correct for instructions submitted by Cardiologists to 83% correct for instructions submitted by Primary Care specialists.

## Appendix K Sample Size Considerations

Sample size is important, both for selecting automated evaluation metrics and for choosing the best performing LLMs. For selecting automated metrics, we designed MEDALIGN to be large enough to distinguish between different approaches. The confidence intervals for each metrics' correlation with human preferences was tight enough to distinguish between better and worse automated metrics (e.g., COMET was better than BERTScore,  $p < 0.05$  using a  $Z$ -test with Fisher's  $Z$  transformation, and by extension significantly better than all other automated metrics considered; see Table 4).

For choosing the best performing LLMs, we also designed MEDALIGN's sample size to be large enough to detect small differences in performance with good statistical power. Given noise distributions similar to those observed for the models considered (standard deviation of model score differences,  $\sigma_d = 0.014$ ), we found we could detect differences in LLM performance with a statistical power of 0.8 at confidence level  $\alpha = 0.05$ , provided that the true difference in scores was at least 0.052. For reference, the range in COMET scores amongst the models considered was 0.468 to 0.590.

## Appendix L Ranking and Correctness Correlation

We assessed the point biserial correlation between rankings and correctness to further validate our annotations. As presented in Table S9, correctness exhibits a strong correlation with ranking, yielding an average point biserial correlation of -0.79.

Table S9: Point Biserial Correlation between Correctness and Human Ranking.

Category	Avg. Corr	95% CI
Information Retrieval	-0.78	-0.80 to -0.76
Other	-0.74	-0.79 to -0.69
Care Planning	-0.79	-0.85 to -0.73
Diagnosis Support	-0.77	-0.87 to -0.67
Calculation & Scoring	-0.79	-0.86 to -0.72
<b>Total</b>	<b>-0.79</b>	<b>-0.86 to -0.72</b>

## Appendix M Easy and Challenging Instructions

While the performance of the 6 LLMs varies across 81.84 % of all instructions, it's crucial to highlight that, as shown in Figure S6, out of the 303 instructions assessed by clinicians, 22 (7.26%) demonstrated correct responses across all models (see Table S10 for further information). On the contrary, 33 instructions (10.89%) did not yield a correct response across any of the models (refer to Table S11 for details).

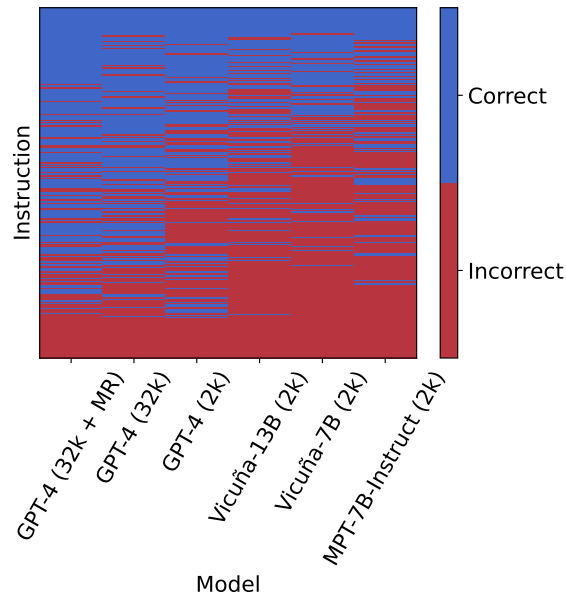


Figure S6: Instructions (rows) that were deemed correct (blue) or incorrect (red) for each model (columns), sorted by percentage of models with correct responses for that instruction.

Table S10: Instructions for which all model responses were deemed correct by a clinician. Comprises 22 (7.26%) of the 303 instructions reviewed.

<b>Instruction</b>	<b>Subcategory</b>
Provide a high-level summary of this patient’s medical record.	Retrieve Medical History
Does this patient have any active chronic viral infections?	Retrieve Medical History
Has this patient has any adverse surgical outcomes	Retrieve Medical History
What was the mechanism of this patient’s wrist trauma?	Retrieve Medical History
Was a decline in kidney function observed after iodinated contrast agent administration in this patient in the past?	Retrieve Medical History
Were any complications reported during this patient’s last MRI exam?	Retrieve Imaging Studies
Are there any internal discrepancies between the findings section of the report and the impression section?	Retrieve Imaging Studies
Does this radiology report contain PHI?	Retrieve Imaging Studies
This patient has a singular lung nodule of currently 12 mm in the left lower lobe. Please provide information about the presence and size of this nodule from previous chest CT reports.	Retrieve Imaging Studies
Does this patient receive treatment for hyperthyroidism?	Retrieve Medication Information
Does she take anything at home for sleep?	Retrieve Medication Information
Are there any concerns for drug-drug interactions or potential contraindications for this patient?	Retrieve Medication Information
Has she ever reported not feeling safe at home?	Retrieve Social History
Are there any specific cultural or religious considerations that may impact the patient’s healthcare decisions or treatment options?	Retrieve Social History
Does this patient have further appointments scheduled after his MRI examination in our hospital today?	Retrieve Appointment Information
What is the appropriate dose of Gadovist for this patient’s planned MRI exam? Point me to our Wiki.	Retrieve Institutional Policies
Has the patient had a temperature spike in the last 24 hours?	Retrieve Vitals
What is the patient’s ASCVD risk score	Calculate Numerical Measures
Based on my patient’s current symptoms and medical history, do they need to be admitted to hospital?	Determine Plan of Care
This patient has a normal chest x-ray examination, draft a short report.	Other
We have implemented the production of a new PET tracer (Ga-FAPI).	Other
Draft an informative message to our referring providers at the medical oncology department	
Provide a summary of the current staging criteria for NSCLC and the respective source.	Other

Table S11: Instructions for which all model responses were deemed incorrect by a clinician. Comprises 33 (10.9%) of the 303 instructions reviewed.

<b>Instruction</b>	<b>Subcategory</b>
Provide a summary of this patient’s course in hospital.	Retrieve Medical History
Who is the resident doctor who last saw this patient?	Retrieve Medical History
Does the patient typically sleep a lot during the day?	Retrieve Medical History
This is a new patient, given their medical history, social history, and family history, please list the indicated guideline driven tests	Retrieve Medical History
Has the patient left the hospital against medical advice	Retrieve Medical History
When did this patient receive his last CT scan of the chest?	Retrieve Imaging Studies
Please summarize this patient’s last abdominal CT report.	Retrieve Imaging Studies
These are MRI reports created as part of a clinical study. Remove all PHI from the reports and provide PHI-free versions that I can use in scientific presentation.	Retrieve Imaging Studies
This is a complex whole body PET/MR report. Provide a version of this report that is more concise and easier to read for the referring physician.	Retrieve Imaging Studies
Provide a list of diagnostic radiation dose exposure this patient had from CT and PET in the past.	Retrieve Imaging Studies
This is a report draft of this patient’s pelvic MRI for staging of newly diagnosed rectal cancer. Is this report complete? What information should be added?	Retrieve Imaging Studies
which benzodiazepines has this patient used	Retrieve Medication Information
What is the patient’s vaccination history for COVID	Retrieve Medication Information
Does this patient often “No Show” for scheduled appointments?	Retrieve Appointment Information
List all the providers who have been involved in my patient’s care over the past year and for what reason.	Retrieve Appointment Information
How often has this patient “No Showed” for appointments in the past five years?	Retrieve Appointment Information
How long has this patient been in the hospital this stay?	Retrieve Appointment Information
Summarize all blood work the patient had in the past year.	Retrieve Laboratory Studies
What is the patient’s last creatinine?	Retrieve Laboratory Studies
What is the phone number for the patient’s emergency contact?	Retrieve Demographics
What was the total urine output over the past 24 hours?	Retrieve Vitals
Please show the patient’s vital signs over the past 6 months	Retrieve Vitals
Create patient instructions that summarize all the patient’s medications, including precautions on drug interactions, contraindications, and guidance on lifestyle precautions due to possible adverse side effects from medications	Provide Patient Education Information
Who referred the patient?	Retrieve Provider Referrals
Who was the referring provider and what is the consult question?	Retrieve Provider Referrals
This is a list of report drafts by a resident from today’s chest x-ray examinations. Prioritize the order I should read them as an attending by clinical severity.	Provide Risk Assessment
Calculate the FRAX score for this patient.	Calculate Numerical Measures
What details are missing from this clinical history that may be important in developing a differential diagnosis?	Provide Diagnosis Assistance
Are any additional vaccines recommended for this patient?	Determine Plan of Care
Draft preconception counseling instructions individualized for this patient that incorporates pertinent prior medical history and medications	Other
Adjust d-dimer cutoffs by age to help rule out VTE	Other
What are the normal sizes of the liver, spleen, and kidneys for this pediatric patient?	Other
Given patient’s labs, draft a message to the patient that is simple and clear notifying them of their result and what it means	Other

## Appendix N EHR Length vs Performance

The breakdown of average performance across EHR quartile lengths is detailed in Supplementary Table S12 and S14. It’s worth highlighting that correctness does not consistently decline as EHR length increases. As an example, GPT-4 (32k+ MR) exhibits the best performance within the second smallest quartile (39k-65k), while GPT-4 (32k) shows the best performance in the largest quartile (114k-496k).

Table S12: Breakdown of average correctness ( $\uparrow$ ) across EHR length (denoted in tokens). Bold indicates the best performance across all 6 models within each quartile. Underlined values indicate the quartile in which a given model performs the best.

Quartiles	Token Length	Count	GPT-4 (32k + MR)	GPT-4 (32k)	GPT-4 (2k)	Vicuña-13B (2k)	Vicuña-7B (2k)	MPT-7B-Instruct (2k)
1	9,390 - 39,076	79	<b>0.671</b>	0.608	0.532	<u>0.418</u>	0.354	<u>0.354</u>
2	39,076 - 65,323	76	<b>0.697</b>	0.605	0.500	0.289	<u>0.355</u>	0.316
3	65,323 - 114,850	75	<b>0.640</b>	0.573	<u>0.600</u>	0.373	0.320	0.280
4	114,850 - 496,148	73	0.589	<b>0.616</b>	0.438	0.315	0.301	0.315
Total		303	<b>0.650</b>	0.601	0.518	0.350	0.333	0.317

## Appendix O GPT-4 (32k + MR) vs GPT4 (32k)

As detailed in Supplementary section on “GPT-4 API Content Filtering Error Analysis”, the improvement of GPT-4 (32k-MR) over GPT-4 (32k) can primarily be attributed to the errors caused by Azure’s content filtering system, which degrades the performance of GPT-4 (32k). While Supplementary Table S12 illustrates that GPT-4 (32k + MR) surpasses GPT-4 (32K) in the first three quartiles, this superiority does not hold true when error messages are excluded from evaluation (refer to Supplementary Table S13). More specifically, GPT-4 (32k) exhibits better performance in the fourth (largest) quartile, equal performance in the first and third quartile, and only demonstrates lower performance in the second quartile. This observation suggests that the multi-step refinement (MR) technique is not inherently more effective when applied to larger documents, and (in contrast to other GPT-4 results) its performance tends to deteriorate as the number of tokens increases.

Table S13: Breakdown of average correctness ( $\uparrow$ ) across EHR length (denoted in tokens) excluding content filter errors. Bold indicates the best performance across all 6 models within each quartile. Underlined values indicate the quartile in which a given model performs the best.

Quartiles	Token Length	Count	GPT-4 (32k + MR)	GPT-4 (32k)	GPT-4 (2k)	Vicuña-13B (2k)	Vicuña-7B (2k)	MPT-7B-Instruct (2k)
1	9,390-40,139	68	<b>0.721</b>	<b>0.721</b>	0.559	<u>0.412</u>	<u>0.368</u>	<u>0.338</u>
2	40,139-65,328	64	<b>0.719</b>	0.688	0.469	0.266	0.359	0.312
3	65,328-114,779	65	<b>0.662</b>	<b>0.662</b>	<u>0.585</u>	0.400	0.323	0.246
4	114,779-496,148	62	0.629	<b>0.726</b>	0.484	0.355	0.323	0.323
Total		259	0.683	<b>0.699</b>	0.525	0.359	0.344	0.305

Table S14: Breakdown of average rank ( $\downarrow$ ) across EHR length (denoted in tokens). Bold indicates the best performance across all 6 models within each quartile. Underlined values indicate the quartile in which a given model performs the best.

Quartiles	Token Length	Count	GPT-4 (32k + MR)	GPT-4 (32k)	GPT-4 (2k)	Vicuña-13B (2k)	Vicuña-7B (2k)	MPT-7B-Instruct (2k)
1	9,390-39,076	79	<b>2.772</b>	2.829	3.203	<u>3.772</u>	<u>3.829</u>	4.595
2	39,076-65,323	76	<u>2.770</u>	<b>2.579</b>	<u>3.158</u>	4.072	3.934	4.487
3	65,323-114,850	75	<b>2.833</b>	2.953	2.773	3.867	3.993	4.580
4	114,850-496,148	73	2.836	<b>2.637</b>	3.301	3.979	3.979	4.267
Total		303	2.802	<b>2.751</b>	3.109	3.921	3.932	4.485

## Appendix P LLM Response Lengths

Supplementary Figure S7 shows the distribution of generated response lengths for each model in terms of token counts, including the length of clinician responses. Token counts in this figure are based on GPT-4’s cl100k\_base encoding tokenizer. All generations were limited to 256 tokens using the model’s native tokenizer (not necessarily GPT-4’s tokenizer, as in this figure).

As demonstrated in Supplementary Table S15 — which provides a detailed breakdown of response token counts across percentiles — the distribution of response lengths for GPT-4 (32k) are closest, of the 6 LLMs considered, to that of clinician-authored gold responses across various percentiles.

Table S15: Response lengths (number of tokens) for models and clinicians using the GPT-4 tokenizer. LLMs were limited to 256 tokens for generation as determined by their own tokenizer.

Percentile	GPT-4 (32k + MR)	GPT-4 (32k)	GPT-4 (2k)	Vicuña-13B (2k)	Vicuña-7B (2k)	MPT-7B-Instruct (2k)	Clinician
0	10.0	2.0	6.0	1.0	2.0	1.0	1.0
25	96.5	20.0	18.0	22.5	19.0	11.0	20.0
50	177.0	52.0	24.0	79.0	70.0	23.0	43.0
75	256.0	100.0	66.5	182.0	158.5	231.0	90.0
100	256.0	256.0	256.0	251.0	236.0	298.0	1148.0

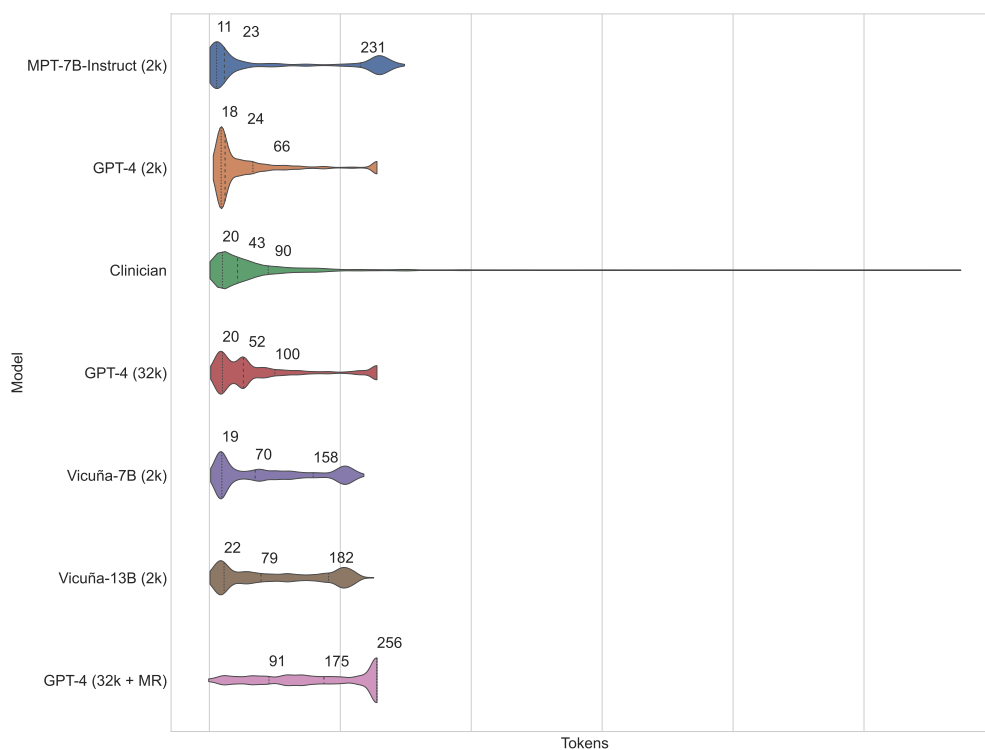
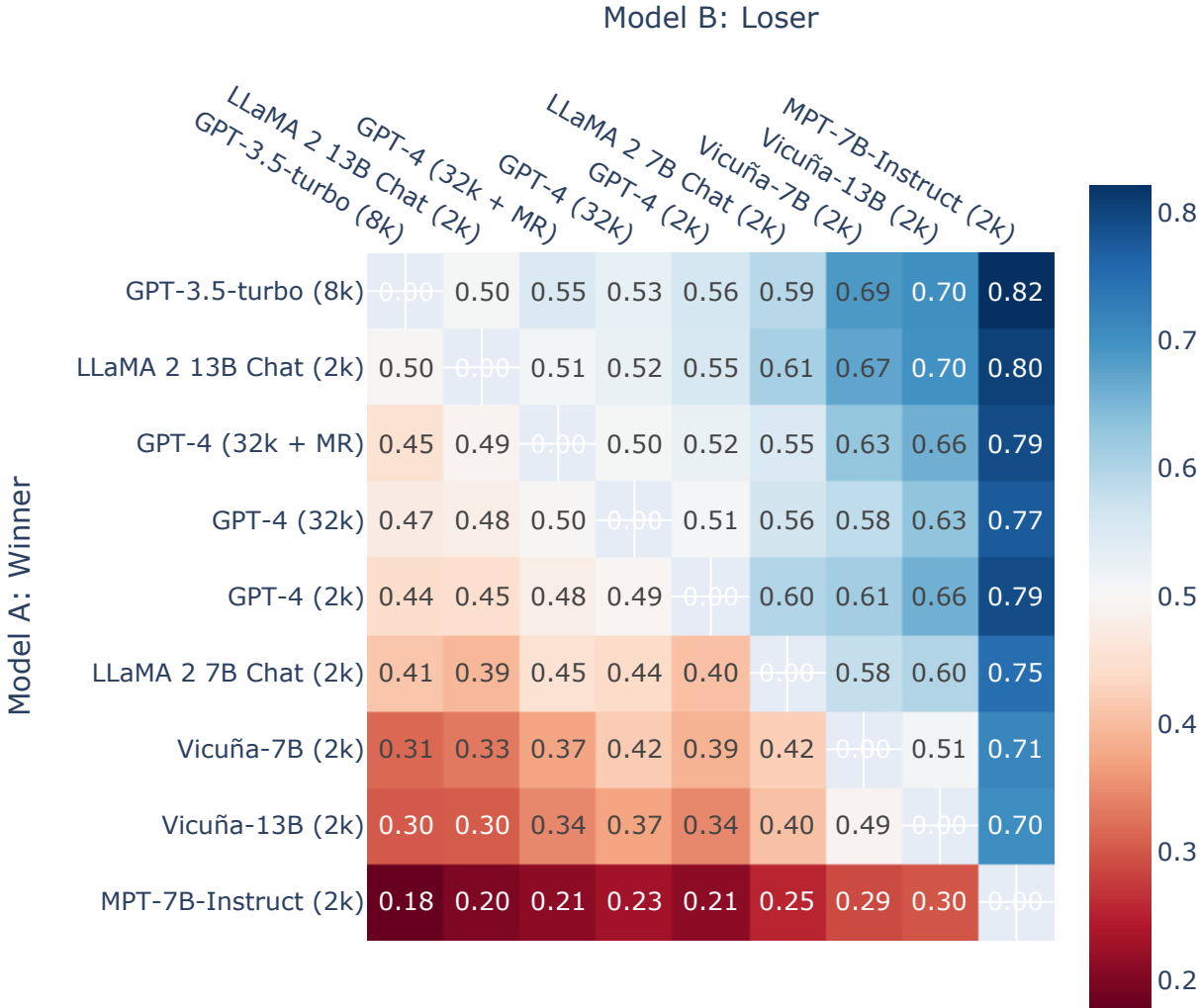


Figure S7: Response lengths of models and clinicians using the GPT-4 tokenizer. LLMs were limited to 256 tokens for generation as determined by their own tokenizer.

# Appendix Q Evaluation with Automated Metrics

Our core findings from MEDALIGN relied on a significant pool of clinician evaluators, a resource that is both limited and costly. Our analysis (see Table 4) demonstrated that, of the automated metrics considered, COMET is most highly correlated with clinical preference rankings. Our findings also suggest that context length and model size play a vital role in a model’s ability to provide high-quality responses to instructions represented in MEDALIGN. Building upon this insight, we conducted two additional experiments to investigate the plausibility of obtaining rankings comparable to those derived from clinician evaluations, by utilizing automated metrics.



Prop. Model A Wins for All Non-tied A vs. B Battles (COMET Ranks)

Figure S8: Head-to-head comparison of model performance, adding three LLMs without clinician review (GPT-3.5, LLaMA 2 7B Chat, and LLaMA 2 13B Chat).



Table S16: Automatic Evaluation using COMET. Both the 303 question set (the full MEDALIGN dataset) and the 259 question set (excluding instruction-EHR pairs that triggered Azure’s content filter) are considered.

Model	Context	WR ↑ 303	WR ↑ 259
GPT-3.5	8192	<b>0.61</b>	0.60
LLaMA 2 13B Chat	4096	<b>0.61</b>	0.60
GPT-4 (MR)	32768 <sup>†</sup>	0.58	0.57
GPT-4	32768	0.56	<b>0.64</b>
GPT-4	2048*	0.56	0.56
LLaMA 2 7B Chat	4096	0.51	0.50
Vicuña-7B	2048	0.43	0.42
Vicuña-13B	2048	0.40	0.39
MPT-7B-instruct	2048	0.23	0.22

### Q.1 Evaluating LLMs with COMET

In addition to our original 6 models (evaluated by clinicians), we introduced three LLMs — LLaMA 2 13B Chat, LLaMA 2 7B Chat, and GPT-3.5-turbo (snapshot: Azure 2023-03-15-preview) — and evaluated their performance using COMET.

As indicated by Supplementary Table S16 and Supplementary Figure S8, initial analysis of the win rates for all models across the total 303 questions suggests that GPT-3.5 and LLaMA 2 exhibit the best performance. However, upon closer examination (excluding “error” messages triggered by GPT-4’s content filter, leaving 259 instructions), GPT-4 (32k) retains its ranking as the top-performing model. Additionally, win rate also tends to be higher in models with larger context windows. These results largely align with our previous analysis. Thus, these findings support the use of COMET as a proxy to measure technical progress on MEDALIGN when access to organizational infrastructure and clinician labeling is unavailable.

### Q.2 COMET Win Rate vs Human Win Rate

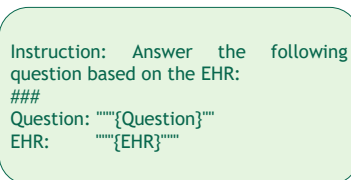
We assessed the absolute error between the win rates generated from human evaluations (using rankings) and COMET scores. While certain models exhibit a larger absolute error, this discrepancy consistently remains below 10%. Moreover, the average absolute error in win rates across all LLMs is just 4.33% (see Table S17).

Table S17: Human-evaluated Win Rate vs. COMET-evaluated Win Rate using 303 instructions.

Model	Context	WR ↑ Human	WR ↑ COMET	\Delta WR
GPT-4 (MR)	32768 <sup>†</sup>	0.66	<b>0.63</b>	0.03
GPT-4	32768	<b>0.68</b>	0.6	0.08
GPT-4	2048*	0.6	0.6	0
Vicuña-13B	2048	0.40	0.45	0.05
Vicuña-7B	2048	0.40	0.48	0.08
MPT-7B	2048	0.27	0.25	0.02

## Appendix R Experiments Prompt

All experiments conducted in this study, including the generation of responses post-clinical evaluation (LLaMA 2 and GPT3.5), were performed employing the same prompt (see Figure S9). Subsequent de-identified generated outputs are presented in Table S3



```
Instruction: Answer the following
question based on the EHR:
###
Question: ""{Question}""
EHR: ""{EHR}""
```

Figure S9: Standard prompt used to generate LLM responses to questions and instructions grounded on matched EHR.

## Appendix S Compute Environment

Experiments are performed in a local on-prem university compute environment using 24 Intel Xeon 2.70GHz CPU cores, 8 Nvidia V100 GPUs, 4 Nvidia A100 GPUs, and 1.48 TB of RAM. All compute environments supported HIPAA-compliant data protocols.