

<https://doi.org/10.1038/s41746-025-01589-z>

# Expert of Experts Verification and Alignment (EVAL) Framework for Large Language Models Safety in Gastroenterology



Mauro Giuffrè<sup>1,12</sup>, Kisung You<sup>2,12</sup>, Ziteng Pang<sup>3,12</sup>, Simone Kresevic<sup>1,4</sup>, Sunny Chung<sup>1</sup>, Ryan Chen<sup>3</sup>, Youngmin Ko<sup>3</sup>, Colleen Chan<sup>5</sup>, Theo Saarinen<sup>6</sup>, Milos Ajcevic<sup>4</sup>, Lory S. Crocè<sup>7</sup>, Guadalupe Garcia-Tsao<sup>1</sup>, Ian Gralnek<sup>8</sup>, Joseph J. Y. Sung<sup>9</sup>, Alan Barkun<sup>10</sup>, Loren Laine<sup>1</sup>, Jasjeet Sekhon<sup>5</sup>, Bradly Stadie<sup>3,13</sup> & Dennis L. Shung<sup>1,11,13</sup> ✉

Large language models generate plausible text responses to medical questions, but inaccurate responses pose significant risks in medical decision-making. Grading LLM outputs to determine the best model or answer is time-consuming and impractical in clinical settings; therefore, we introduce EVAL (Expert-of-Experts Verification and Alignment) to streamline this process and enhance LLM safety for upper gastrointestinal bleeding (UGIB). We evaluated OpenAI's GPT-3.5/4/4o/o1-preview, Anthropic's Claude-3-Opus, Meta's LLaMA-2 (7B/13B/70B), and Mistral AI's Mixtral (7B) across 27 configurations, including zero-shot baseline, retrieval-augmented generation, and supervised fine-tuning. EVAL uses similarity-based ranking and a reward model trained on human-graded responses for rejection sampling. Among the employed similarity metrics, Fine-Tuned CoBERT achieved the highest alignment with human performance across three separate datasets ( $\rho = 0.81\text{--}0.91$ ). The reward model replicated human grading with 87.9% of cases across temperature settings and significantly improved accuracy through rejection sampling by 8.36% overall. EVAL offers scalable potential to assess accuracy for high-stakes medical decision-making.

Large language models (LLMs) have demonstrated a remarkable ability to generate relevant text in response to clinical questions<sup>1,2</sup>. However, the inherent variability and occasional inaccuracy of these models can limit their application in high-stakes situations such as clinical decision-making in patient care<sup>3–8</sup>. The issue of Artificial Intelligence (AI) safety becomes a critical concern when LLMs are used for medical advice<sup>9</sup>, as preliminary studies have shown that these models may generate inaccurate recommendations for patients and healthcare providers in gastroenterology and hepatology<sup>10</sup>. Although techniques

such as few-shot prompting<sup>11</sup>, retrieval-augmented generation<sup>12</sup>, and supervised fine-tuning have been employed to improve model performance, the criteria to evaluate performance metrics (e.g., accuracy) remain inconsistent across studies. Moreover, verifying model performances is time and resource-intensive, requiring extensive manual review from medical experts<sup>13</sup>. Ensuring AI safety in LLMs for medical advice requires a clear definition of appropriate performance metrics, which make it difficult to evaluate LLMs and establish an appropriate regulatory framework<sup>14</sup>.

<sup>1</sup>Section of Digestive Diseases, Department of Medicine, Yale School of Medicine, New Haven, USA. <sup>2</sup>Department of Mathematics, Baruch College, The City University of New York, New York, USA. <sup>3</sup>Department of Statistics and Data Science, Northwestern University, Chicago, USA. <sup>4</sup>Department of Engineering and Architecture, University of Trieste, Trieste, Italy. <sup>5</sup>Department of Statistics and Data Science, Yale University, New Haven, USA. <sup>6</sup>Department of Statistics, University of California Berkeley, Berkeley, USA. <sup>7</sup>Department of Medical, Surgical, and Health Sciences, University of Trieste, Trieste, Italy. <sup>8</sup>Rappaport Faculty of Medicine Technion Israel Institute of Technology, Haifa, Israel. <sup>9</sup>Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore, Singapore. <sup>10</sup>Division of Gastroenterology, McGill University, Montreal, Canada. <sup>11</sup>Department of Biomedical Informatics and Data Science, Yale School of Medicine, New Haven, USA. <sup>12</sup>These authors contributed equally: Mauro Giuffrè, Kisung You, Ziteng Pang. <sup>13</sup>These authors jointly supervised this work: Bradly Stadie and Dennis L. Shung. ✉e-mail: [dennis.shung@yale.edu](mailto:dennis.shung@yale.edu)

In the context of generative AI safety, establishing a reliable ground truth is essential. Evidence-based medicine (EBM) is the prevailing paradigm for clinical practice to define a consensus for medical practice by emphasizing systematic literature reviews and formal evidence-based decision rules to inform clinical decision-making<sup>15</sup>. This approach has been consolidated in systematic reviews, meta-analyses, and evidence synthesis in clinical practice with an estimated number of over 2,700 published clinical guidelines<sup>16</sup>. Within this framework, the accuracy of generative AI systems can be defined as the degree to which its outputs align with the recommendations outlined in established clinical practice guidelines and disease-specific protocols.

Existing studies involving LLMs application in clinical practice seek to pool the responses of board-certified clinical practitioners to crowd-source the appropriate response to clinical questions. This process is time consuming, heterogeneous across practitioners, and may not reflect the best specialized knowledge for evidence-based management of diseases. To overcome these limitations, we define our reference standard using free-text responses from lead or senior guideline authors - the so-called “expert-of-experts”. These responses provide the elusive “golden labels” that can be used to enable the automated ranking of various LLM configurations and to facilitate the identification of those most aligned with expert-level guidance.

We propose expert-of-experts verification and alignment (EVAL) framework, which comprises two complementary tasks operating at different levels of evaluation. The first task provides a scalable solution at the model level, using unsupervised embeddings to automatically evaluate and rank different LLM configurations based on how closely their responses align with expert-generated answers. The unsupervised embedding approach works by converting both LLM outputs and expert answers into mathematical representations (vectors), allowing semantic similarity comparison without requiring manual labeling or supervision. These vector representations capture the meaning of text in high-dimensional space, where distance metrics quantify the degree of alignment between LLM and expert responses. The second task operates at the individual answer level, using a reward model trained on expert-graded LLM responses to score and filter out inaccurate outputs automatically across multiple temperature thresholds, thus accounting for different levels of randomness and diversity. This two-level approach allows us to both identify the most reliable LLM configurations but also ensure that individual outputs meet clinical quality standards.

To illustrate the utility of EVAL, we applied the framework to the management of upper gastrointestinal bleeding (UGIB), a common and costly condition. UGIB affects up to 116 per 100,000<sup>17</sup> individuals and carries a mortality rate of up to 11%<sup>18</sup>. Robust national and international clinical guidelines provide evidence-based recommendations for management across the pre-endoscopic, endoscopic, and post-endoscopic phases of clinical care<sup>19–24</sup>. Adherence to guideline-based recommendations is variable and low despite efforts to knowledge dissemination, with an estimated adherence rate that ranged from 14.3% to 95.7% across 20 guideline-recommended measures and only 30% of practitioners ever having used a guideline-recommended risk stratification score<sup>25–27</sup>. Adherence to UGIB guideline recommendations is poor in clinical practice, and there is potential for LLMs to serve as a clinical decision support to improve guideline implementation. The implementation of LLMs as a clinical decision support tool<sup>27,28</sup>. Our study benchmarks the EVAL framework across three datasets: 13 expert-generated questions on UGIB, 40 multiple-choice questions (MCQs) derived from the self-assessments test of the American College of Gastroenterology (ACG), and 117 real-world questions asked by physician trainees to LLMs in simulation scenarios on UGIB diagnosis and management.

In summary, The EVAL framework aims to provide a scalable solution to enhance AI safety for provider-facing LLMs by simultaneously identifying robust model configurations and verifying that individual responses align with established, guideline-based recommendations. This dual approach ultimately aims to improve the quality and safety of LLM-enhanced clinical decision support.

## Results

### Model Ranking by Similarity Metrics

In terms of model ranking by similarity metrics (Table 1), Claude-3-Opus in the baseline configuration achieved the best performance in both Term Frequency-Inverse Document Frequency (TF-IDF) ( $0.252 \pm 0.002$ ) and Sentence Transformers ( $0.579 \pm 0.003$ ), while SFT-GPT-4o demonstrated the highest similarity using the Fine-Tuned Contextualized Late Interaction over BERT (ColBERT) scoring ( $0.699 \pm 0.012$ ). With ranking by TF-IDF metric, Claude-3-Opus baseline showed statistically significant differences ( $p < 0.01$ ) compared to all other models and configurations, with only RAG-GPT-o1 showing a less stringent statistical significance ( $p < 0.05$ ). For the Sentence Transformers metric, Claude-3-Opus baseline showed no statistically significant differences when compared to its RAG configuration ( $0.578 \pm 0.003$ ) and SFT-GPT-4o ( $0.554 \pm 0.003$ ), while all other model configurations demonstrated statistically significant differences ( $p < 0.01$ ). The Fine-Tuned ColBERT evaluation revealed no statistically significant differences between the best model (i.e., SFT-GPT-4o) and several highly similar configurations such as baseline GPT-o1 ( $0.683 \pm 0.009$ ), baseline GPT-4o ( $0.669 \pm 0.011$ ), RAG-Claude-3-Opus ( $0.680 \pm 0.006$ ), RAG-GPT-4 ( $0.679 \pm 0.006$ ), RAG-GPT-o1 ( $0.687 \pm 0.004$ ), SFT-GPT 3.5 ( $0.673 \pm 0.009$ ), SFT-GPT-4 ( $0.691 \pm 0.014$ ), RAG-SFT-GPT4 ( $0.683 \pm 0.010$ ) and RAG-SFT-GPT4o ( $0.681 \pm 0.015$ ).

It is important to note that similarity metrics, particularly Fine-Tuned ColBERT, are primarily designed as ranking tools, with their raw output values mainly indicating relative performance rather than absolute scores. Given Fine-Tuned ColBERT’s superior correlation with human evaluation (as demonstrated later in the manuscript) compared to TF-IDF and Sentence Transformers, we focused our visualization efforts on ColBERT scores. To enhance visualization clarity while preserving the ranking information, we applied a logit transformation to the Fine-Tuned ColBERT scores in Fig. 1, as this transformation maintains the monotonic relationship between scores while providing better visual differentiation between high-performing models.

### Model ranking by human grading and multiple-choice questions

Regarding human evaluation metrics, SFT-GPT-4o achieved the highest performance in both expert-generated questions (88.5%) and ACG-MCQ evaluation (87.5%), while RAG-GPT-o1 demonstrated superior performance in real-world questions (88.0%) as reported in Table 1 and depicted in Fig. 2. For expert-generated questions, no statistically significant differences were observed between the accuracy of the best model and RAG-GPT-4 (84.6%), RAG-GPT-4o (87.7%), RAG-Claude-3-Opus (86.2%), SFT-GPT-3.5 (80.8%), SFT-GPT-4 (84.6%), RAG-SFT-GPT-4 (81.5%), and RAG-SFT-GT4o (83.1%). At the same time the best model for expert-generated questions showed statistically significant higher accuracy when compared to RAG-GPT-o1 (76.9%,  $p < 0.05$ ) and with all the other model configurations demonstrated statistically significant differences ( $p < 0.01$ ). Similarly, in ACG-MCQ evaluation, no statistically significant differences were observed between the accuracy of the best model and baseline GPT4o (72.5%), RAG-GPT-4 (80%), RAG-GPT-4o (82.5%), RAG-Claude-3-Opus (75%), RAG-GPT-o1 (77.5%), SFT-GPT-3.5 (72.5%), SFT-GPT-4 (85%), RAG-SFT-GPT-4 (80%), and RAG-SFT-GT4o (82.5%). At the same time the best model for ACG-MCQs showed statistically significant higher accuracy when compared to baseline Claude-3-Opus (65%,  $p < 0.05$ ), baseline GPT-o1 (60%,  $p < 0.05$ ), and with all the other model configurations demonstrated statistically significant differences ( $p < 0.01$ ). For real-world questions, no statistically significant differences were observed between the accuracy of the best model and RAG-GPT-4 (80.3%), RAG-GPT-4o (82.1%), SFT-GPT-4 (82.9%), SFT-GPT-4o (84.6%), RAG-SFT-GPT-4 (81.2%), and RAG-SFT-GPT4o (82.1%). The best model for real-world questions showed statistically significant higher accuracy when compared to RAG-Claude-3-Opus (76.9%,  $p < 0.05$ ), with all the other model configurations demonstrated statistically significant differences ( $p < 0.01$ ).

**Table 1 | Model Ranking Comparison across similarity-based metrics, human grading, and performance of multiple-choice questions (MCQs) dataset**

Model Configuration	Ranking by Similarity Metrics			Ranking by Human Grading and Multiple-Choice Questions		
	TF-IDF Average ( $\pm$ SD)	Sentence Transformers Average ( $\pm$ SD)	Fine-Tuned Colbert Score Average ( $\pm$ SD)	Expert-Generated Questions <i>N</i> (%)	ACG-MCQs Performance <i>N</i> (%)	Real-World Questions <i>N</i> (%)
<b>Baseline configuration</b>						
Llama-2-7B	0.210 (0.002)**	0.514 (0.003)**	0.603 (0.010)**	35 (26.9%)**	8 (20%)**	39 (33.3%)**
Llama-2-13B	0.210 (0.002)**	0.525 (0.002)**	0.633 (0.013)**	49 (37.7%)**	12 (30%)**	41 (35.0%)**
Llama-2-70B	0.228 (0.002)**	0.547 (0.004)**	0.633 (0.007)**	65 (50.0%)**	13 (32.5%)**	45 (38.5%)**
Mistral-7B	0.199 (0.002)**	0.543 (0.003)**	0.634 (0.008)**	66 (50.8%)**	16 (40%)**	53 (45.3%)**
Claude-3-Opus	0.252 (0.002) <sup>BM</sup>	0.579 (0.003) <sup>BM</sup>	0.672 (0.007)*	95 (73.1%)**	26 (65%)*	80 (68.4%)**
GPT-3.5	0.199 (0.001)**	0.499 (0.001)**	0.639 (0.009)**	66 (50.8%)**	21 (52.5%)**	73 (62.4%)**
GPT-4	0.192 (0.001)**	0.499 (0.001)**	0.642 (0.007)**	82 (63.1%)**	22 (55%)**	82 (70.1%)**
GPT-4o	0.242 (0.002)**	0.559 (0.001)**	0.669 (0.011) <sup>NS</sup>	90 (69.2%)**	29 (72.5%) <sup>NS</sup>	84 (71.8%)**
GPT-o1	0.221 (0.004)**	0.555 (0.005)**	0.683 (0.009) <sup>NS</sup>	95 (73.1%)**	24 (60%)*	87 (74.4%)*
<b>Retrieval augmented generation configuration</b>						
Llama-2-7B	0.223 (0.002)**	0.555 (0.003)**	0.648 (0.009)**	80 (61.5%)**	11 (27.5%)**	45 (38.5%)**
Llama-2-13B	0.218 (0.002)**	0.540 (0.003)**	0.662 (0.011)*	91 (70.1%)**	23 (57.5%)**	45 (38.5%)**
Llama-2-70B	0.232 (0.001)**	0.565 (0.003)**	0.662 (0.008)**	88 (67.7%)**	22 (55%)**	46 (39.3%)**
Mistral-7B	0.223 (0.001)**	0.544 (0.002)**	0.660 (0.008)**	88 (67.7%)**	23 (57.5%)**	52 (44.4%)**
Claude-3-Opus	0.243 (0.003)**	0.578 (0.003) <sup>NS</sup>	0.680 (0.006) <sup>NS</sup>	112 (86.2%) <sup>NS</sup>	30 (75%) <sup>NS</sup>	90 (76.9%)*
GPT-3.5	0.199 (0.002)**	0.499 (0.001)**	0.653 (0.007)**	83 (63.8%)**	18 (45%)**	61 (51.3%)**
GPT-4	0.225 (0.001)**	0.559 (0.001)**	0.679 (0.006) <sup>NS</sup>	110 (84.6%) <sup>NS</sup>	32 (80%) <sup>NS</sup>	94 (80.3%) <sup>NS</sup>
GPT-4o	0.234 (0.002)**	0.571 (0.002)**	0.670 (0.006)*	114 (87.7%) <sup>NS</sup>	33 (82.5%) <sup>NS</sup>	96 (82.1%) <sup>NS</sup>
GPT-o1	0.239 (0.004)*	0.563 (0.004)**	0.687 (0.004) <sup>NS</sup>	100 (76.9%)*	31 (77.5%) <sup>NS</sup>	103 (88.0%) <sup>BM</sup>
<b>Supervised fine-tuning configuration</b>						
Llama-2-7B	0.216 (0.001)**	0.525 (0.002)**	0.630 (0.011)**	27 (20.8%)**	18 (45%)**	28 (23.9%)**
Llama-2-13B	0.223 (0.001)**	0.529 (0.002)**	0.646 (0.016)**	43 (33.1%)**	13 (32.5%)**	31 (26.5%)**
Llama-2-70B	0.226 (0.002)**	0.545 (0.001)**	0.649 (0.007)**	79 (60.8%)**	16 (40%)**	85 (72.6%)**
Mistral-7B	0.197 (0.003)**	0.527 (0.002)**	0.634 (0.008)*	66 (50.8%)**	17 (42.5%)**	37 (31.6%)**
GPT-3.5	0.223 (0.002)**	0.559 (0.002)**	0.673 (0.009) <sup>NS</sup>	105 (80.8%) <sup>NS</sup>	29 (72.5%) <sup>NS</sup>	79 (59.8%)**
GPT-4	0.215 (0.002)**	0.540 (0.003)**	0.691 (0.014) <sup>NS</sup>	110 (84.6%) <sup>NS</sup>	34 (85%) <sup>NS</sup>	97 (82.9%) <sup>NS</sup>
GPT-4o	0.219 (0.003)**	0.554 (0.003) <sup>NS</sup>	0.699 (0.012) <sup>BM</sup>	115 (88.5%) <sup>BM</sup>	35 (87.5%) <sup>BM</sup>	99 (84.6%) <sup>NS</sup>
<b>Retrieval augmented generation and supervised fine-tuning configuration</b>						
GPT-4	0.217 (0.003)**	0.538 (0.006)**	0.683 (0.010) <sup>NS</sup>	106 (81.5%) <sup>NS</sup>	32 (80%) <sup>NS</sup>	95 (81.2%) <sup>NS</sup>
GPT-4o	0.213 (0.003)**	0.535 (0.004)**	0.681 (0.015) <sup>NS</sup>	108 (83.1%) <sup>NS</sup>	33 (82.5%) <sup>NS</sup>	96 (82.1%) <sup>NS</sup>

This table compares the performance of different LLM configurations using three evaluation approaches: automated similarity metrics (TF-IDF, Sentence Transformers, and ColBERT scores), human expert validation (expert-generated and real-world questions), and standardized testing (ACG-MCQs). Models are evaluated in four configurations (Baseline, RAG, SFT, and Combined RAG-SFT), with statistical significance noted as BM (Best Model), NS (Not Significant from best), \* $p < 0.05$ , \*\* $p < 0.01$ . Higher scores indicate better performance across all metrics. Abbreviations: LLM Large Language Model, RAG Retrieval Augmented Generation, SFT Supervised Fine-Tuning, ACG-MCQs American College of Gastroenterology Multiple Choice Questions, TF-IDF Term Frequency-Inverse Document Frequency, ColBERT Contextualized Late Interaction over BERT.

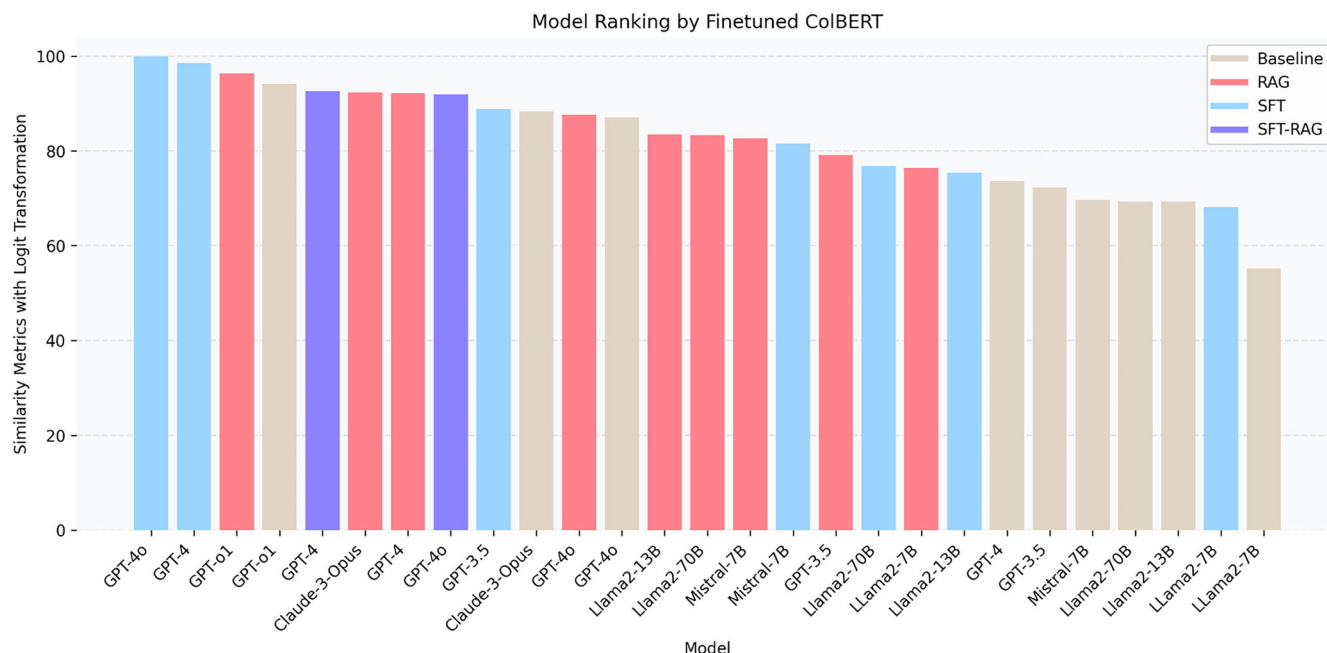
### Alignment between similarity metrics and human performance

To assess which similarity metrics best reflected model performance across our three validation datasets, we explored the correlation between their scores and the accuracy by human grading and performance on ACG-MCQs using Spearman correlation coefficients. The Fine-Tuned ColBERT metric demonstrated the strongest correlation with human evaluation across all three datasets, showing high correlation coefficients with expert-generated questions ( $\rho = 0.91$ ,  $p < 0.001$ ), ACG-MCQs performance ( $\rho = 0.86$ ,  $p < 0.001$ ), and real-world questions ( $\rho = 0.81$ ,  $p < 0.001$ ). Sentence Transformers showed moderate correlations with expert-generated questions ( $\rho = 0.59$ ,  $p < 0.01$ ), ACG-MCQ performance ( $\rho = 0.47$ ,  $p < 0.05$ ), and real-world questions ( $\rho = 0.44$ ,  $p < 0.05$ ). TF-IDF demonstrated the weakest correlation, with a marginally significant correlation only with expert-generated questions ( $\rho = 0.38$ ,  $p < 0.05$ ), while correlations with ACG-MCQ performance ( $\rho = 0.30$ ,

$p = 0.13$ ) and real-world questions ( $\rho = 0.28$ ,  $p = 0.16$ ) were not statistically significant.

### Evaluation of reward model alignment to human-grading

The human-grading evaluation accuracy for each model used in the reward model training and validation across multiple temperature threshold is reported in Supplementary Fig. 1. The reward model produced a true label (i.e., the same grade produced by human graders) in 87.9% of cases across all temperature values for RAG-GPT-4. In the two regimens where the LLM output quality is easy to distinguish (i.e., lower temperatures with more deterministic outcomes vs. higher temperatures with less deterministic outcomes) the reward model produced true labels in 90.0% (positive regime, temperature  $< 1.2$ ) and 99.2% (negative regime, temperature  $> 1.6$ ) of cases (Fig. 3). In the mixed regime (i.e., temperature values between 1.2 and 1.6), where the distinction between good and bad LLM-generated answers may



**Fig. 1 | Model performance ranking based on Fine-Tuned ColBERT similarity scores.** The figure shows the ranking of different LLM configurations based on their similarity to expert-generated responses, as measured by Fine-Tuned ColBERT scores after logit transformation. Models are grouped by configuration type

(Baseline, RAG, SFT, and SFT-RAG). The logit transformation was applied to enhance visualization while maintaining the relative ranking. Abbreviations: RAG Retrieval Augmented Generation, SFT Supervised Fine-Tuning, ColBERT Contextualized Late Interaction over BERT.

result in less obvious and the classification task results less performant, the reward model produced true labels in 76.2% of cases. For temperatures  $< 1.2$  (positive regime) the reward model provides true labels for 90% of correct answers and 67% of inaccurate answers. For temperatures  $> 1.6$  (negative regime), the reward model provides true labels for 94.1% of correct answers and 100% of inaccurate answers. In the mixed regime (temperature values between 1.2 and 1.6), the reward model produced true labels for 68.8% of correct answers and 97.1% of inaccurate answers.

In the external validation using the SFT-GPT-4o model, the reward model produced a true label in 81.8% of cases across all temperature values, with slightly different performance in the positive regime when compared to the internal validation. In particular, in the positive regime (temperature < 1.2), it achieved 72.6% accuracy for correct answers and 89.3% for inaccurate answers. In the negative regime (temperature > 1.6), it showed a similarly strong performance with 90.9% accuracy for correct answers and 99.6% for inaccurate answers. However, in the mixed regime (temperature values between 1.2 and 1.6), true labels were achieved in 75.3% of correct answers and 96.4% of inaccurate answers.

We performed a sensitivity analysis to detect the different levels of alignment for RAG-GPT-4 and SFT-GPT-4o across all temperature thresholds and alignment with human-grading on real-world questions for each model are reported in Supplementary Table 1 and Supplementary Table 2 respectively.

### Rejection sampling across multiple temperature thresholds

Rejection sampling was employed to enhance the accuracy of LLM responses by leveraging the alignment observed in the reward model analysis. To evaluate its effectiveness, we compared human-graded accuracy with and without rejection sampling, using  $K = 5$  candidate responses for each query. Across all regimes, including a large portion of temperature that LLM model already has a high accuracy, rejection sampling improves the overall accuracy by 9.39% in answers produced by RAG-GPT-4 and 8.36% in answers produced by SFT-GPT-4o (Table 2). The improvement in accuracy produced by the rejection sampling of the positive regime was 1.14% for answers produced by RAG-GPT-4 and 1.12% for answers produced by SFT-GPT-4o. In the mixed regime (temperature 1.2–1.6), where

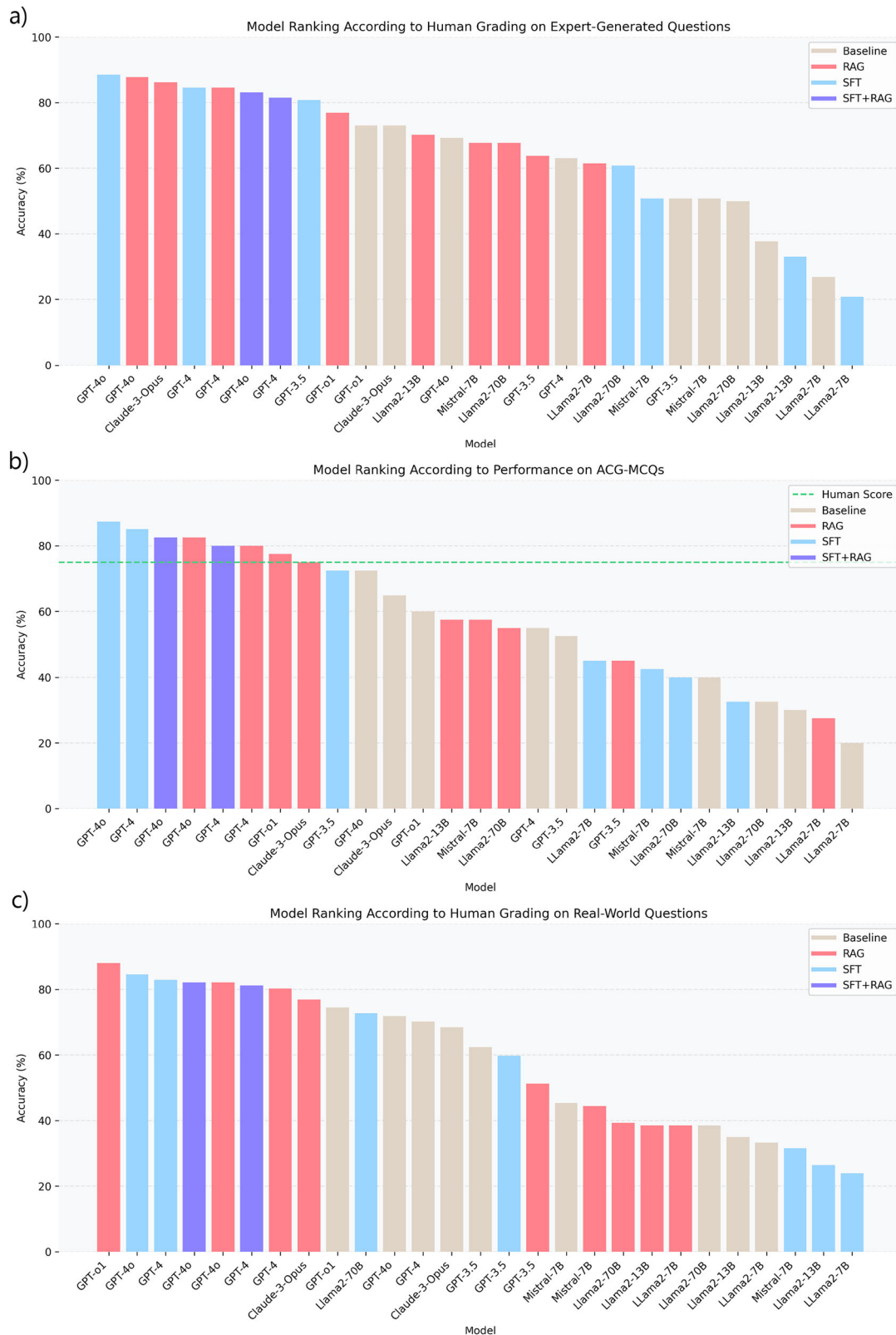
classification is more challenging, rejection sampling provides a significant improvement of 7.65% for RAG-GPT-4 (increasing accuracy from 51.0%–54.9%) and of 23.60% for SFT-GPT-4o (increasing accuracy from 64.4%–79.6%). In the negative regime (temperature > 1.6), rejection sampling drastically improves accuracy by 98.35% (increasing accuracy from 12.1% to 24.0%) in answers generated by RAG-GPT-4 and by 121.43% (increasing accuracy from 4.2% to 9.3%) in answers generated by SFT-GPT-4o. These findings highlight the ability of rejection sampling to improve performance in more difficult regimes, particularly at higher temperatures where the model’s baseline accuracy is low.

## Discussion

We present EVAL, a novel framework that leverages expert-of-expert free text responses to identify the best-performing LLM configurations and a trained reward model to identify high-quality responses from several LLM configurations. We demonstrate benchmark performance for accuracy across an expert-generated dataset, a multiple-choice question dataset, and a real-world question dataset focused on the management of UGIB.

AI safety in deploying LLMs in clinical medicine can encompass many categories, but for clinical practice impacts most practically the task of diagnosis using published clinical cases<sup>29-31</sup> and the task of management as measured by performance on multiple-choice questions featured in clinical exams<sup>32,33</sup>. LLM configurations used to retrieve information from clinical guidelines for clinical decision support have focused on simple retrieval<sup>34-36</sup>, but strategies to optimize the use of LLMs for the task of clinical decision support are important in mitigating the risk of using these systems in clinical care. Our approach is rooted in the paradigm of evidence-based medicine and can be used across multiple domains to improve the performance of LLMs when deployed for clinical decision support in high-risk, time-constrained medical settings.

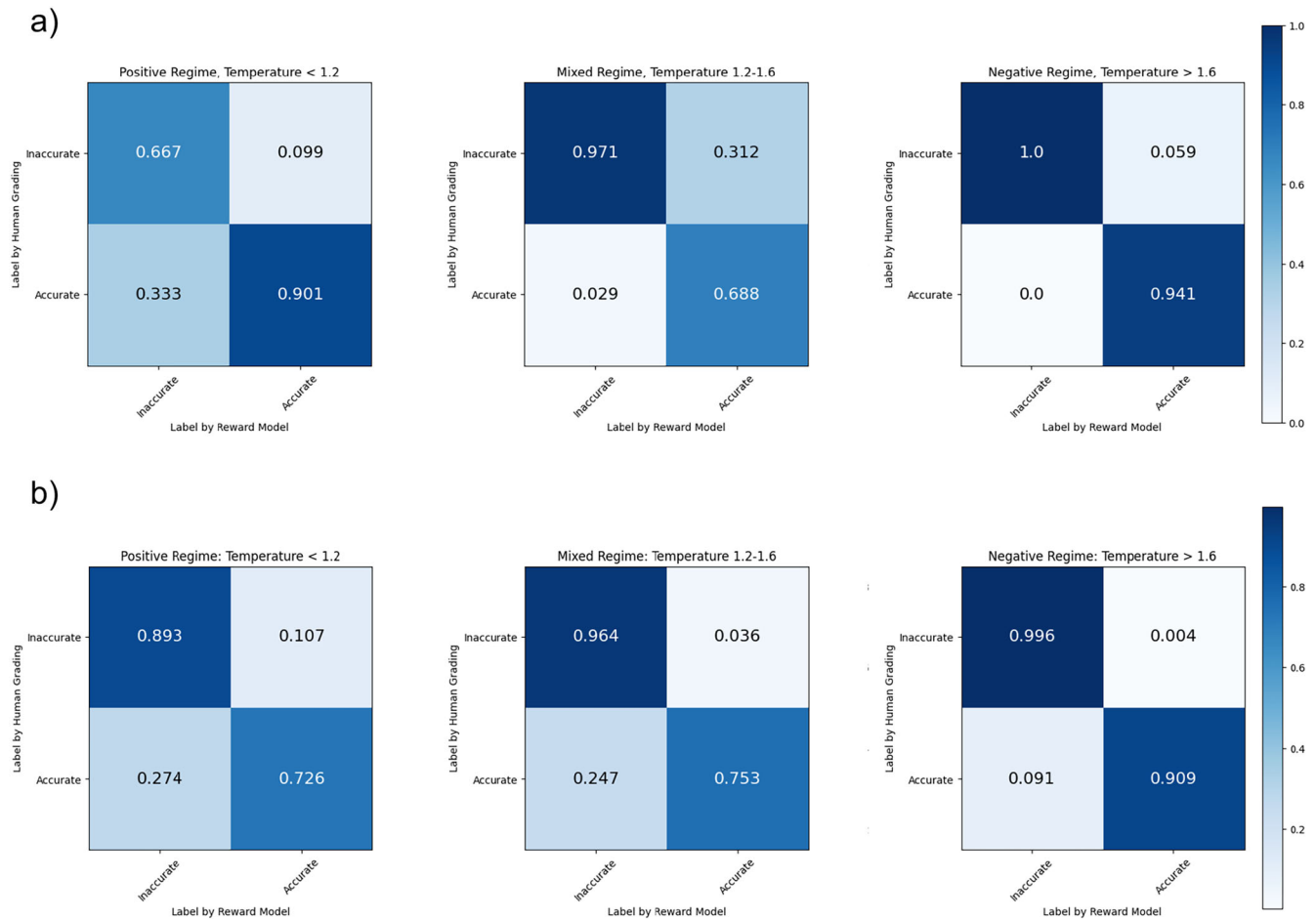
Our study is the first to use unsupervised embeddings and reward models to select the best performing LLM configurations at the model and at the answer level. Unsupervised similarity metrics based on a high-quality comparator (i.e., expert-of-experts golden labels) using the embedding representation, along with pre-trained reward models to screen for high-quality LLM responses, demonstrate potential as a less resource-intensive



**Fig. 2 | Model ranking according to human grading and performance on ACG-MCQs.** The figure presents model performance rankings across three different human evaluation approaches: **a** Expert-generated questions; **b** ACG-MCQs performance; and **c** Real-world questions. Models are grouped by configuration type (Baseline, RAG, SFT, and SFT + RAG), with advanced GPT models consistently

performing well across all evaluation metrics. Notably, enhanced configurations (RAG, SFT, SFT + RAG) generally outperformed baseline models. Abbreviations: ACG-MCQs American College of Gastroenterology Multiple Choice Questions, RAG Retrieval Augmented Generation, SFT Supervised Fine-Tuning.





**Fig. 3 | Confusion matrix comparing labels by reward model and human grading.** The two confusion matrices compare labels according to human grading vs. labels provided by the reward model in the three regimes (i.e., temperature ranges). **a** Internal validation of the reward model with answers generated by the RAG-GPT-4 configuration; **b** External validation of the reward model with answers generated

by SFT-GPT-4o, which was the best model selected according to human grading and embedding similarity metrics. The reward model was able to detect most of the inaccurate answers in the context of higher temperature settings. Abbreviations: RAG Retrieval Augmented Generation; SFT Supervised Fine-Tuning.

**Table 2 | Rejection Sampling for automated grading**

Settings	Overall Temperature 0–2	Positive Regime Temperature <1.2	Mixed Regime Temperature 1.2–1.6	Negative Regime Temperature >1.6
<b>RAG-GPT-4 (Internal validation)</b>				
Baseline	0.511	0.880	0.510	0.121
With Rejection Sampling	0.559	0.890	0.549	0.240
Improvement (%)	9.39%	1.14%	7.65%	98.35%
<b>SFT-GPT-4o (External validation)</b>				
Baseline	0.529	0.893	0.650	0.043
With Rejection Sampling	0.598	0.903	0.796	0.093
Improvement (%)	8.36%	1.12%	23.60%	121.43%

This table illustrates the impact of implementing rejection sampling (with  $K = 5$ ) on the accuracy of the reward model for automated grading across different temperature regimes.

approach to screen LLM configurations. We believe that this approach has value for healthcare systems, clinical providers, and patient advocacy groups to choose wisely in an increasingly crowded space of different LLMs with various customizations. When medical entities (corporate, hospital, or individual teams) need to choose between different model configurations, a scalable method that does not require manual human-grading can help to save time and mitigate risks when thinking through the implementation of LLMs for clinical decision-making. We tested 27 different model

configurations across three different datasets. Our findings highlight that RAG, SFT, and combined approaches can significantly improve performance over baseline LLM configurations, which is consistent with the results of other studies testing different LLM configurations in healthcare applications<sup>12,34–36</sup>. However, we note that there was similar accuracy with either RAG or SFT among multiple proprietary models. Interestingly, while one might expect that combining RAG and SFT would yield superior performance compared to either approach alone, our results indicate this was

not a consistent pattern. This observation may be explained by information redundancy—when the domain-specific knowledge provided through RAG overlaps substantially with the knowledge already encoded in the fine-tuned model parameters through SFT, the marginal benefit of combining both approaches diminishes. Additionally, SFT induces parametric changes that can alter the model's interpretation mechanisms for specialized medical text, potentially creating interference when subsequently processing retrieved external context from RAG. This phenomenon resembles catastrophic forgetting<sup>37</sup> in continual learning scenarios, where fine-tuning on one objective can degrade performance on previously learned tasks. The optimal configuration appears to depend on which approach better aligns with the specific knowledge representation requirements of a given clinical domain and question type. In addition to identifying the LLM configuration with the highest quality responses, the reward model may be useful in mitigating risk across LLM hyperparameter settings such as temperature. Higher temperatures could be beneficial for reasoning over complex clinical cases<sup>38</sup>, but also lead to higher risk of more hallucinations, potentially deviating from guideline recommendations in harmful ways. Our preliminary findings suggest that a reward model could be used to reject inaccurate responses at higher temperatures ( $>1.2$ ), leading to a partial rescue for clinical accuracy.

Finally, we present a set of UGIB databases with labels and a benchmark performance of our approach that can be used to test other approaches to evaluating LLM configurations for accuracy in the high-stakes realm of clinical decision support for evidence-based medical practice. We believe this provides a valuable and novel contribution towards the field of LLM safety testing in medicine.

The real-world efficacy of EVAL is demonstrated with the improvement in accuracy over the baseline model in a real-world question dataset generated by clinical providers within medical simulation for the management of acute upper gastrointestinal bleeding. No other study, to our knowledge, has evaluated available LLM configurations on real clinician questions in the context of clinical decision making. EVAL also has the potential to automate comparisons of LLMs and identify the optimal configurations for accuracy. EVAL uses an unsupervised embedding to measure similarity to expert-of-expert free text responses confirmed with multiple-choice question dataset, and then leverages a trained reward model to provide automated estimates of LLM output accuracy. The trained reward model can also be used to identify optimal temperature thresholds and improve the performance at other temperature thresholds with rejection sampling. Our results in the real-world question dataset suggests that despite training reward models on high-quality data, a gap in accuracy persists between reward models and human-graded accuracy.

Limitations of our approach include the following: we require a pooling of free text responses from expert physicians and existence of high-quality clinical guidelines, our approach may not be able to fully account for heterogeneity across different guidelines, and the real-world questions were derived from real physicians in simulation workflows rather than actual clinical workflow. We present a narrow use case to showcase our framework, focusing only on the management of patients with UGIB. Nonetheless, our approach is flexible and can be readily applied to other conditions that have both expert responses and associated clinical guideline text. Another consideration is our use of United States clinical guidelines for training and European/Asia-Pacific guidelines for testing, aligned with our expert panel's geographical distribution. While this approach validates cross-system generalizability, it may introduce subtle biases, though UGIB management principles remain largely consistent across international guidelines. Regarding broader generalizability, we cannot definitively claim the reward model would maintain performance on entirely different clinical questions. Different medical conditions may present unique challenges: less standardized guidelines, more complex decision trees, or nuanced clinical judgments that are harder to evaluate systematically. Additionally, the real-world questions were generated by providers within medical simulation on standardized patient cases and only approximate live clinical care. While medical simulation is well-established as an environment for testing medical

technologies, particularly those with potential risks to patient safety, real-world questions when deployed in clinical practice are the real test for LLM performance. We do not directly capture the feedback of clinical provider users to the LLM output, which may inform how the output may influence their clinical decision within the clinical scenario. Future studies should consider mechanisms to collect provider feedback so that their expressed preference for LLM responses and quantify downstream impact of how they were used in their clinical decision-making.

Our findings suggest that AI safety can be optimized within an evidence-based medicine framework, where clinical guidelines and expert guidance can be codified to evaluate LLM outputs and reject inaccuracies. Further work to scale AI safety solutions across other domains of medicine is necessary to ensure that answers to high-stakes medical issues are factually accurate, reliable, and reflect the current standard of care.

## Methods

### Large language model configurations

We tested the following large language model architectures based on availability for clinical use: GPT-3.5-Turbo, GPT-4-Turbo, GPT-4o, GPT-o1-preview, Claude-3-Opus, LLaMA-2-7B, LLaMA-2-13B, LLaMA-2-70B, and Mistral-7B. We tested models at the zero-shot baseline, with Retrieval Augmented Generation (RAG) using clinical guidelines, after Supervised Fine-Tuning (SFT) using clinical guidelines, and RAG with a fine-tuned model. Of note, we could not fine-tune GPT-o1 and Claude-3-Opus due to company restrictions on accessing model weights.

To create the external knowledge dataset used for RAG and SFT, we collected six guideline documents for UGIB (related to variceal and non-variceal bleeding) created by major Northern American, European, and Asia-Pacific societies<sup>19–24</sup>. Following our previously published protocol<sup>12</sup>, we reformatted the original documents from raw PDF formats to ones suitable for LLMs, as described elsewhere<sup>12</sup>. This involved converting all information, both text and non-text, into a textual format, creating a coherent structure across all guidelines, and dividing each document into three macro sections: pre-endoscopic, endoscopic, and post-endoscopic management.

For retrieval augmented generation (RAG)<sup>39</sup>, the reformatted guidelines were integrated according to each model's context window size. RAG is a technique that combines retrieval of relevant documents with generation, enabling the model to produce more accurate and contextually appropriate responses. For example, OpenAI's GPT-3.5-turbo can take an input context of up to 4096 tokens, roughly equal to 800 English words. Due to this constraint, each clinical guideline was split into smaller sections, or "chunks," of text at the paragraph level. When a user inputs a query to RAG-GPT-3.5-Turbo, it first searches the most relevant text among the chunks by similarity search using cosine similarity and selects the chunk with the highest similarity. The same chunking strategy was used for LLaMA-2-7B, LLaMA-2-13B, LLaMA-2-70B, and Mistral-7B. On the other hand, OpenAI's GPT-4-Turbo, GPT-4o, and GPT-o1-preview have a context window of up to 128,000 tokens, whereas Anthropic's Claude-3-Opus has a context window of up to 200,000 tokens allowing for chunking at the document level. In these cases, we provided three chunks: one containing the Northern American Guidelines, one with European Guidelines, and one with Asia-Pacific Guidelines.

Supervised fine-tuning was performed using low-rank adaptation (LoRA)<sup>40,41</sup>, which updates a small fraction of the model's parameters, significantly reducing the computational cost and memory usage compared to traditional fine-tuning methods. We employed LoRA to fine-tune GPT-3.5-Turbo, GPT-4-Turbo, GPT-4o, LLaMA-2-7B, LLaMA-2-13B, LLaMA-2-70B, and Mistral-2-7B on the reformatted clinical guidelines. We performed human-guided chunking at the paragraph level, obtaining 96 chunks in total. Train/test split was not performed randomly but was designed to ensure complete information about each management part in training to avoid loss of key information. We used the United States clinical guidelines as the training dataset, and the European/Asia-Pacific guidelines as the testing dataset. Technical details related to the fine-tuning process are reported in the Supplementary Materials.

**Table 3 | List of Expert-Generated Questions for Upper Gastrointestinal Bleeding Management**

Direct content retrieval	
1	Which risk stratification score should I use to assess for very-low-risk patients with UGIB, and what threshold should I use to discharge them from the ED?
2	At what hemoglobin level should I transfuse red blood cells for patients presenting with acute UGIB?
3	Should I use erythromycin as a pre-endoscopic therapy?
4	How should I use epinephrine in endoscopic therapy for patients with NVUGIB?
5	When should I consider pre-emptive TIPS therapy for patients with acute UGIB from portal hypertensive bleeding?
6	How should I manage a patient with rebleeding after initial endoscopic therapy for a bleeding ulcer (Forrest IIa, treated with epinephrine and hemoclips)?
7	How should I manage a patient who had rebleeding after initial endoscopic therapy for a bleeding ulcer, had repeat endoscopic therapy and now is bleeding again? Should I recommend surgery or interventional radiology and why?
8	Should Proton Pump Inhibitor therapy be given to all patients presenting with UGIB even before endoscopy?
9	What is the best time for endoscopy for patients with UGIB? Does this change with variceal bleeding?
Analysis of clinical context	
1	A 30 year-old woman with no significant past medical history presents to the emergency department with an episode of melena. She reports some epigastric discomfort for the past week but denies any history of peptic ulcer disease, alcohol abuse, or use of NSAIDs. She denies any dizziness, weakness, chest pain, or shortness of breath. Her vital signs are within normal limits: blood pressure 120/80 mmHg, pulse 70 bpm, respiratory rate 16 breaths per minute, and temperature 98.6 °F. On physical examination, she appears well, abdomen is soft and non-tender, with no signs of peritoneal irritation or organomegaly. Her initial labs show a hemoglobin of 12 g/dL, normal liver function tests, and normal coagulation profile. She has a Glasgow-Blatchford score of 1. How should this patient be managed in the first 12 h? Should she undergo red blood cell transfusion or upper endoscopy within 24 h?
2	A 65 year-old man with a history of chronic NSAID use for arthritis presents to the emergency department with sudden onset of melena and mild epigastric pain. He denies any other symptoms such as dizziness or weakness. His vital signs are stable: blood pressure 130/80 mmHg, pulse 75 bpm, respiratory rate 18 breaths per minute, and temperature 98.4 °F. His initial labs show a hemoglobin of 10 g/dL (down from his baseline of 14 g/dL), normal liver function tests, and normal coagulation profile. He is admitted for further evaluation and management. The EGD reveals a gastric ulcer with active oozing (Forrest Ib). Endoscopic therapy is successful in achieving hemostasis using a combination of epinephrine injection and application of hemoclips. Should we prescribe PPI? If so, what is the recommended dosage and therapy duration?
3	A 75 year-old man with a previous stroke and atrial fibrillation on apixaban presents to the emergency department with hematemesis and melena. His vital signs are stable: blood pressure 130/80 mmHg, pulse 80 bpm (irregular), respiratory rate 18 breaths per minute, and temperature 98.2 °F. His initial labs show a hemoglobin of 9 g/dL (down from his baseline of 14 g/dL), normal liver function tests, and prolonged coagulation profile due to the apixaban. He is admitted for further evaluation and management. EGD reveals a bleeding duodenal ulcer with active oozing (Forrest Ib). Endoscopic therapy is successful in achieving hemostasis using a combination of thermal therapy and epinephrine injection. Following the procedure, he is started on a high-dose PPI therapy. How should this patient be managed after endoscopy? When should we restart apixaban?
4	A 50 year-old woman with a history of cirrhosis secondary to alcohol use disorder decompensated by ascites presents to the emergency department with acute onset hematemesis. On exam she has dried blood around her mouth, has icteric sclera, no asterixis and moderate abdominal distension with a fluid wave. She denies any other symptoms such as dizziness or weakness. Her vital signs are: blood pressure 110/75 mmHg, pulse 90 bpm, respiratory rate 16 breaths per minute, and temperature 98.6 °F. Her initial labs show a hemoglobin of 7.5 g/dL, ALT 45 (IU/L), AST 103 (IU/L), Total Bilirubin 3.4 mg/dL, and Alkaline Phosphatase 137 (IU/L), INR 1.3, and Albumin 2.9 (g/dL). She is admitted for further evaluation and management. How should this patient be managed?

The questions encompass two main categories: direct content retrieval (i.e., extraction of straight-to-the-point information from clinical guidelines text) and analysis of clinical context (i.e., extraction and interpretation of text from clinical guidelines to answer a clinical case).

### Benchmark datasets and human-grading

To ensure methodological rigor in our framework evaluation across multiple datasets, we implemented a standardized documentation structure to address the following four items: the question dataset (which encompasses the methodological approach to dataset construction and question development), the answer generation process (which delineates the systematic implementation of LLMs for response generation), the answer review criteria (which explicates the comprehensive evaluation protocol employed for response assessment), and the task (which specifies the precise validation objective within our framework's evaluation schema). Each dataset is systematically analyzed through these four methodological dimensions. Before proceeding, it is important to highlight that human-evaluation of the accuracy of LLM-generated answers is based on the following criteria: (1) the answer was entirely accurate and free from any inaccuracies, (2) the answer directly addressed the question posed, and (3) the answer was comprehensive, providing a complete response that covered all critical aspects of the question.

The first benchmarking dataset was the expert-generated UGIB questions. We created a 13-question expert-generated dataset written in conjunction with the expert-of-experts who were senior authors (in North America, Europe, and Asia-Pacific regions) of clinical guidelines for UGIB (L.L., A.B., G.G.T., I.G., J.S.) focused on areas of high value and relevance to the care of patients with UGIB. These key topics encompassed the full spectrum of UGIB care, from initial risk assessment and pre-endoscopic management through to post-procedural care (e.g., risk stratification, transfusion thresholds, or resuming of anticoagulant medication). The

questions were separated into two types of question-related tasks: direct content retrieval ( $n = 9$ ) and analysis of clinical context ( $n = 4$ ) in the form of clinical cases (Table 3). These cases were specifically designed to test the ability to integrate multiple guideline recommendations in realistic clinical contexts.

We also invited those five expert-of-experts to independently provided free-text answers (i.e., “golden-labels”) to each question, collected on the Qualtrics Platform. Each answer was stored in a separate dataset, with the number of characters and word for each question. Each expert answer is reported in the Supplementary Files.

Using these expert-curated questions, we also generated responses using all LLM configurations at a temperature setting of 0.8<sup>42</sup>, producing ten answers per question for each configuration for a total of 3510 responses. These same questions were previously used to collect responses from five different model configurations (i.e., baseline PaLM, baseline GPT-3.5, baseline GPT-4, RAG-GPT-3.5, RAG-GPT-4) across multiple temperature thresholds (0.0 to 2.0, with 0.2 increments), creating a dataset of 8580 answers. We generated an additional dataset ( $n = 1430$ ) using only the best-performing model configuration, following the same temperature range pattern. In all cases, through heuristic prompt engineering, we constrained LLM response lengths to match the maximum word count of the corresponding expert answers, ensuring comparable response formats.

Two independent gastroenterologists blindly evaluated the accuracy of the responses generated at temperature 0.8, comparing them against clinical guidelines and expert answers. In cases of disagreement, a third expert reviewer served as a tiebreaker (disagreement requiring a tiebreaker



happened in 6.6% of cases). Four medical experts independently graded the responses generated across different temperature thresholds, and majority voting was used to resolve any disagreements.

The expert responses (“golden labels”) were used to develop and evaluate different text similarity approaches. The LLM-generated responses at temperature 0.8 were used as a validation benchmark to evaluate which similarity technique (fine-tuned ColBERT, Sentence Transformers, and TF-IDF) best correlated with actual model performance. The historical temperature-varying dataset ( $n = 8580$ ) served for training and internal validation, while the additional dataset from the best-performing model ( $n = 1430$ ) was used for external validation of the reward model.

The second benchmarking dataset was obtained from the American College of Gastroenterology (ACG) Multiple-Choice Questions (MCQs). Among all self-assessment board preparation tests published by the ACG, only 40 MCQs strictly focused on the management of patients with UGIB. To establish a benchmark for human performance, we calculated the pooled percentage of correct answers from previous practicing ACG physician test-takers at varying career stages, which averaged 75% for these specific questions. This dataset cannot be released due to the proprietary nature of the MCQs.

Each LLM configuration was tested using a zero-shot approach, where models were instructed to provide only the letter corresponding to the correct answer among the available choices, without any additional explanation or context. All responses were generated using a temperature setting of 0.8.

Two independent reviewers evaluated the number of correct responses for each LLM configuration, comparing them against the reference answers.

This dataset served as a validation benchmark to evaluate which similarity technique (fine-tuned ColBERT, Sentence Transformers, and TF-IDF) best correlated with actual model performance.

The third benchmarking dataset was obtained from real-world questions from the Simulation Scenario. In particular, we compiled a dataset of 117 questions from 82 physician trainees across 29 sessions involving 5 standardized UGIB scenarios, conducted in medical simulation settings between 2023–2024 (IRB protocol number #2000034521). The complete list of scenarios and related questions is provided in the Supplementary Materials. The simulation scenarios were designed as part of a clinical trial evaluating the LLM interface (named GUT-GPT) effectiveness in clinical decision support, which was conducted in accordance with the ethical principles outlined in the Declaration of Helsinki<sup>43</sup>. Each clinical case-question pair is reported in the Supplementary Files.

Each LLM configuration was tested using a heuristic prompting approach, necessary due to the unpredictable nature of trainee questions.

The prompts were structured to include complete clinical case analysis, providing all relevant context (including patient demographics, laboratory findings, and clinical presentation) and requesting both case-specific information and management recommendations based on the trainee’s specific query. This approach allowed the models to address both direct management questions and requests for case-specific information (e.g., age, laboratory values, etc.). All responses were generated using a temperature setting of 0.8.

Two independent gastroenterologists blindly evaluated the accuracy of responses for each LLM configuration against established clinical guidelines. In cases of disagreement, a third expert reviewer served as a tiebreaker (disagreement requiring a tiebreaker happened in 9.5% of cases).

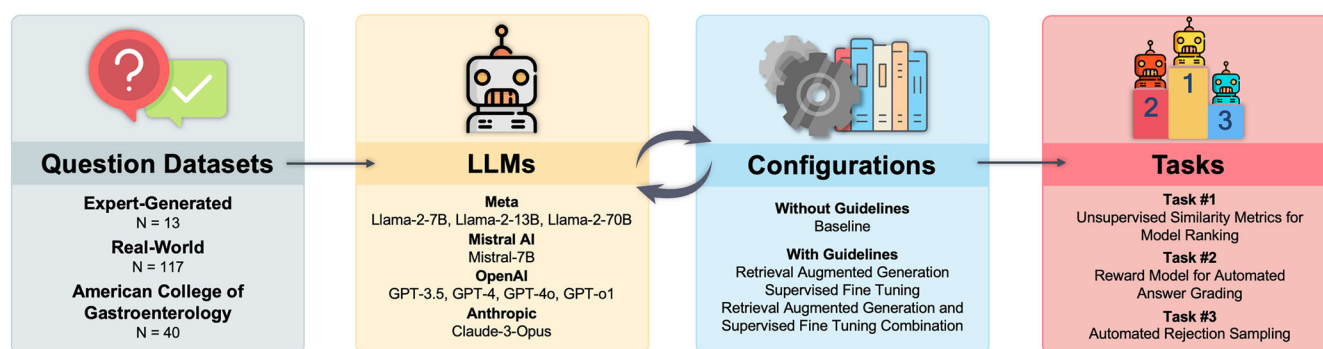
This dataset served as a validation benchmark to evaluate which similarity technique (fine-tuned ColBERT, Sentence Transformers, and TF-IDF) best correlated with actual model performance. This dataset was also used for a supplementary analysis of the reward model alignment with human-grading.

### Unsupervised similarity metrics alignment with expert-of-expert golden labels

The EVAL framework provides a scalable solution for AI safety in clinical settings through complementary approaches operating at two levels: at the model level, using unsupervised embeddings to automatically evaluate and rank different LLM configurations based on expert-generated answers (“golden labels”), and at the answer level, employing a reward model to screen individual responses for accuracy against guideline-based recommendations, as illustrated in Fig. 4.

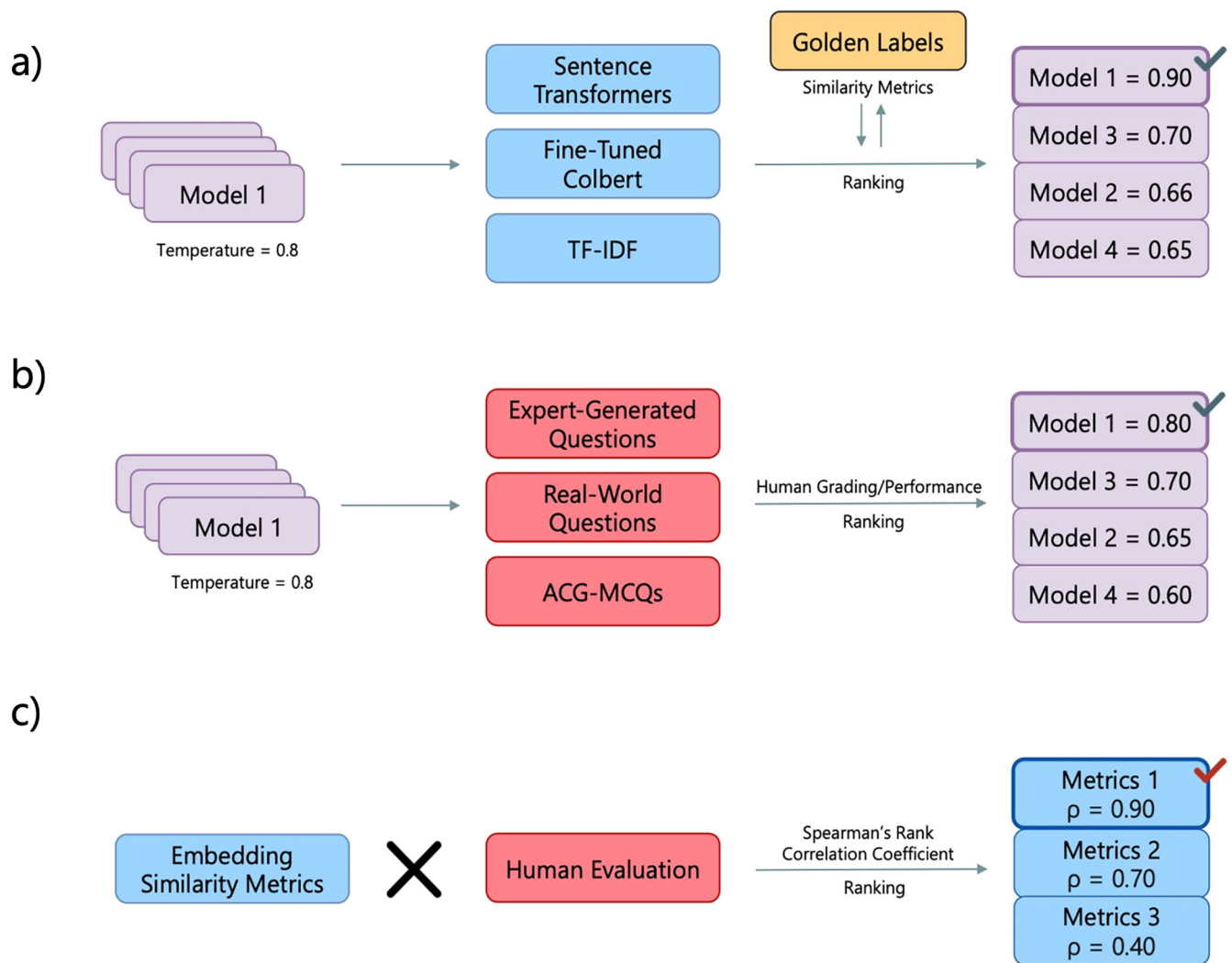
We evaluated three different similarity metrics to quantify the alignment between LLM-generated responses and expert-provided answers: Contextualized Late Interaction over BERT (ColBERT), Sentence Transformers, and TF-IDF as summarized in Fig. 5.

We used ColBERT<sup>44</sup> to quantify the alignment between responses generated by LLMs and responses by experts (Fig. 5). We chose ColBERT for its ability to handle the variability of responses within a relatively small semantic space, and its unique token-level comparison approach. Unlike traditional embedding methods that create a single vector representing an entire text (paragraph-level embedding or “early aggregation”), ColBERT preserves the meaning of individual words or tokens separately and compares these individual representations between texts before making a final similarity decision (token-level embedding or “late interaction”). This approach allows for more precise matching of specific clinical terms and concepts in context, rather than simply comparing overall text meanings. To



**Fig. 4 | EVAL framework summary.** The EVAL framework consists of three interconnected components. The first component comprises the Question Datasets: expert-generated questions ( $N = 13$ ), real-world questions ( $N = 117$ ), and American College of Gastroenterology questions ( $N = 40$ ). The second component shows the LLM configurations, which combines different LLM architectures (Meta’s Llama-2-7B/13B/70B, Mistral AI’s Mistral-7B, OpenAI’s GPT-3.5/4/4o/o1, and Anthropic’s Claude-3-Opus) with various configurations (without guidelines as baseline, with

guidelines through Retrieval Augmented Generation, Supervised Fine Tuning, and a combination of Retrieval Augmented Generation and Supervised Fine Tuning). These LLMs and configurations are then evaluated through three distinct tasks: Task #1 uses unsupervised similarity metrics for model ranking, Task #2 employs a reward model for automated answer grading, and Task #3 implements automated rejection sampling to ensure response quality and safety.



**Fig. 5 | Evaluation and validation framework for embedding similarity metrics.** This figure illustrates a comprehensive framework for evaluating the alignment of responses generated by large language models (LLMs) with expert-defined Golden Labels (i.e., free-text answers from the experts). **a** Step 1 - Embedding Similarity Metrics: Model ranking by comparing the similarity of LLM-generated answers to the Golden Labels using TF-IDF, Sentence Transformers, and Fine-Tuned ColBERT. Fine-tuning was performed to maximize the cosine similarity between the embeddings of the “golden labels” and their corresponding paragraphs while minimizing similarity with unrelated paragraphs. This step enhances the model’s ability to differentiate between relevant and irrelevant responses. **b** Step 2 - Model Performance Evaluation: model responses were assessed by human experts, who

graded them for accuracy using expert-generated datasets, real-world questions, and the American College of Gastroenterology Multiple-Choice Questions (ACG-MCQs). Models were then ranked based on their performance and accuracy scores. **c** Step 3 - Selection of the Best Embedding Similarity Metrics: the average similarity values for each model were correlated with human performance evaluations using Spearman’s rank correlation coefficient. This process identified the similarity metrics with the highest correlation coefficients, underscoring their utility in assessing model response quality. Abbreviations: ACG-MCQs American College of Gastroenterology Multiple Choice Questions, TF-IDF Term Frequency-Inverse Document Frequency, ColBERT Contextualized Late Interaction over BERT.

enhance precision in distinguishing between high-quality and lower-quality responses, we fine-tuned the ColBERT embeddings as follows: for each expert label, we created triplets consisting of the label itself, a closely matching paragraph, and a non-matching paragraph from a set of clinical guidelines. We used Bidirectional Encoder Representations from Transformers (BERT)<sup>45</sup> embeddings for each triplet component. The matching paragraphs were chosen based on their high relevance to the expert label, while the non-matching paragraphs were selected based on their slight, but not complete, irrelevance (an example is provided in Supplementary Table 3). The objective function for fine-tuning maximized the cosine similarity between the embeddings of the expert label and the matching paragraph while minimizing the similarity between the expert label and the non-matching paragraphs. This is achieved using pairwise softmax cross-entropy loss, which effectively pushes the model to enhance the distinction between relevant and irrelevant responses regarding embedding proximity. Fine-tuned ColBERT can produce a more refined separation between

relevant and irrelevant text snippets. To account for the plurality of opinions from multiple experts, we evaluated this by calculating the average similarity score across multiple sets of embeddings generated from a variety of responses to different questions. This score reflects the overall alignment of the model’s generated responses with expert-provided answers (details in Supplementary Materials.) To validate model ranking accuracy, we compared the ranking of the Fine-Tuned Colbert to the accuracy rankings of each LLM configuration for the expert-generated answer dataset and the performance on ACG-MCQs. For better visualization of the relative gap between the Colbert score from different models, we provide the transformation of first normalizing the Colbert raw score with its maximum attainable score and then applying the logit function. To showcase the performance of our Fine-Tuned Colbert method, we provide the following two baselines: Sentence Transformer<sup>46</sup>, a common existing LLM-based method for textual similarity, and TF-IDF<sup>47</sup>, which is a classical method based on word and document statistics.

For the Sentence Transformers-based similarity metrics, we use the publicly available pre-trained embedding model, all-MiniLM-L6-v2, from Sentence Transformer<sup>38</sup> to calculate embeddings for answers and then use the cosine similarity to calculate the score between a pair of answer embeddings. The model is a pre-trained BERT model further finetuned by paired sentences optimized for producing high similarity scores for paired sentences. It's oftentimes a decent approach for similarity tasks and thus serves as a well-suited baseline to be compared with our model.

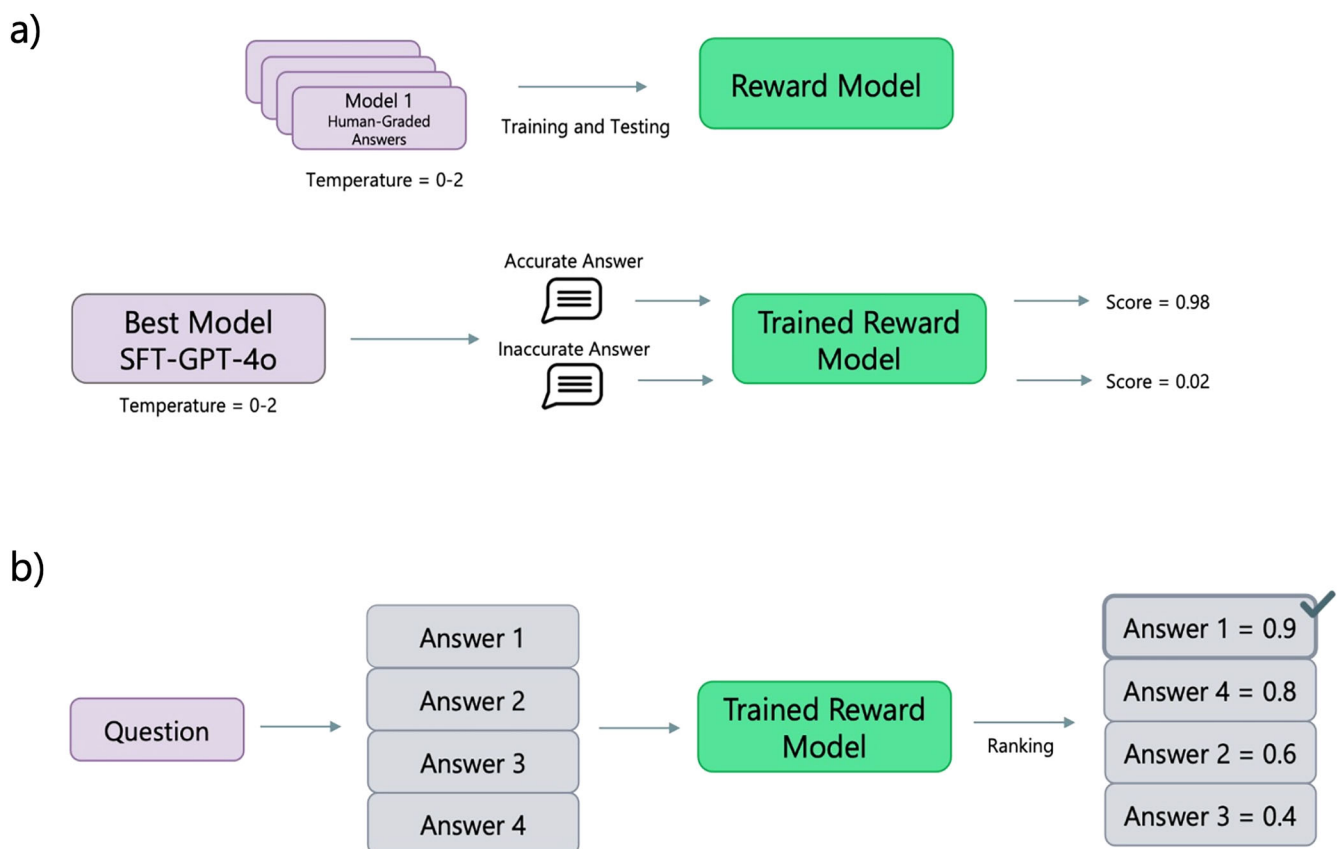
For the TF-IDF-based similarity metric, we follow the standard practice of calculating the feature vector and then compare feature vectors with cosine similarity, which falls under the similar framework of our Colbert method, with the difference being TF-IDF uses pre-defined statistics instead of our highly specialized data-driven Colbert. Specifically, for each pair of LLM output and expert response, we calculate the TF-IDF score by multiplying the term frequency and inverse document frequency. In this context, the document is either one LLM output or one expert answer. The term frequency, TF, is the number of times a given term appears in the document. The inverse document frequency, IDF, is the ratio of one plus the total number of documents divided by one plus the number of documents having the term, then take the log and add one again. The several constant value ones are in place for normalizing and avoiding the divided by zero issues and is the standard common approach<sup>48</sup>. Lastly, we calculate the cosine similarity between the calculated TF-IDF score to serve as the final similarity score.

For each similarity method, we performed pairwise *t*-tests comparing the highest-scoring model configuration against all other configurations individually. Similarly, we conducted pairwise *t*-tests for human-graded accuracies across the three evaluation sets (expert-generated questions, real-world questions, and ACG MCQs), comparing the best-performing configuration against all others. For all statistical comparisons, we considered a two-tailed *p*-value < 0.05 as statistically significant. To determine which similarity metric best aligned with human evaluation, we calculated Spearman rank correlation coefficients between the average scores from each method and the model accuracies determined by human grading. This analysis allowed us to identify which of the three proposed methods showed the strongest alignment with both human-graded accuracy and performance on ACG MCQs.

### Reward model to screen for high-quality LLM responses

One concern of deploying probabilistic large language models in clinical settings is the presence of hallucinations—seemingly plausible but inaccurate information<sup>49</sup>. It is not uncommon for models to output answers that contain factual inaccuracies or “misread” the guidelines, or to be *confidently incorrect* in giving factually incorrect information without any indication of uncertainty. This part of our framework that addresses the issue of hallucinations is represented graphically in Fig. 6.

As a solution to the best model selection, we employ an alternative approach by training an additional Reward Model to serve as a substitute for



**Fig. 6 | Reward model training, testing, and validation and application with automated rejection sampling.** This figure illustrates a two-step framework for optimizing the accuracy and reliability of responses generated by large language models (LLMs), with clear stages for reward model training and application. **a** Step 1 - Reward Model Training and Validation: previously graded answers from the expert-generated questions were utilized for training and testing the reward model. The reward model assigns accuracy scores to the generated answers (e.g., 0.98 for accurate responses and 0.02 for inaccurate ones). Validation was performed using human-graded answers from the best-performing model, determined through

Fine-Tuned ColBERT ranking. This process ensured that the reward model could accurately evaluate the quality of new question-answer pairs, thereby validating its grading accuracy. **b** Step 2 - Application with Automated Rejection Sampling: For each question, the LLM generates multiple candidate answers (K answers). These answers are passed through the trained reward model, which assigns accuracy scores and ranks the responses. The answer with the highest score is selected as the final output. This filtering mechanism increases the reliability of the model by systematically rejecting less accurate responses, thereby ensuring only the most accurate answers are retained.

human feedback. A reward model is an LLM tasked with approximating part of the traditional environment in a reinforcement learning problem. The reward model takes in text and returns a score. The objective of this reward model is to assess the level of congruence between a model's response and human preferences. In simpler terms, a reward model is a type of model that takes a pair of inputs (prompt and response) and produces an output in the form of a reward or score. The primary difficulty in constructing such a model lies in obtaining a dataset of high quality. The subjective evaluation of good and bad varies among individuals, making it unfeasible to quantify. Previous evidence suggests that a dataset containing between 1000 and 10000 high-quality question-answer pairs is sufficient for training a reward model in moderately complex domains<sup>50,51</sup>. For larger or more nuanced topics, a dataset exceeding 50000 pairs may be necessary<sup>52</sup>.

To train our reward model, which we will refer to as the Grader Model (GM), the LLM receives data in the following format: [Question, Answer, Score]. The GM's task is to take a specific [Question, Answer] pair and map it to the answer's score. Scores are provided by a human evaluator who reads the response and assigns it a numerical ranking of 0 or 1 based on the accuracy. To train this model, we replace the LLM's traditional head, which outputs the log probability of the next word, with a value head that predicts the score of [Question, Answer] pair. Since the answers are classified as either Good (Score = 1) or Bad (Score = 0), the value head outputs the probability that the answer is good. The model is trained using cross entropy (classification) loss and gradient descent to improve score accuracy.

We used the previously graded dataset ( $n = 7150$ ) obtained from multiple LLM configurations (i.e., baseline PaLM, baseline GPT-3.5, baseline GPT-4, RAG-GPT-3.5 with American Guidelines, RAG-GPT-3.5 with American, European and Asia-Pacific Guidelines) to train the Reward Model, which was then internally validated to the previous state-of-the-art model (i.e., RAG-GPT-4 with American, European and Asia-Pacific Guidelines;  $n = 1430$ ). The Reward Model performance was externally validated using the new state-of-the-art model (i.e., SFT-GTP-4o;  $n = 1430$ ) that was selected according to the highest similarity metrics according to Fine-Tuned Colbert.

The reward model was trained using Meta's OPT-350M, a 350 million parameters decoder-only LLM. The use of a smaller RM such as Meta's OPT-350M aligns with findings indicating that compact models are sufficient for tasks where the dataset quality is prioritized over model scale, as smaller models demonstrate robust generalization and efficiency without significant performance trade-offs in preference learning or alignment tasks, provided they are trained on high-quality, curated datasets<sup>46,53,54</sup>. The reward model output is binary: "Good" (Score = 1) or "Bad" (Score = 0). Alignment to human-experts was evaluated as the number of true labels (i.e., the number of answers for which the reward model produced the same label with human grading). The results were interpreted by breaking down the temperatures into three regimes, *positive* (temperature < 1.2), *negative* (temperature > 1.6), and *mixed* (temperature between 1.2 and 1.6) according to the model's graded performance. These thresholds were chosen such that the *positive* regime has over 80% graded accuracy and the *negative* regime has <20% graded accuracy. The reward model was then applied to the best model according to ColBERT ranking and validated the grading accuracy on this new dataset of question-answer pairs. As a sensitivity analysis, we reported alignment across all temperature thresholds in the Supplementary Materials. In addition, we tested the alignment of the reward model with human grading on the real-world questions for all models at the fixed temperature of 0.8, with results being reported in the Supplementary Materials. The reward model is publicly available on Hugging Face ([https://huggingface.co/ZachariahPang/medical\\_reward\\_model](https://huggingface.co/ZachariahPang/medical_reward_model)).

### Automated rejection sampling

Extending the reward model pipeline, we can incorporate the reward function directly into the answer pipeline by using a rejection sampling approach. For each question, the LLM agent generates  $K$  candidate answers. These  $K$  answers are evaluated by the reward model, and only the top-

scoring answer is sent forward. This serves as a form of self-filtering, allowing the reward model to capture and filter out suboptimal answers before they reach the end user. In this way, rejection sampling enhances the model's overall output quality by rescuing from suboptimal answers. To evaluate the rejection sampling approach, we used the same curated dataset for reward model alignment described in the previous section. Human-graded accuracy was compared across multiple  $K$  values (1, 3, 5, 7, and 10), as reported in the Supplementary Table 4. The results demonstrated a consistent improvement in accuracy with increasing  $K$ . However, larger  $K$  values also demand significantly more computational resources. We selected  $K = 5$  for the main analysis as it provides a practical balance between computational efficiency and improved accuracy. Detailed trends in accuracy with and without rejection sampling, as well as the impact of varying  $K$ , are included in the Supplementary Materials to illustrate the trade-offs and performance improvements.

### Data availability

Expert-generated questions are available in Table 3 of the manuscript, while expert free-text answers and real-world clinical questions can be found in the supplementary files.

### Code availability

Code can be provided based on personal requests. Please contact the corresponding author. The reward model has been uploaded on Hugging Face at the following link: [https://huggingface.co/ZachariahPang/medical\\_reward\\_model](https://huggingface.co/ZachariahPang/medical_reward_model).

Received: 16 August 2024; Accepted: 25 March 2025;

Published online: 03 May 2025

### References

1. Singhal, K. et al. Large language models encode clinical knowledge. *Nature* **620**, 172–180 (2023).
2. Peng, C. et al. A study of generative large language model for medical research and healthcare. *NPJ Digit. Med.* **6**, 210 (2023).
3. Giuffrè, M., You, K. & Shung, D. Evaluating ChatGPT in medical contexts: the imperative to guard against hallucinations and partial accuracies. *Clin. Gastroenterol. Hepatol.* <https://doi.org/10.1016/j.cgh.2023.09.035> (2023).
4. Giuffrè, M. & Shung, D. L. Scrutinizing chatGPT applications in gastroenterology: a call for methodological rigor to define accuracy and preserve privacy. *Clin. Gastroenterol. Hepatol.* <https://doi.org/10.1016/j.cgh.2024.01.024> (2024).
5. Hager, P. et al. Evaluation and mitigation of the limitations of large language models in clinical decision-making. *Nat. Med.* **30**, 2613–2622 (2024).
6. Fraser, H. et al. Comparison of diagnostic and triage accuracy of Ada Health and WebMD Symptom Checkers, CHATGPT, and physicians for patients in an emergency department: clinical data analysis study. *JMIR Mhealth Uhealth*. **11**, e49995 (2023).
7. Wilhelm, T. I., Roos, J. & Kaczmarczyk R. Large language models for therapy recommendations across 3 clinical specialties: comparative study. *J. Med. Internet Res.* **25**, e49324 (2023).
8. Soroush, A., Giuffrè, M., Chung, S. & Shung, D. L. Generative Artificial Intelligence in Clinical Medicine and Impact on Gastroenterology. *Gastroenterology* <https://doi.org/10.1053/j.gastro.2025.03.038> (2025).
9. Lee, P., Bubeck, S. & Petro, J. Benefits, limits, and risks of GPT-4 as an AI chatbot for medicine. *N. Engl. J. Med.* **388**, 1233–1239 (2023).
10. Giuffrè, M. et al. Systematic review: The use of large language models as medical chatbots in digestive diseases. *Aliment. Pharmacol. Ther.* **60**, 144–166 (2024).
11. Ge, Y., Guo, Y., Das, S., Al-Garadi, M. A. & Sarker, A. Few-shot learning for medical text: a review of advances, trends, and opportunities. *J. Biomed. Inf.* **144**, 104458 (2023).



12. Kresevic, S. et al. Optimization of hepatological clinical guidelines interpretation by large language models: a retrieval augmented generation-based framework. *NPJ Digit. Med.* **7**, 102 (2024).
13. Shah, N. H., Entwistle, D. & Pfeffer, M. A. Creation and adoption of large language models in medicine. *JAMA* **330**, 866 (2023).
14. Meskó, B. & Topol, E. J. The imperative for regulatory oversight of large language models (or generative AI) in healthcare. *NPJ Digit. Med.* **6**, 120 (2023).
15. Guyatt, G. Evidence-based medicine. *JAMA* **268**, 2420 (1992).
16. Institute of Medicine (US) Committee on Standards for Developing Trustworthy Clinical Practice Guidelines; Graham, R. et al., editors. *Clinical Practice Guidelines We Can Trust*. Washington (DC): National Academies Press (US); 2011. 2, Background and Key Stakeholders in Guidelines Development and Use. Available from: <https://www.ncbi.nlm.nih.gov/books/NBK209534/>.
17. Zheng, N. S., Tsay, C., Laine, L. & Shung, D. L. Trends in characteristics, management, and outcomes of patients presenting with gastrointestinal bleeding to emergency departments in the United States from 2006 to 2019. *Aliment Pharm. Ther.* **56**, 1543–1555 (2022).
18. Rosenstock, S. J. et al. Improving quality of care in peptic ulcer bleeding: nationwide cohort study of 13,498 consecutive patients in the danish clinical register of emergency surgery. *Am. J. Gastroenterol.* **108**, 1449–1457 (2013).
19. Gralnek, I. M. et al. Endoscopic diagnosis and management of nonvariceal upper gastrointestinal hemorrhage (NVUGIH): European society of gastrointestinal endoscopy (ESGE) guideline – update 2021. *Endoscopy* **53**, 300–332 (2021).
20. Laine, L., Barkun, A. N., Saltzman, J. R., Martel, M. & Leontiadis, G. I. ACG clinical guideline: upper gastrointestinal and ulcer bleeding. *Am. J. Gastroenterol.* **116**, 899–917 (2021).
21. Abraham, N. S. et al. American college of gastroenterology-Canadian association of gastroenterology clinical practice guideline: management of anticoagulants and antiplatelets during acute gastrointestinal bleeding and the perendoscopic period. *Am. J. Gastroenterol.* **117**, 542–558 (2022).
22. de Franchis, R. et al. Baveno VII – renewing consensus in portal hypertension. *J. Hepatol.* **76**, 959–974 (2022).
23. Kaplan, D. E. et al. AASLD Practice Guidance on risk stratification and management of portal hypertension and varices in cirrhosis. *Hepatology* **79**, 1180–1211 (2024).
24. Sung, J. J. et al. Asia-Pacific working group consensus on non-variceal upper gastrointestinal bleeding: an update 2018. *Gut* **67**, 1757–1768 (2018).
25. Barkun, A. N. et al. Effectiveness of disseminating consensus management recommendations for ulcer bleeding: a cluster randomized trial. *Can. Med. Assoc. J.* **185**, E156–E166 (2013).
26. Lu, Y., Barkun, A. N. & Martel, M. Adherence to guidelines: a national audit of the management of acute upper gastrointestinal bleeding. The REASON registry. *Can. J. Gastroenterol. Hepatol.* **28**, 495–501 (2014).
27. Liang, P. S. & Saltzman, J. R. A national survey on the initial management of upper gastrointestinal bleeding. *J. Clin. Gastroenterol.* **48**, e93–e98 (2014).
28. Prosenz, J., Stättermayer, M. -S., Riedl, F. & Maieron, A. Adherence to guidelines in patients with non-variceal upper gastrointestinal bleeding (UGIB) – results from a retrospective single tertiary center registry. *Scand. J. Gastroenterol.* **58**, 856–862 (2023).
29. McDuff, D. et al. Towards accurate differential diagnosis with large language models. *Nature* <https://doi.org/10.1038/s41586-025-08869-4> (2025).
30. Saab, K. et al. Capabilities of gemini models in medicine. *arXiv* <https://doi.org/10.48550/arXiv.2404.18416> (2024).
31. Bedi, S. et al. Testing and Evaluation of Health Care Applications of Large Language Models: A Systematic Review. *JAMA*. **333**, 319–328 (2025).
32. Nori, H. et al. Capabilities of GPT-4 on medical challenge problems (2023).
33. Gilson, A. et al. How does ChatGPT perform on the United States Medical Licensing Examination (USMLE)? the implications of large language models for medical education and knowledge assessment. *JMIR Med. Educ.* **9**, e45312 (2023).
34. Ferber, D. et al. GPT-4 for information retrieval and comparison of medical oncology guidelines. *NEJM AI* **1**, 6 (2024).
35. Unlu, O. et al. Retrieval-augmented generation-enabled GPT-4 for clinical trial screening. *NEJM AI* **1**, 7 (2024).
36. Zakka, C. et al. Almanac — retrieval-augmented language models for clinical medicine. *NEJM AI* **1**, 2 (2024).
37. Luo, Y. et al. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *arXiv* <https://doi.org/10.48550/arXiv.2308.08747> (2023).
38. Reimers, N. & Gurevych, I. Sentence-BERT: Sentence embeddings using siamese BERT-networks. *arXiv* <https://doi.org/10.48550/arXiv.1908.10084> (2019).
39. Lewis, P. et al. Retrieval-augmented generation for knowledge-intensive NLP Tasks. *arXiv* <https://doi.org/10.48550/arXiv.2005.11401> (2020).
40. Dettmers, T., Pagnoni, A., Holtzman, A. & Zettlemoyer, L. QLoRA: Efficient finetuning of quantized LLMs. *arXiv* <https://doi.org/10.48550/arXiv.2305.14314> (2023).
41. Hu, E. J. et al. LoRA: Low-rank adaptation of large language models. *arXiv* <https://doi.org/10.48550/arXiv.2106.09685> (2021).
42. Giuffrè, M. et al. Su1979 GUTGPT: novel large language model pipeline outperforms other large language models in accuracy and similarity to international experts for guideline recommendation management of patients with upper gastrointestinal bleeding. *Gastroenterology* **166**, S-889–S-890 (2024).
43. Rajashekar, N. C. et al. Human-algorithmic interaction using a large language model-augmented artificial intelligence clinical decision support system. In *Proc. CHI Conference on Human Factors in Computing Systems* 1–20 (ACM, New York, NY, USA, 2024).
44. Khattab, O. & Zaharia, M. ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. *arXiv* <https://doi.org/10.48550/arXiv.2004.12832> (2020).
45. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv* <https://doi.org/10.48550/arXiv.1810.04805> (2018).
46. Stiennon, N. et al. Learning to summarize from human feedback. *arXiv* <https://doi.org/10.48550/arXiv.2009.0132> (2020).
47. Sparck Jones, K. A statistical interpretation of term specificity and its application in retrieval. *J. Document.* **28**, 11–21 (1972).
48. Pedregosa, F. et al. *Scikit-learn: Machine Learning in Python*. <https://scikit-learn.org/stable/> (2012).
49. Dhuliawala, S. et al. Chain-of-verification reduces hallucination in large language models. *arXiv* <https://doi.org/10.48550/arXiv.2309.11495> (2023).
50. Nath, S. et al. Leveraging domain knowledge for efficient reward modelling in RLHF: a case-study in E-commerce opinion summarization. *arXiv* <https://doi.org/10.48550/arXiv.2402.15473> (2024).
51. Wang, Z. et al. HelpSteer2: Open-source dataset for training top-performing reward models. *arXiv* <https://doi.org/10.48550/arXiv.2406.08673> (2024).
52. Ouyang, L. et al. Training language models to follow instructions with human feedback. *arXiv* <https://doi.org/10.48550/arXiv.2203.02155> (2022).
53. Rafailov, R. et al. Direct preference optimization: your language model is secretly a reward model. *arXiv* <https://doi.org/10.48550/arXiv.2305.18290> (2023).
54. Bai, Y. et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv* <https://doi.org/10.48550/arXiv.2204.05862> (2022).

## Acknowledgements

DLS is supported by NIH NIDDK grant DK125718.

## Author contributions

M.G., K.Y., Z.P., S.K., B.S., and D.L.S.: Conceptualization; methodology; validation; investigation; data curation; writing—original draft; writing—review and editing; visualization; supervision. S.C., R.C., Y.K., C.C., T.S., M.A., L.S.C., G.G., I.G., J.J.Y.S., A.B., L.L., J.S.: Conceptualization; methodology; writing—review and editing.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41746-025-01589-z>.

**Correspondence** and requests for materials should be addressed to Dennis L. Shung.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025