



**UNIVERSITÀ  
DEGLI STUDI  
DI TRIESTE**

**UNIVERSITÀ DEGLI STUDI DI TRIESTE**

**XXXVIII CICLO DEL DOTTORATO DI RICERCA IN**

**APPLIED DATA SCIENCE AND ARTIFICIAL INTELLIGENCE**

**VALIDATION OF GUIDELINE-ALIGNED LARGE  
LANGUAGE MODELS FOR SAFE CLINICAL DECISION  
MAKING IN DIGESTIVE DISEASES**

Settore scientifico-disciplinare: MED/12

DOTTORANDO  
**MAURO GIUFFRÈ**

COORDINATORE  
**PROF. PAULI FRANCESCO**

SUPERVISORE DI TESI  
**PROF. LUCA BORTOLUSSI**

CO-SUPERVISORE DI TESI  
**PROF. LORY SAVERIA CROCÈ**

**ANNO ACCADEMICO 2024/2025**



# Abstract

Large language models (LLMs) promise major gains for clinical decision support in gastroenterology and hepatology, but safe adoption requires more than clever prompting. This dissertation develops and validates a translational pipeline that (i) renders clinical guidance machine-readable, (ii) embeds expert oversight, and (iii) layers an automated safety stack. Across systematic evidence synthesis, expert-benchmarked tests, and society-curated question banks, engineered, guideline-grounded systems achieve expert-adjacent performance on defined tasks while preserving clear boundaries for human supervision.

A systematic review of baseline LLM performance (Chapter 2) showed wide accuracy (6.4–91.4%) across 18 studies, with a median near 50% for larger sets, driven by inconsistent question design, evaluator expertise, and non-standard grading. Foundational models without domain adaptation pose unacceptable safety risks.

Chapter 3 introduced a guideline-grounded retrieval-augmented generation (RAG) framework using European Association for the Study of the Liver (EASL) hepatitis C virus (HCV) guidelines. Ablations showed that converting guidelines into LLM-friendly formats—cleaning text, converting tables to structured lists, and enforcing schemas—plus principled prompts raised accuracy from 43% (baseline GPT-4 Turbo) to 99% overall; table questions from 28% to 96%, clinical scenarios from 20% to 100%.

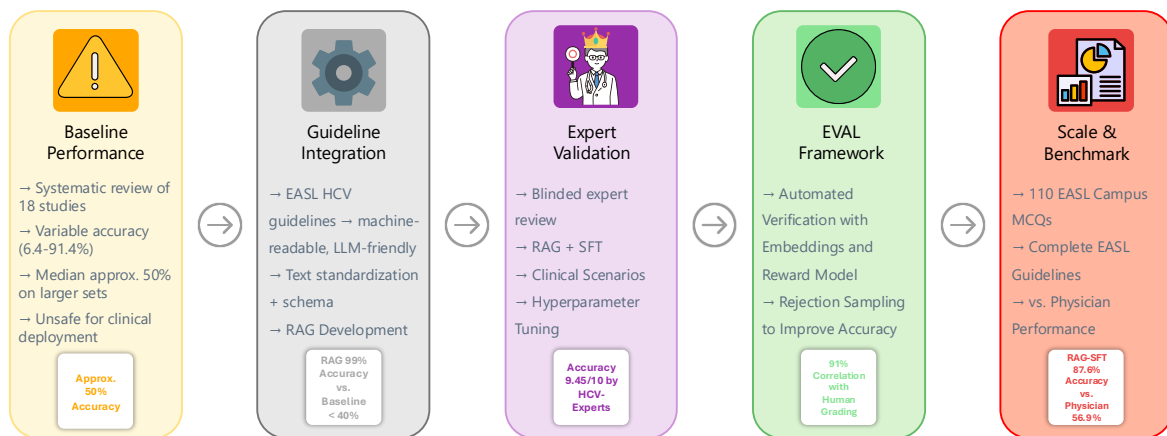
For external validation (Chapter 4), international hepatology experts, including guideline authors, blindly graded multiple configurations. Supervised fine-tuning (SFT) combined with RAG achieved mean scores 9.45/10 (authors) and 8.7/10 (independent clinicians), exceeding baselines (6.4/10,  $p < 0.001$ ). Temperatures 0–0.8 and top-p 0–0.5 minimized hallucinations. On direct-acting antiviral (DAA) regimen selection, median fully correct recommendations rose from 24% to 76% (RAG-Top10).

To scale oversight (Chapter 5), the Expert-of-Experts Verification and Alignment (EVAL) framework used fine-tuned Contextualized Late Interaction over BERT (ColBERT) embeddings (Spearman  $\rho = 0.91$  with human judgment) for model-level ranking. A reward model trained on expert labels reproduced grades at 87.9% accuracy; rejection sampling reached 98% in high-temperature regimes, enabling automated filtering.

Chapter 6 tested generalization on 110 EASL Campus multiple-choice questions (MCQs). Optimized RAG+SFT achieved 87.6% accuracy—31 points above pooled physicians (56.9%,  $p < 0.001$ )—across liver tumors (95.0% vs 50.7%), viral hepatitis (80.0% vs 55.1%), and cirrhosis (80.0% vs 58.2%).

LLMs can approach expert-level performance on guideline-referenced tasks when knowledge is machine-encoded, expert judgment shapes representation and evaluation, and automated verification enforces continuous safety. The work offers a reproducible blueprint for evidence-based AI in gastroenterology while delineating limits that require ongoing human oversight.

## From Baseline Unreliability to Verified Clinical Deployment



**Key Principle:** Accuracy alone is insufficient for clinical deployment. Safe LLM decision support requires **guideline alignment** (machine-readable EBM), **expert validation** (domain expertise shaping evaluation), and **continuous verification** (automated safety monitoring). This reproducible blueprint enables evidence-based AI in gastroenterology while preserving essential human oversight.



# Table of Contents

<b>Preface</b>	<b>1</b>
<b>Chapter 1 — Introduction: Foundations of Large Language Models in Digestive Diseases: Opportunities, Risks, and Privacy</b>	<b>6</b>
1.1 The Importance of Large Language Models in Healthcare	7
1.2 The Transformer Architecture: Technical Foundations of Modern LLMs	8
1.2.1 The Self-Attention Mechanism	8
1.2.2 Encoder–Decoder Structure	9
1.2.3 From Transformers to Large Language Models	9
1.2.4 Computational Advantages	9
1.3 Challenges in Defining Accuracy for LLMs	9
1.4 How to Infuse Domain Knowledge	10
1.4.1 In-Context Learning	10
1.4.2 Supervised Fine-Tuning with Reinforcement Learning from Human Feedback	15
1.5 Possible Applications in Gastroenterology and Hepatology	17
1.5 Ethical Considerations	19
1.7 Privacy and Liability	20
1.8 Overall Aims and Objectives	20
<b>Chapter 2 — Baseline: What General-Purpose Large Language Models Can (and Cannot) Do in Digestive Diseases: A Systematic Review of Performance Evaluation</b>	<b>23</b>
2.1 Chapter Overview	24
2.2 Materials and Methods	25
2.2.1 Search Strategy and Selection Criteria	25
2.2.2 Outcome Assessment	25
2.2.3 Data Extraction	26
2.2.4 Bias Evaluation	26
2.2.5 Data Synthesis and Statistical Analysis	26
2.3 Results	27
2.3.1 Study Characteristics	27
2.3.2 Question Generation	28
2.3.3 Question Clustering	29
2.3.4 Question Evaluators	30
2.3.5 Answers Grading Systems	30
2.3.6 Accuracy Evaluation	32
2.3.7 Study Bias Evaluation	33
2.4 Chapter’s Deliverables	33
<b>Chapter 3 — Build: Making Hepatitis C Virus Guidelines LLM-Friendly to Align Responses to Evidence-Based Medicine Using Retrieval Augmented Generation</b>	<b>35</b>
3.1 Chapter’s Overview	36
3.2 Materials and Methods	37
3.2.1 Guidelines Selection	37
3.2.2 Standardized prompts creation	38
3.2.3 Ablation study: customized LLM framework	39
3.2.4 Primary outcome	41
3.2.5 Secondary outcome	41
3.2.6 Statistical analysis	41
3.3 Results	42
3.3.1 Accuracy Analysis	42
3.3.2 Text Similarity Analysis	43
3.4 Chapter’s Deliverables	44
<b>Chapter 4 — Validate: Leading Expert-Blinded Validation of Retrieval Augmented Generation and Supervised Fine-Tuning for HCV Management</b>	<b>47</b>
4.1 Chapter’s Overview	48

4.2 Materials and Methods	49
4.2.1 Question and Clinical Cases Generation	49
4.2.2 Qualitative Answer Evaluation of Open-Ended Questions	50
4.2.3 Quantitative Answer Evaluation of Open-Ended Questions	51
4.2.4 Quantitative Evaluation of Clinical Scenario Treatment Recommendation	52
4.2.5 Model Configurations and Answer Generation	52
4.2.6 Statistical Analysis	55
4.3 Results	56
4.3.1 Optimal Number and Strategy of Text Chunking	56
4.3.2 Optimal Hyperparameters Tuning	56
4.3.3 Inter-grader Agreement	57
4.3.4 Evaluation of Accuracy and Clarity by Human Graders	57
4.3.5 Model Performance on Simulated Clinical Scenarios	58
4.4 Chapter’s Deliverables	60
<b>Chapter 5 — Verify: The Expert-of-Experts Verification and Alignment (EVAL) Framework for Safe, Aligned and Automatically Verified Outputs</b>	<b>63</b>
5.1 Chapter’s Overview	64
5.2 Materials and Methods	65
5.2.1 Large Language Model Configurations	65
5.2.1.1 Retrieval Augmented Generation	65
5.2.1.2 Supervised Fine-Tuning	66
5.2.2 Benchmark Datasets and Human-Grading	67
5.2.3 Unsupervised Similarity Metrics Alignment with Expert-of-Expert Golden Labels	70
5.2.4 Reward Model to Screen for High-Quality LLM Responses	74
5.2.5 Automated Rejection Sampling	76
5.3 Results	77
5.3.1 Model Ranking by Similarity Metrics	77
5.3.2 Model Ranking by Human Grading and Multiple-Choice Questions	81
5.3.3 Alignment between Similarity Metrics and Human Performance	81
5.3.4 Evaluation of Reward Model Alignment to Human-Grading	82
5.3.5 Rejection Sampling Across Multiple Temperature Thresholds	84
5.4 Chapter’s Deliverables	85
<b>Chapter 6 — Scale: From Single-Disease Prototyping to Society-Level Benchmarking: Guideline-Grounded Testing on European Association for the Study of the Liver Multiple-Choice Benchmark</b>	<b>89</b>
6.1 Chapter’s Overview	90
6.2 Materials and Methods	91
6.2.1 Dataset Creation and Model Comparison	91
6.2.2 Large Language Model Configuration	91
6.3 Results	93
6.4 Chapter’s Deliverables	94
<b>Chapter 7 — Reflect: Discussion and Conclusions: What Holds, What Does Not, and Where to Go Next</b>	<b>97</b>
7.1 Discussion	98
7.2 Conclusions	103
<b>Bibliography</b>	<b>104</b>



# Preface

Three years ago, I embarked on a doctoral journey that, as a physician, I knew would be demanding. It has been so; yet—as often happens—the greater the challenge, the more meaningful the achievement. This work sits at the intersection of clinical gastroenterology and computational methods, shaped by a question that became unavoidable after the public release of ChatGPT. Within weeks of launch, ChatGPT reportedly reached ~ **100 million active users in about two months**, a diffusion curve rarely seen for any information technology. That milestone did more than capture public imagination: it reframed expectations for everyday reasoning tasks and made a simple question pressing for healthcare—if general-purpose large language models (LLMs) can transform routine knowledge work so quickly, *what would it take to make them reliable, safe, and useful for clinical decision support?*

From the vantage point of a gastroenterologist, the promise is clear: LLMs appear to retrieve knowledge, synthesise evidence, and communicate complex plans in plain language. But healthcare is unlike general information work. It demands verifiable reasoning, traceability to authoritative sources, attention to subgroup nuance, and accountability for patient outcomes—alongside robust privacy safeguards and governance. Early experiments made the tension concrete: models that seemed persuasive in casual use could hallucinate, over-generalise, or omit contraindications in clinical scenarios. The question became not “*Can a model answer?*” but “*When, how, and under what guardrails should a model be allowed to answer at the bedside?*” This thesis is my attempt to turn that question into methods, evaluations, and practical guidance.

This dissertation is organised into seven chapters—beginning with *Chapter 1 (Introduction)*—that trace a coherent arc from early feasibility to rigorous validation and society-level scaling, reflecting a progression of thought and discovery throughout my PhD journey. The Introduction frames the landscape and states the central stance that binds the chapters together: make authoritative knowledge computable, validate with experts, and verify at scale under real-world constraints. It outlines what modern LLMs are, why digestive-disease workflows present unique challenges (clinical fidelity, medico-legal risk, privacy), and how prompt-engineering, retrieval-augmented generation or supervised fine-tuning can begin to constrain model behaviour—while also clarifying why naive, off-the-shelf use is insufficient in medicine.

Building directly on that foundation, *Chapter 2 (Baseline)* presents a systematic review of baseline LLM performance in gastroenterology and hepatology. Across multiple studies, accuracy on domain-specific questions is variable and often unsatisfactory; heterogeneity in question design, evaluator expertise, and grading rubrics further limits cross-study comparability. For clinical decision support, the conclusion is unambiguous: **baseline LLMs should not be presumed safe or sufficient for**

**healthcare applications.** This negative result—carefully documented—serves a positive purpose: it motivates the engineering and evaluation choices made in subsequent chapters.

With the problem space delineated, *Chapter 3 (Build)* proposes a scalable pathway to **embed evidence-based medicine (EBM)** into LLM workflows. Using chronic hepatitis C (HCV) guidelines as a testbed, we reformatted recommendations into an **LLM-friendly, machine-readable representation**, with explicit decisions about text standardisation, semantic structuring, retrieval granularity, and prompt schemata. The goal was not to “improve the model” in the abstract, but to **stabilise reasoning over complex recommendations**, including contraindications, stage-dependent treatments, and exceptions. Internal validation showed improved factual adherence, interpretability, and explainability, laying the engineering backbone for external scrutiny.

Design, however, is not evidence. *Chapter 4 (Validate)* turns design into externally credible results. We engaged senior European hepatology experts—the authors of the very guidelines used in Chapter 3—to perform **expert-blinded evaluations** of method and outputs. Their adjudication refined edge-case handling (e.g., treatment in comorbidity subgroups), clarified exception logic, and disciplined reporting formats. Beyond specific edits, the chapter demonstrates a process: **co-development with domain experts** to ensure that what works technically also holds clinically, and to make explicit where human oversight is essential.

With foundations and face-validity in place, the next step is to **make verification scalable and durable.** *Chapter 5 (Verify)* introduces an **automatic, rubric-based evaluation framework** for medical accuracy, developed with international experts in upper gastrointestinal bleeding. Crucially, experts authored **free-text reference answers**—rich in pathophysiology and management nuance—to carefully designed prompts. Programmatic comparisons against these references preserve clinical fidelity while enabling **large-scale benchmarking.** The framework naturally supports **post-deployment safeguards:** monitoring for performance drift, incident review when model outputs conflict with guidelines, and human-in-the-loop escalation pathways. Verification, here, is not a one-off score; it is a **mechanism for ongoing assurance.**

Finally, to examine how far these ingredients travel, *Chapter 6 (Scale)* applies the pipeline to a **clinical-society benchmark** (EASL multiple-choice questions). We evaluate multiple configurations—baseline, guideline-grounded RAG, and variations of reasoning strategies—against physician test-takers. Under carefully engineered, guideline-aligned conditions, models can **match or even surpass physicians on specific tasks**, while analysis clarifies the **boundaries where human oversight remains indispensable.** The chapter also discusses benchmarking caveats—leakage risks, and generalisability—emphasising why verification must accompany deployment.

The arc closes with *Chapter 7 (Reflect)*, which synthesises cross-chapter insights into a unified view of **clinical translation, limitations, and governance**. It distils what the evidence supports today versus what requires prospective, multi-centre validation; and sets a forward-looking research agenda for trustworthy, guideline-aligned AI in digestive diseases.

Collectively, these chapters trace a progression from **initial scepticism** about baseline models to **concrete strategies** for aligning LLMs with domain expertise: represent guidelines in computable form; engage expert validators; verify at scale with clinically meaningful rubrics; and surround deployment with privacy, safety, and governance guardrails. This progression also reflects my personal learning curve as a physician-researcher: how to translate clinical instincts (attention to exception logic, subgroup nuance, and safety) into computational artefacts (templates, retrieval design, grading rubrics) that others can reproduce and critique.

A few themes recur by design:

- **Clinical fidelity over surface correctness.** The thesis privileges alignment with authoritative guidance and explicit justification over persuasive prose.
- **Representation matters.** Making guidelines machine-readable—carefully and transparently—proved more impactful than ad-hoc prompt engineering alone.
- **Expert involvement is not optional.** Where models approach “expert-adjacent” performance, it is because experts shaped both the knowledge representation and the evaluation.
- **Verification is an ongoing practice.** Scores are snapshots; trustworthy systems require monitoring, incident response, and clear escalation pathways.
- **Privacy and governance set the envelope.** Choices about data, deployment, and evaluation are bounded by patient confidentiality and institutional accountability.

This dissertation is intended to be read **front-to-back**. The **Introduction** provides the conceptual primer and motivation. *Chapter 1* defines the problem and exposes failure modes. *Chapter 2* establishes the limitations of off-the-shelf models. *Chapter 3* presents the guideline-grounded engineering approach. *Chapter 4* supplies expert validation. *Chapter 5* introduces scalable verification and safety instrumentation. *Chapter 6* tests generalisation on a society benchmark and draws boundaries for responsible use. Figures mark key decision points (e.g., retrieval granularity), tables condense grading rubrics and error taxonomies, and appendices (where included) provide prompt templates and reference examples. Limitations are stated plainly: dependence on guideline quality, risks of distribution shift, constraints of retrospective datasets, and the persistent need for human adjudication. These are not afterthoughts; they **define the conditions under which AI assistance can be useful without being unsafe**.

In closing, while this dissertation collects individual studies, I hope it also functions as a **roadmap**: a way to combine clinical insight and computational design so that AI systems in digestive diseases are **guideline-aligned, expert-validated, and verifiably safe**. The rapid adoption of tools like ChatGPT set high expectations for what AI can do; the pages that follow are about **what AI should do, and how**, when the setting is the care of real patients.



# Chapter 1

## **Introduction: Foundations of Large Language Models in Digestive Diseases: Opportunities, Risks, and Privacy**

**“The map is not the territory”**

*– Alfred Korzybski*

## 1.1 The Importance of Large Language Models in Healthcare<sup>1</sup>

Large Language Models (LLMs) are transformer-based neural networks with billions of parameters trained on very large text corpora from diverse sources (articles, books, the internet).<sup>1</sup> LLMs are autoregressive, which means they predict the next part of the text based on the previous few words, using probabilities to determine the most likely continuation.<sup>1</sup> This architecture allows them to take in prompts, or inputs, in natural language and mimic human-like responses while also being able to accomplish a wide variety of complex tasks that were previously intractable.<sup>2</sup>

LLMs have the potential to overcome the productivity paradox in healthcare, characterized by the lack of generalized improvement in healthcare delivery by artificial intelligence (AI) based approaches, due to their user-friendliness and the fast pace of improvement cycles.<sup>3</sup> LLMs have already shown real potential in acting as virtual nurses for chronic disease support, proficiency in clinical note-taking, detection of drug adverse reactions, prediction of cancer metastasis, and evaluating the social determinant of health.<sup>4</sup>

These areas are relevant to Gastroenterology and Hepatology because LLMs can offer significant advancements in personalized patient care, by allowing closer automated follow-ups in patients with chronic conditions (e.g., diuretics management in patients with decompensated liver cirrhosis), managing drug adverse reactions (e.g., immunosuppressants in inflammatory bowel disease), retrieval of relevant note information and determine follow-up (e.g., imaging of pancreatic cysts based on their dimension or colonoscopy intervals based on polyps characteristics), and definition of health policies (e.g., increase chronic hepatitis C access to treatment in at-risk populations). However, a key challenge is the absence of established accuracy metrics, which are crucial to ensure the safe and effective application and deployment of LLMs in healthcare.

Beyond text-only applications, Large Multimodal Models (LMMs) are poised to transform digestive diseases workflows because it can operate across truly multimodal pipelines—endoscopic video, radiologic imaging, tabular electronic health records (EHR) data, and unstructured clinical notes—within a unified platform that spans documentation, decision support, and patient education.<sup>5</sup> Early evidence supports benefits mainly for lower-risk administrative and informational tasks (clinical documentation, billing, scheduling, literature summarization and patient education), whereas there is *no current*

---

<sup>1</sup>This chapter reproduces in full (text and images) the following article: *Giuffrè M., Kresevic S., Pugliese N., et al. Optimizing large language models in digestive disease: strategies and challenges to improve clinical outcomes. Liver Int. 2024 Sep;44(9):2114-2124. doi: 10.1111/liv.15974.* The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit.

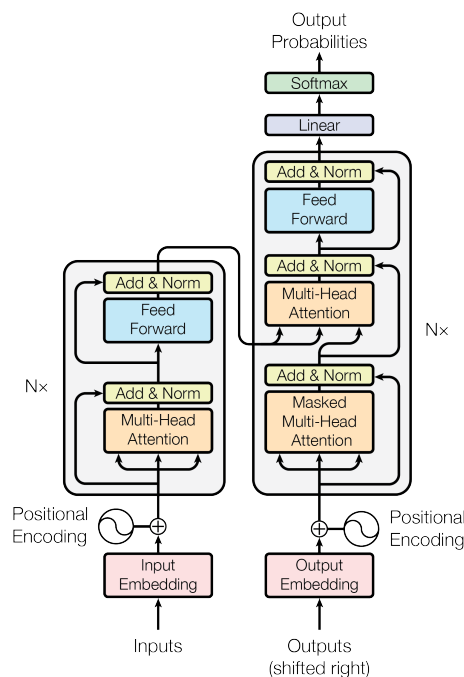
evidence that LLMs reliably performs complex clinical-reasoning tasks that would change diagnostic or treatment decisions and affect outcomes.<sup>5</sup>

## 1.2 The Transformer Architecture: Technical Foundations of Modern LLMs

To understand the capabilities and limitations of LLMs discussed throughout this thesis, it is essential to briefly examine their underlying technical architecture. Modern LLMs are built upon the Transformer architecture, introduced by Vaswani and et al. in their seminal work "Attention Is All You Need".<sup>6</sup> This architecture represented a paradigm shift in neural network design for sequence processing tasks, moving away from recurrent neural networks (RNNs) that had previously dominated the field.<sup>7,8</sup>

### 1.2.1 The Self-Attention Mechanism

The core innovation of the Transformer is the self-attention mechanism, which allows the model to weigh the importance of different words in a sequence when processing each word. Unlike RNNs that process text sequentially, self-attention enables parallel processing and captures dependencies between words regardless of their distance in the text. This mechanism computes attention by transforming input text into three representations—queries, keys, and values—and then calculating how much each word should "attend to" every other word through scaled dot-product operations (Figure 1.1).



**Figure 1.1: The Transformer Model Architecture.** The Transformer consists of an encoder (left) and decoder (right), each with  $N$  stacked layers. The encoder processes input through multi-head self-attention and feed-forward sub-layers with residual connections. The decoder adds a third cross-attention sub-layer connecting to encoder outputs, with masked self-attention ensuring autoregressive generation. Multi-head attention enables parallel processing across different representation subspaces by computing scaled dot-product attention with queries, keys, and values. This parallel architecture captures long-range dependencies more efficiently than sequential recurrent networks, enabling the scalability of modern billion-parameter Large Language Models. *Figure 1.1 is taken from Vaswani et al. "Attention is all you need"<sup>6</sup> – the authors have granted permission to reproduce the figures from the paper solely for scholarly works upon proper attribution provided by the authors.*

### **1.2.2 Encoder-Decoder Structure**

The original Transformer consists of an encoder that processes input sequences and a decoder that generates output sequences. Both components are composed of stacked layers (typically six in the base model), each containing multi-head attention sub-layers and feed-forward networks, with residual connections and layer normalization. The multi-head attention mechanism runs multiple attention operations in parallel, allowing the model to simultaneously focus on different aspects of the input, such as syntactic structure, semantic relationships, or contextual dependencies.<sup>6</sup>

### **1.2.3 From Transformers to Large Language Models**

While the original Transformer was designed for machine translation, modern LLMs like GPT (Generative Pre-trained Transformer) and similar models adapt this architecture primarily using the decoder component in an autoregressive manner—predicting one token at a time based on all previous tokens. The architecture incorporates positional encodings to maintain information about word order, since unlike RNNs, the attention mechanism itself has no inherent notion of sequence position. These models scale the original architecture dramatically, expanding from millions to billions or even trillions of parameters, and are pre-trained on vast text corpora before being fine-tuned for specific tasks.

### **1.2.4 Computational Advantages**

A key advantage of the Transformer architecture is its computational efficiency: self-attention layers connect all positions with a constant number of operations and enable significantly more parallelization compared to recurrent architectures that require sequential processing. This efficiency allowed the original Transformer model to achieve state-of-the-art translation quality while requiring only a fraction of the training time compared to previous models—a scalability property that has proven essential for training today's massive LLMs.

Understanding this technical foundation helps contextualize both the remarkable capabilities of LLMs in generating human-like medical text and their fundamental limitations. The attention mechanism's ability to capture long-range dependencies makes LLMs effective at maintaining coherent medical reasoning across extended clinical scenarios, yet the statistical nature of these operations means models may generate plausible-sounding but factually incorrect outputs—a phenomenon particularly problematic in high-stakes medical applications, as discussed in the following sections.

## **1.3 Challenges in Defining Accuracy for LLMs**

Gastroenterologists and hepatologists have sought to assess the accuracy of foundational LLM models in responding to many clinical and endoscopic questions.<sup>9–26</sup> Preliminary reports indicate a notable variability in accuracy, ranging from 25% to 90%, with an average near 50% for studies with larger question datasets.<sup>10,12,13,18,21,22</sup> This variability is influenced by several factors, including the lack of standardized reporting of both methodologies (e.g., model-related specifications, question generation

techniques, output evaluation criteria) and results. Creating common standards is still ongoing, as reported in a recent call to action, creating a cross-discipline initiative known as ChatGPT and Artificial Intelligence Natural Large Language Models for Accountable Reporting and Use (CANGARU) that is supposed to promote consensus on disclosure and guidance for reporting LLM use in academic research.<sup>27</sup> Despite the current absence of standardized reporting, one of the critical areas yet to be addressed is related to a widely accepted definition of output “accuracy”. In the context of LLMs, some digestive diseases researchers have explored a binary dimension of accuracy<sup>10,12,13,18,21,22</sup> (completely accurate versus all else), while others adopt a scalar approach (e.g., using Likert scales).<sup>9,19</sup> In addition, it appears that the definition of accurate outputs is often bound to grader expertise, without a set of well-established criteria and rules dictated by a juxtaposition to a given gold standard, such as the recommendations reported in current clinical guidelines.

Despite the ongoing efforts to define standards in reporting and accuracy, the essential concern remains in the context of the real-world application of foundational LLMs, whose performance in general medicine and high-yield medical specialties remains unacceptable, where accuracy is not just a metric but a matter of patient safety with decisions that can have life-or-death consequences.<sup>28–33</sup>

The variability in LLM accuracy and the absence of standardized metrics highlight a key issue within the digital health domain. We focus our discussion pragmatically on the application of strategies that have been used to improve LLM outputs rather than the definition of accuracy.

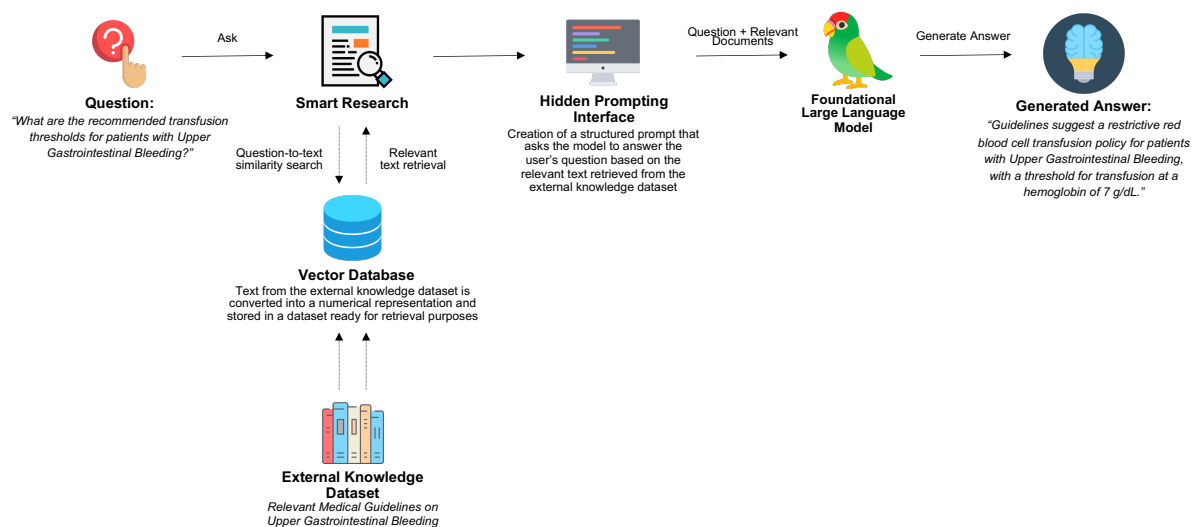
## **1.4 How to Infuse Domain Knowledge**

Confronted with the critical challenge of the unacceptable accuracy of foundational LLMs, as experts in digestive diseases, we must urgently seek strategies to refine LLMs' accuracy while ensuring that their outputs align with the expert knowledge. To bridge this gap, we must participate in efforts to challenge and evaluate foundational LLMs rigorously. Current initiatives in the LLM industry include red teaming, which are communities of experienced domain experts that can help LLM companies with risk assessment and guide them towards mitigation efforts. These efforts are resource-intensive and limited to specific companies and their LLM products; more accessible approaches that are critical for providers to integrate digestive disease domain knowledge into foundational LLMs include in-context learning and supervised fine-tuning (SFT) followed by reinforcement learning with human feedback (RLHF).

### **1.4.1 In-Context Learning**

The first approach uses a technique called Retrieval Augmented Generation (RAG).<sup>34</sup> RAG can be performed with any LLM and is a variant of prompt engineering, which is a sequence of task-specific instructions, or prompts, to optimize LLMs output.<sup>35</sup> A “prompt” is a sequence of text that outlines the

task for LLMs to perform<sup>1</sup>, with several types of prompting strategies (Table 1.1) that can be used together or individually for optimally completing the desired task.<sup>36</sup> A core part of RAG is the context that is provided as part of the prompt engineering strategy. In gastroenterology and hepatology, the relevant text for clinical decision-making may be the relevant national guidelines, which could be provided as part of the prompt engineering for the model to be preferentially accessed with each new query (Figure 1.2). For example, consider a common scenario in gastroenterological practice, such as upper gastrointestinal bleeding (UGIB).<sup>37</sup> If a query about transfusion thresholds is posed (i.e., “Should I transfuse a patient with suspected UGIB and hemoglobin of 9.5 g/dL?”), the new prompt, incorporating the clinical question visible to the user, is queried together with the clinical guidelines (hidden from the user), and the relevant portion of the text is retrieved to create the output (e.g., “The guidelines suggest a restrictive policy of red blood cell transfusion threshold at hemoglobin of 7 g/dL in patients with UGIB”).



**Figure 1.2: Retrieval Augmented Generation.** The foundational model retrieves information from an external information dataset and binds its response generation to the relevant information contained in the context and the instruction provided by the user. The retrieval strategy of relevant information depends on its length compared to the maximum context window of the model. If the relevant information length is longer than the context window, it must be segmented into chunks, tokenized (split into small foundational units), embedded (converted into numerical vector representations), and stored in a vector database. The chunk with the highest similarity to the user’s query is selected for each new interaction, and the response is provided according to the query and the retrieved chunk. Chunking is unnecessary if the relevant information length is shorter than the context window. Instead, the entire relevant context is passed to the model in one pass, ensuring a comprehensive and uninterrupted text generation.

Prompt Type
<p><b>Prefix Prompt:</b> The LLM is given a direct instruction or question as a "prefix" to the information the user is seeking (i.e., giving a command or asking a question that the LLM should answer).</p> <p><b>Query:</b> Context [X]: A patient presents with melena and hematemesis. Given context [X], determine the meaning of the abbreviation [UGIB].</p> <p><b>Answer:</b> The abbreviation [UGIB] stands for Upper Gastrointestinal Bleeding in the context [X].</p>
<p><b>Cloze Prompt:</b> The LLM is presented with a sentence or a paragraph with a blank (cloze) that needs to be filled in.</p> <p><b>Query:</b> Context [X]: A patient with cirrhosis is experiencing vomiting of blood. The abbreviation [ ] stands for [Variceal Hemorrhage] in the context [X].</p> <p><b>Answer:</b></p>

The abbreviation [VH] stands for [Variceal Hemorrhage] in the context [X].
<b>Anticipatory Prompt:</b> The LLM is expected to anticipate or predict something based on the context provided (it is often used when the LLM is supposed to understand nuance or read between the lines).
<p><b>Query:</b> Context [X]: A patient presents with melena. Abbreviation [Y]: NSAID What does the abbreviation [Y] mean in context [X]?</p> <p><b>Answer:</b> In the context of a patient presenting with melena, the abbreviation "NSAID" stands for Nonsteroidal Anti-Inflammatory Drugs.</p>
<b>Chain-of-Thought Prompting:</b> The LLM is asked to explain step by step how to reach a given conclusion (it is helpful when it is required to define the reasoning behind an answer or when the path to the answer is as important as the answer itself).
<p><b>Query:</b> Context [X]: A patient presents with hematemesis. Abbreviation [Y]: EVs Explain step-by-step the meaning of the abbreviation [Y] in context [X].</p> <p><b>Answer:</b> Step-by-Step Explanation:</p> <ol style="list-style-type: none"> <li>1. Identify the Symptom: Hematemesis indicates bleeding in the upper gastrointestinal (GI) tract, which includes the esophagus, stomach, and the initial part of the small intestine.</li> <li>2. Abbreviation Analysis: The abbreviation "EVs" in the medical context often stands for Esophageal Varices.</li> <li>3. Understanding Esophageal Varices (EVs): Esophageal Varices are swollen veins in the lining of the lower esophagus near the stomach. They typically result from increased pressure in the portal vein, known as portal hypertension, and are most often associated with liver diseases like cirrhosis.</li> <li>4. Connection to Hematemesis: When EVs rupture, they can cause significant bleeding within the esophagus, which can lead to hematemesis. The blood vomited may be bright red or have a coffee-ground appearance, depending on how long the blood has been in the stomach.</li> </ol>
<b>Heuristic Prompt:</b> The LLM is asked to apply a rule or a set of rules to generate a response (it is useful where there is a set of established guidelines or principles that should be followed to reach a solution). Unlike chain-of-thought prompts which explain the reasoning or logic behind the output, heuristic prompts use a set of predefined rules to guide the LLM to perform a task.
<p><b>Query:</b> Rule: Upon assessing a patient with UGIB, calculate the GBS. When a patient with suspected upper gastrointestinal bleeding has a GBS <math>\leq 1</math>, he/she should be considered for outpatient management or early discharge given the low risk of serious bleeding or need for intervention. Should a patient with suspected upper gastrointestinal bleeding be discharged with a GBS = 0?</p> <p><b>Answer:</b> Yes, the patient should be considered for outpatient management or early discharge. A GBS of 0 falls within the threshold (GBS <math>\leq 1</math>) indicating a low risk of serious bleeding or need for intervention. It suggests that the patient is at low risk for complications and may not require in-hospital treatment or urgent endoscopic intervention.</p>
<b>Ensemble Prompt:</b> The LLM is asked to perform a task combining multiple prompts using majority voting (i.e., most common output) for aggregated outputs.
<p><i>Heuristic Prompt</i></p> <p><b>Query:</b> Rule: Patients with suspected UGIB and a GBS score of <math>\leq 1</math> are generally considered low risk and may be candidates for outpatient management or early discharge. Scenario: A patient with suspected UGIB has a GBS score of 1. Should this patient be considered for outpatient management or early discharge?</p> <p><b>Answer:</b> The rule suggests that the patient should be considered for outpatient management or early discharge.</p> <p><i>Chain-of-Thought Prompt:</i></p> <p>Query: A patient with suspected UGIB has a GBS score of 1. Should this patient be considered for outpatient management or early discharge? Explain your decision step by step:</p> <p><b>Answer:</b></p> <ol style="list-style-type: none"> <li>1. GBS Score Interpretation: A score of 1 in GBS indicates a low risk of serious bleeding or a need for intervention.</li> <li>2. Patient Stability: The patient's vital signs and overall clinical stability are essential. A stable patient with a low GBS score is a good candidate for outpatient management.</li> </ol> <p>The patient can be discharged if he or she is clinically stable, has adequate support and follow-up, and has no significant comorbid conditions.</p> <p><b>Outcome:</b> Both the heuristic and chain-of-thought approaches lead to the conclusion that a patient with suspected UGIB and a GBS score of 1, who is clinically stable and has adequate support for outpatient management, can be considered for early discharge. Therefore, following this ensemble approach, the patient in this scenario can be safely managed outside the hospital setting, with appropriate follow-up care planned.</p>

**Table 1.1:** Prompting in Large Language Models with Gastroenterology-Related Examples on upper gastrointestinal bleeding. Abbreviations: UGIB, upper gastrointestinal bleeding; GBS, Glasgow-Blatchford Score.

Preliminary data suggests that in-context learning can produce more accurate outcomes when applied to digestive disease (RAG-enhanced GPT-4, 79% vs. foundational GPT-4, 50.5%) in the context of colorectal cancer screening follow-up.<sup>38</sup> Despite the encouraging results, in-context learning has several challenges: (1) it is dependent on the ability of LLMs to process large amounts of new data (defined as the context window), (2) it does not compensate for changes in the base foundational model (e.g., update from GPT-3.5 to GPT-4), (3) it can be affected by context semantics and word repetition, (4) there is no standardized prompt architectures available for benchmarking individual approaches, and (5) it cannot consistently handle information stored in tables and figures.

**Challenge 1: LLMs do not process large amounts of data similarly, and in-context learning may perform differently when the underlying model is updated.**

The ability of LLMs to process prompts is as follows: the entered text is broken down into smaller fundamental units of analysis called tokens. For the tokenized text, an embedding vector is obtained that can be further processed<sup>28</sup>. An embedding is a mathematical representation of arbitrary types of data as high-dimensional vector designed to capture meaningful characteristics in a lower dimensional space.<sup>39</sup> When embeddings are generated for texts in the case of LLMs, they are meant to capture the semantic and syntactic information and condense it into a form that allows for efficient retrieval of relevant information.<sup>39</sup> Embeddings are accessed in different ways using prompt engineering strategies, leading to a field of expertise in designing the optimal prompt architectures to maximize the usage of embeddings. Each LLM operates with a maximum number of tokens for its context window, which is defined as the maximum acceptable number of tokens in a single forward interaction.<sup>40</sup> For example, the initial version of GPT-3.5 Turbo has a maximum context window of 4096 (~3000 words) tokens, whereas the advanced GPT-4 Turbo is capped at 128000 (~96000 words) tokens. This limitation indicates that if the information needed for RAG exceeds the token limit, it cannot be used fully, and therefore the original text must undergo a process of chunking, defined as the segmentation of the original text into smaller portions based on a set of rules, such as dividing by paragraphs or creating segments at the sentence level.<sup>41</sup> Each text segment is then tokenized, converted into an embedding vector, and stored in a vector database. For each new query, the system selects the text segment that most closely matches the entered query based on quantitative estimation of textual similarity.<sup>41</sup> Using more advanced models with larger context windows might negate the need for chunking with a single guideline text. However, due to text length, chunking is still required when applying RAG across multiple guidelines. While having a larger context window may solve the information retrieval problem, it may result in additional costs for proprietary models for large-scale deployment. For example, if we consider deploying RAG-enhanced OpenAI's GPT-4 Turbo with the Hepatitis C Virus guidelines from the European (approximately 33000 words and 41000 tokens) Association for the Study of The Liver, a 10-message interaction chat would cost approximately 4 USD, which could limit large-scale deployments. This limitation calls for the strategic

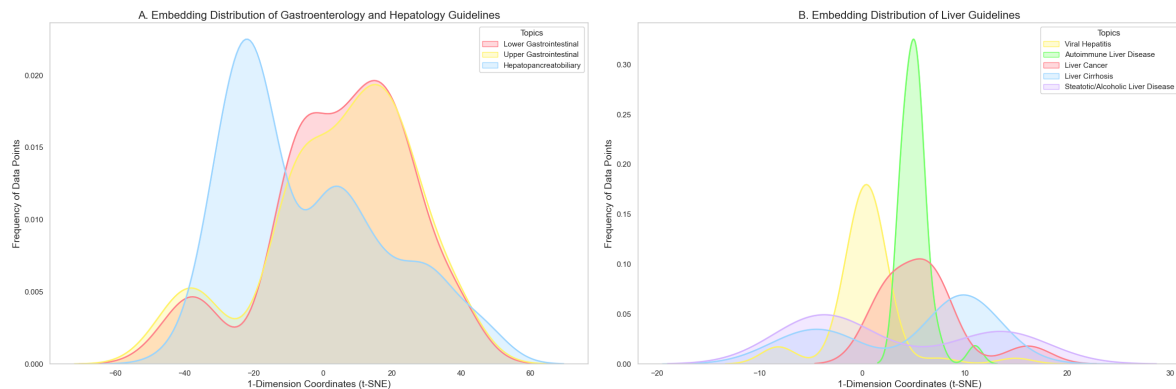
imperative of judicious token utilization to maintain fiscal prudence within the clinical decision support systems.

**Challenge 2: Semantics and word distribution affect LLM outputs.** The necessity for chunking can lead to less accurate responses since it may fetch less or irrelevant text sections, particularly when the embedding space operates within a confined semantic field where specialized terms are densely interconnected, as described in Figure 1.3. The effectiveness of RAG hinges on the quality of these embeddings, as they determine how well the retrieval mechanism can match a user's query with the appropriate segment of text within the national guidelines or other clinical resources. Advanced embedding techniques that consider context and semantic meaning rather than simple word frequency can significantly enhance the accuracy of the retrieved information, thereby improving the model's output for clinical decision-making. With the idea of creating a multi-guideline digestive disease RAG-enhanced LLM, when seeking information about “esophageal variceal hemorrhage”, chunking could mistakenly prioritize text segments based on word frequency and proximity and may provide answers on “gastroesophageal reflux” or “esophagitis” rather than focusing to specific terms on “esophageal variceal hemorrhage” that are critical for an accurate result. Besides, in the case of queries addressing knowledge contained in separate text segments, the retrieved chunks may not contain the whole information needed to address the query (e.g., complex clinical cases). Additionally, in the case of RAG with guidelines containing contrasting information (e.g., different recommendations regarding usage of pre-endoscopic infusion of erythromycin in UGIB among different guidelines), chunking retrieval may result in answers related to only one pertinent text source and not report all possible recommendations. Furthermore, the success of RAG is contingent on the choice of effective text embedding models that can accurately capture semantic similarities across various texts, introducing an additional layer of complexity to its successful implementation.

**Challenge 3: No benchmarks or standardized prompt architectures exist for gastrointestinal disease.** The absence of a unified framework for prompt engineering means that each study or application may use a different approach to elicit responses from LLMs. This variability can lead to inconsistencies in performance and difficulty in comparing results across different studies or implementations. For instance, in the context of digestive disease, only two studies<sup>12,18</sup> have clearly employed prompt-engineering strategies to guide the LLM to perform a certain task (e.g., define the next steps in patient management or provide a brief clinical summary). Without standardized prompts and tasks, it is challenging to gauge whether variations in the LLMs' performance are due to differences in the underlying model or the way the prompts are structured.

**Challenge 4: Information contained in figures and tables is not consistently and robustly available.** An often-neglected aspect of RAG-based approaches is the optimal formatting of guidelines

to ensure high-quality, coherent text within and across multiple documents from various medical societies. Additionally, the capabilities of recent LLMs in processing non-text sources of information, such as graphical tables or flowcharts commonly found in clinical guidelines, remain to be fully explored and utilized, with a recent study in the context of digestive disease reporting accuracy of 16% in parsing information from graphical tables and flowcharts with GPT-4 Vision.<sup>42</sup>



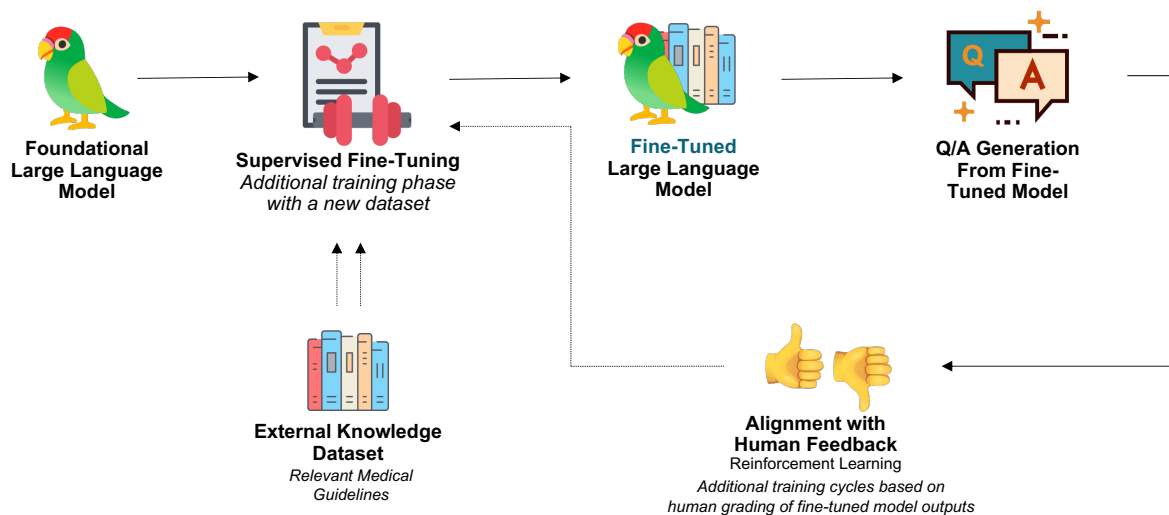
**Figure 1.3: Visualization of embedding space distributions for text contained in digestive disease guidelines.** These visual representations have been generated by employing a dimensionality reduction technique, specifically a one-dimensional t-distributed Stochastic Neighbor Embedding (t-SNE), executed on the embedding vectors derived from the textual content of clinical guidelines provided by the European Society of Gastrointestinal Endoscopy (ESGE), the European Association for the Study of the Liver (EASL), and the United European Gastroenterology (UEG). The embeddings were calculated using the 'text-ada-002' model from OpenAI, facilitating an advanced, nuanced analysis of the multi-dimensional semantic and syntactic relationships within the clinical guideline texts. A. The guidelines are collectively labeled and subsequently categorized into three principal domains: Lower Gastrointestinal, Upper Gastrointestinal, and Hepatopancreatobiliary. Notably, this representation exhibits a discernible overlap between the Lower and Upper Gastrointestinal guidelines, indicating a potential intersection of topics within these domains. B. Exclusively dedicated to visualizing Liver Guidelines without incorporating the broader spectrum of digestive diseases. While this focused approach permits a detailed examination of liver-specific content, overlaps between topics within this category are evident. These overlaps highlight the intrinsic complexity of the subject matter and underscore the imperative for precise differentiation within the vectorial embedding space to ensure accurate information retrieval and minimize the potential for ambiguity in clinical decision support.

#### 1.4.2 Supervised Fine-Tuning with Reinforcement Learning from Human Feedback

While RAG directs LLMs to use only certain parts of its architecture to fit the domain-specific question, SFT with RLHF aims at inducing persistent changes in the LLM's internal architecture with alignment from human feedback<sup>1</sup>. SFT involves transfer learning<sup>43,44</sup>, a technique that enables the leveraging of a pre-existing model to solve new problems more efficiently by adjusting the weights and parameters (e.g., changing the number of inner layers or adjusting the learning rate) of an already-trained LLM model, using a high-fidelity domain-specific dataset (Figure 1.4). Subsequently, SFT must be followed by alignment using RLHF, where reward modeling (RM) and reinforcement learning (RL) based on human feedback are used to align the model toward human preferences<sup>45</sup>. The desired phenomenon is that responses with low rewards from human users are less likely to be repeated, incentivizing the LLM to produce outputs that more closely align with human expectations.

While SFT with RLHF can provide fundamental tailoring to specific knowledge domains, there are several challenges: it is computationally expensive, inherently complex, and requires laborious human expert annotation. For open-source LLMs such as Meta’s Llama-2, the primary cost is computational infrastructure, specifically in terms of Graphics Processing Units (GPUs) or, more accurately, GPU hours. As larger models cannot be fine-tuned on standard local setups and require access to High-Performance Computing (HPC) or specialized GPU hosting services, the cost of computing can be significant. For instance, the fine-tuning of medAlpaca<sup>46</sup>, a fine-tuned model based on the Llama-2 7B architecture, required between 72.5 and 387.5 GPU hours at an estimated total cost of 3,000 USD. In contrast, closed-source LLMs like OpenAI’s GPT-4 manage the fine-tuning process on their platforms, which can reduce the direct costs associated with GPU hours. However, users are charged based on the number of tokens processed during the fine-tuning phase (8 USD per million tokens for GPT 3.5 Turbo) and for generating content from the fine-tuned model (6 USD per million tokens for GPT 3.5 Turbo), with the added caveat that the fine-tuned model remains the property of the service provider. This does not include specialized human labor, which includes both engineering expertise and domain expertise. There are no existing studies to our knowledge using LLMs with SFT in gastroenterology and hepatology; however, a recent study published in Nature Medicine suggests that the use of SFT in a clinical summarization task resulted in summaries that were judged to be either equivalent (45%) or superior (36%) to those produced by the medical experts.<sup>47</sup>

RLHF is commonly used to align LLMs to human values, ethics, preferences, and goals.<sup>48,49</sup> For medical applications, an additional challenge is the nature of how the responses should be aligned, not just to specific domain experts but also according to human intentions and values consistent with medical ethics. Specifically, a model can be defined as “aligned” with human values and medical ethics if it fulfills three criteria for helpful, honest, and harmless (HHH).<sup>50</sup> Regarding goals, fine-tuned LLMs can be aligned to address several clinical tasks. For example, LLMs can be fine-tuned to support preliminary diagnosis and patient management by providing initial guidance based on symptoms described by patients.<sup>51</sup> Another goal could be the development of systems capable of summarizing extensive medical notes into concise reports<sup>47</sup>, aiding clinicians in quickly grasping patient histories and treatment plans. These applications aim to enhance the efficiency and accuracy of medical services, support healthcare providers, and potentially improve patient outcomes.



**Figure 1.4: Fine-Tuning.** Initially, the foundational LLM undergoes retraining on domain-specific knowledge, which utilizes transfer learning to adjust the model's parameters using a tailored dataset relevant to the medical field. After the initial fine-tuning, the model enters an alignment phase, which is essential to ensure that the model's outputs are accurate and adhere to human values and intentions. The final stages involve query-answer testing to assess the model's performance and reinforcement learning from human feedback (RLHF). During RLHF, the fine-tuned model is further refined through reward modeling based on human evaluation of the model's outputs, ranging from best to worst or correct/incorrect. This process iteratively adjusts the model's behavior to incentivize the generation of outputs that align more closely with human expectations.

## 1.5 Possible Applications in Gastroenterology and Hepatology

In exploring the application of LLMs in gastroenterology and hepatology, it is beneficial to differentiate between the enhancements offered by SFT with RLHF and those by RAG. While RAG optimizes LLM output by retrieving and incorporating relevant information from medical texts and guidelines, SFT with RLHF goes further by deeply integrating domain-specific knowledge and enabling the LLM to learn from interaction outcomes, a process akin to a medical student learning from both textbooks and clinical feedback.

Across the field of digestive diseases, we are engaged in screening, treatment, and surveillance for a broad range of pathologies for which RAG or SFT with RLHF could improve clinical decision-making and patient outcomes.

**Screening Colonoscopies.** A RAG-enhanced LLM might retrieve the latest guideline on polyp surveillance intervals when prompted with a question about a patient with a specific polyp size and histology. Beyond retrieval, an LLM with SFT and RLHF would learn the specific practice pattern of a group of providers due to information that is not available in the published guidelines. Local best-practice patterns can then be used to tailor the recommendation and respond with how the gastroenterologist would counsel a patient in that specific context, and not simply what the guidelines state. For example, in an underserved area with poor access to colonoscopy, if one were asked about the next surveillance interval for a patient with 3 <10mm tubular adenomas at baseline colonoscopy and 3

<10 mm tubular adenomas at the current colonoscopy, an RAG-enhanced LLM may respond “According to current guidelines, surveillance colonoscopy suggested in 7-10 years”; an LLM with SFT and RLHF may respond with additional information such as “Given the fact that this patient did not enter the screening program at the appropriate age due to barriers to access, a surveillance colonoscopy by guidelines is suggested at 7-10 years, but an earlier test can be considered”. By integrating relevant guidelines into physician feedback, LLMs can derive more sophisticated and personalized patient care plans tailored to the needs of specific community contexts.

**Chronic Hepatitis C (HCV) Treatment.** Chronic HCV is a prime example for which we have very good treatment regimens, but recommendations are often too complex and challenging to navigate. When asked about treatment options for a patient with chronic HCV who has stable HIV co-infection, and managed well on antiretroviral therapy without protease inhibitors, a RAG-enhanced LLM may respond, “Current guidelines suggest that for HCV/HIV coinfecting patients without cirrhosis, a pangenotypic DAA regimen such as glecaprevir/pibrentasvir for 8 weeks could be recommended.” However, given that this patient comes from an area with high rates of HCV reinfection and limited follow-up resources, closer monitoring and patient education on adherence and transmission risk reduction should be emphasized. In such cases, an LLM with SFT and RLHF may add “We may also consider extending the treatment duration to 12 weeks to account for these factors, despite the patient’s non-cirrhotic status, to ensure a higher barrier to treatment failure and enhance the probability of sustained virological response.”

**Pancreatic Cyst Surveillance.** For surveillance, pancreatic cysts require constant monitoring for changes that may suggest worrisome or high-risk features. A RAG-enhanced LLM could retrieve relevant information for imaging follow-up determination based on the cyst's size and appearance, but depending on the availability of radiological or endoscopic ultrasound (EUS)-operator expertise, the recommendation may differ. For example, for a patient with worrisome features and a cyst size > 3cm, a RAG-enhanced LLM may give the following “According to the guidelines, referral to a multidisciplinary group and alternating MRI/EUS every 6 months for three years is recommended.” An LLM with SFT and RLHF may respond: “Given the reduced MRI accessibility of this region to MRI and the risk of losing patients to follow-up, surveillance can be efficiently achieved with EUS given appropriate EUS-MRI agreement on cyst features.”

These examples illustrate how SFT with RLHF can offer additional benefits beyond RAG, by training the LLM to integrate complex, multifaceted medical data reflective of real-world clinical decision-making. While RAG improves the LLM's ability to access and utilize vast information stores, SFT with RLHF teaches it to apply that information in a context-sensitive and nuanced manner, akin to the clinical judgment that gastroenterologists develop through practice.

## 1.6 Ethical Considerations

In the integration of LLMs into patient care, four foundational ethical principles guide our approach: autonomy, ensuring informed self-determination; beneficence, committing to the patient's best interests; non-maleficence, avoiding harm; and justice, guaranteeing equitable access to healthcare advancements.

**Autonomy.** The principle of autonomy emphasizes a patient's right to self-determination. In healthcare, this translates into patients making informed decisions about their treatment options. LLMs could hold the power to sway decisions by framing recommendations that appear authoritative, even when not endorsed by an actual medical expert. This can lead to an over-reliance on LLMs and potentially undermine patient autonomy. It is imperative to ensure that LLMs do not supplant the patient-physician interaction but rather serve to enhance it. Patients must be made aware that the advice provided by LLMs complements, rather than replaces, the nuanced counsel of a healthcare professional.

**Beneficence.** The principle of beneficence requires acting in the best interest of the patient. Companies that develop and deploy LLMs as closed-source models, like OpenAI, operate under commercial imperatives. This reality prompts the question of whether the primary goal is patient welfare or the pursuit of profit. While LLMs can offer extensive knowledge, ensuring that they prioritize patient benefit over other interests is a moral imperative. To uphold this ethical pillar, developers must align model outputs with the ethos of enhancing patient care, ensuring that beneficence remains central to LLM deployment in healthcare.

**Non-maleficence.** The principle of "first, do no harm" is a cornerstone of medical ethics. LLMs in their current state can be susceptible to manipulation or 'jailbreaking,' potentially leading to harm through misinformation. As such, it is crucial to develop and implement robust safeguards against such vulnerabilities. LLM systems must be designed with fail-safes that prevent the propagation of harmful or incorrect medical advice.

**Justice.** Equal treatment and non-discrimination are critical components of justice. Access to LLMs could inadvertently exacerbate disparities in healthcare by providing high-quality information only to those with digital literacy or adequate technology or based on the training data used for LLM development. Efforts must be made to ensure equitable access to these tools so that all patients, regardless of socio-economic status or geographic location, can benefit from the advancements in AI-assisted healthcare.

## **1.7 Privacy and Liability**

The utilization of LLMs as medical chatbots must consider the ethics of utilizing personal health data and the potential liability associated with LLM-generated responses for clinical decision making. Proprietary LLMs like ChatGPT do not ensure comprehensive management of personal data in accordance with the Health Insurance Portability and Accountability Act (HIPAA), potentially jeopardizing data security and confidentiality. Hence, employing ChatGPT or comparable non-compliant systems for the management of sensitive health information is deemed unsafe and not recommended.<sup>52</sup> Another crucial concern pertains to liability, as there is significant variability in accuracy across different models and settings. Within conventional healthcare environments, healthcare professionals and institutions are typically held accountable according to established legal frameworks and professional standards.<sup>53</sup> Nevertheless, incorporating LLMs as intermediaries in the provision of patient care may introduce complexities to this situation. Since LLMs such as ChatGPT are created and managed by technology companies rather than healthcare providers, the boundaries of responsibility become less clear. When a patient adheres to erroneous guidance provided by an LLM and experiences adverse consequences, it becomes difficult to determine the appropriate party responsible for liability. This responsibility may lie with the technology provider, the healthcare professionals who incorporated the LLM into the care process, or a combination of both. Additionally, the matter of consent and informed decision-making must be considered. Patients should be cognizant of including an LLM in their care regimen and comprehend the associated risks. The importance of transparency in ethical practice cannot be overstated. However, it also raises inquiries regarding patients' perceptions of advice provided by LLMs and the level of trust they place in this information when making health-related choices. To address liability concerns, it may be necessary to advocate for regulatory supervision, rigorous testing and validation of LLMs for healthcare applications, explicit standards regarding their involvement in patient care, and robust mechanisms for monitoring outcomes and errors.

## **1.8 Overall Aims and Objectives**

Having outlined the clinical promise, methodological gaps, and ethical constraints of LLMs in digestive diseases, the thesis now translates those insights into a sequenced program of objectives. The overarching aim is to evaluate, optimize, and safeguard the use of LLMs for clinical decision support in digestive diseases – moving from baseline performance assessment to guideline-aligned outputs and benchmarking against human experts. To accomplish this, we proceed in incremental steps that mirror the chapters of the thesis. The first objective is to establish a rigorous baseline by quantifying current LLM performance across representative digestive-disease tasks and characterizing error modes through a systematic review of the literature (Chapter 2). Noting that off-the-shelf models are insufficient for high-stakes medical use, the second objective is to align model outputs with evidence-based medicine by integrating authoritative clinical guidelines via retrieval-augmented generation and principled

prompt design, and by testing “LLM-friendly” guideline formulations to improve interpretability and faithfulness (Chapter 3). The third objective is to translate these methods to point-of-care use—initially in hepatitis C virus (HCV) infection—optimizing retrieval granularity and reasoning strategies and obtaining external validation from internationally recognized experts, including European HCV guideline authors, on realistic clinical cases (Chapter 4). Recognizing the cost and scarcity of expert labeling and the need for robust guardrails, the fourth objective is to safeguard deployment with an expert-informed safety layer that could automatically grade and act as a filter for risky or misleading outputs using similarity to gold standards, and reward modeling (Chapter 5). Finally, the fifth objective is to benchmark the resulting systems against physicians on standardized assessments—including image-based questions—to determine where LLMs meet, exceed, or fall short of human performance (Chapter 6). Within this sequence, the primary aim is reliable, guideline-faithful decision support; the staged objectives operationalize and validate this aim from baseline appraisal to head-to-head human benchmarking.



## **Chapter 2**

### **Baseline: What General-Purpose Large Language Models Can (and Cannot) Do in Digestive Diseases: A Systematic Review of Performance Evaluation**

**“If you can’t measure it, you can’t improve it.”**

*– Peter Drucker*

## 2.1 Chapter Overview<sup>2</sup>

This chapter derives from a registered, PRISMA-guided systematic review of LLMs applied to gastroenterology and hepatology. The review synthesizes published evidence to assess accuracy and potential safety implications of LLM outputs for diagnosis, management, and treatment in digestive diseases, and to identify methodological gaps that must be addressed before safe clinical adoption.

The primary aim is to establish a rigorous baseline of LLM reliability in digestive diseases by quantifying accuracy (reported as proportion of *completely correct* answers) and describing safety-relevant failure modes across published studies.

Specific objectives include:

1. Map the literature: Systematically search 7 databases (plus citation chasing), screen per inclusion criteria, and extract study aims, model versions, prompt strategies, question sets, grading schemes, and evaluator profiles.
2. Quantify accuracy: Report the proportion of *completely correct* answers, allowing comparisons by model/version, topic (GI, liver, pancreas), question type (general vs. clinical cases), and perspective (patient vs. physician).
3. Characterize methodology & bias: Appraise study quality with bias tools; document heterogeneity (question generation, number of items/graders, grading scales), and reasons meta-analysis was not performed.
4. Analyze question space: Build a unified corpus of available questions and visualize clusters via t-SNE by study, topic, type, and perspective.
5. Document safety-relevant behaviors: Describe error modes (including hallucinations) and highlight recommendations that could affect patient safety (e.g., hospitalization bias, imaging/follow-up suggestions).
6. Clarify reporting standards needed: From observed variability, delineate gaps and propose the need for standardized accuracy definitions, evaluator reporting, and question transparency.

---

<sup>2</sup>This chapter (text and images) is adapted from the article: *Giuffrè M., Krešević S., Kisung Y., et al. Systematic review: The use of large language models as medical chatbots in digestive diseases. Aliment Pharmacol Ther. 2024 Jul;60(2):144-166. doi: 10.1111/apt.18058.* The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit. Changes: Methods and Results are reproduced verbatim; Introduction, Background, Discussion and narrative transitions have been reformulated for thesis style; figure/table numbering and layout adapted to the thesis format.

## 2.2 Materials and Methods

The Preferred Reporting Items for Systematic Reviews and Meta-analyses (PRISMA)<sup>54</sup> and Synthesis Without Meta-Analysis (SWiM)<sup>55</sup> were used to guide the reporting in this study. The protocol was registered *a priori* (CRD (PROSPERO ID:42023476782)).

### 2.2.1 Search Strategy and Selection Criteria

A systematic search of the literature was conducted by a medical librarian (A.A.G) in the following databases: Cochrane Library, Google Scholar, Ovid Embase, Ovid MEDLINE, PubMed, Scopus, and Web of Science Core Collection to identify relevant articles published from inception of each database until October 9, 2023 using a combination of keywords and controlled vocabulary for large language models and digestive diseases. Details of the full search strategies are listed elsewhere.<sup>56</sup> The search was not limited by language, publication type, or year. A second medical librarian peer reviewed the search by using the Peer Review of Electronic Search Strategies (PRESS)<sup>57</sup>. Using Citation Chaser<sup>58</sup>, we searched the reference lists of included studies and retrieved articles that cited the included studies with the goal of finding additional relevant studies not retrieved by the initial database search.

The search results from all databases were imported into Endnote 20, then duplicates were removed using the Yale Reference Deduplicator.<sup>59</sup> The deduplicated results were imported into Covidence for screening. Two independent screeners (M.G., J.D.) performed a title and abstract review, with a third screener (D.L.S.) to resolve disagreements. The full texts of the resulting articles were then independently reviewed for inclusion by 2 screeners (M.G., J.D.) and a third screener (J.H.) to resolve disagreements. For manuscripts to be considered for inclusion in this systematic review, they must have undergone peer-review and meet the following criteria: (1) the study must query LLMs related to the disease-related context in the field of gastroenterology and hepatology; (2) language: the article must be written in English and questions/answers from LLMs should be provided and analyzed in English; (3) the study should define the criteria and methodology used to assess the answers provided by LLMs. Manuscripts that did not meet all of the above criteria will be excluded from this systematic review to ensure consistency, relevance, and quality in the data analyzed.

### 2.2.2 Outcome Assessment

We focused on assessing the accuracy of LLMs, defined as the percentage of completely correct answers (i.e., entirely accurate answers, fulfilling all components of the posed question without any factual or interpretational errors). The primary metric for accuracy (the most consistent outcome observed across the studies) was calculated as the percentage of completely correct answers relative to the total number of answers provided. In some studies, the same question was queried multiple times, and accuracy was calculated on the overall number of responses. Therefore, there may be some discrepancies across the overall number of questions in the study and the overall number of answers for accuracy calculation.

### **2.2.3 Data Extraction**

The data extraction was executed by two independent screeners (M.G., J.D.), and any discrepancies between the two were resolved by consulting a third screener (J.H.). We extracted the primary objectives from each study to understand the study's intent and the context in which LLMs were evaluated. We also gathered details about the specific LLMs used in the study, including their characteristics and any pertinent features. Further, we considered how the LLM prompts were engineered, refined, or manipulated for the study's purposes. Attention was given to the number of questions posed to the LLM and the methodologies employed to design or acquire these questions. We also noted whether these questions were reported directly within the main text of the manuscript or if they were relegated to supplementary materials. The number of evaluators participating in each study was recorded, as well as information concerning their background and expertise in the gastro/hep domain or in LLM assessment. An essential part of our extraction process was to understand the grading or assessment system employed by the evaluators. This includes any scales, rubrics, or qualitative methods that were used to evaluate the LLM's outputs.

### **2.2.4 Bias Evaluation**

The quality of the included studies was assessed using the Joanna Briggs Institute (JBI) critical appraisal tools.<sup>60</sup> These tools evaluate the level of methodological rigor in a study and the degree to which the study has considered and accounted for potential biases in its design, implementation, and analysis. According to the JBI guidelines, the evaluation was carried out by two authors (M.G. and J.H.), and disagreement was resolved by a third evaluator (K.Y.).

### **2.2.5 Data Synthesis and Statistical Analysis**

Given the high variability of LLMs version, prompt engineering, and prompt type we decided to descriptively report accuracies if this was part of the study design. Due to the heterogeneity of the studies with inconsistent data availability and the absence of shared outcomes, no statistical tests were performed.

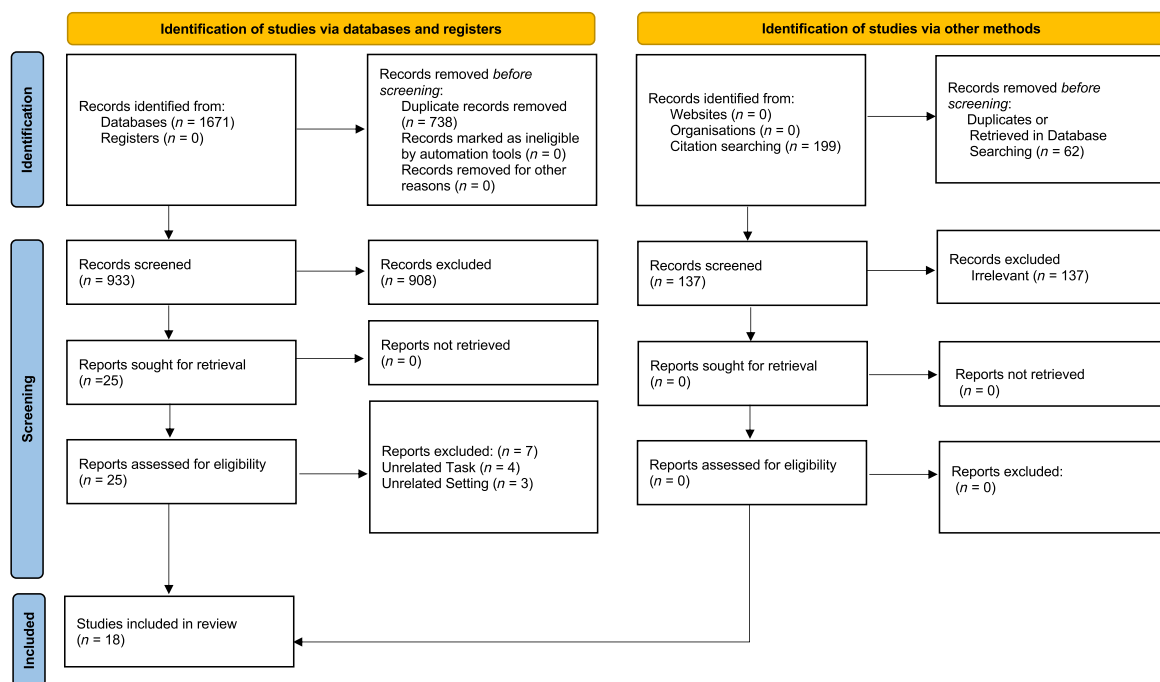
All the available questions from the included manuscript were collected and labeled by topic (gastrointestinal, pancreatic, and liver disease), question type (general question vs. clinical scenario), and perspective (patient vs. physician).

Once stored and labeled, each question undergoes a series of preprocessing steps to ensure uniformity and enhance the embeddings' quality, including lowercasing, removing punctuation and special characters, and word-based tokenization. Given all the studies investigated ChatGPT, we used the text-embedding model "text-embedding-ada-002" provided by OpenAI.<sup>61</sup> This model transforms each preprocessed sentence into a high-dimensional vector, capturing semantic and syntactic meaning. The embedding process is executed under standard parameters as defined by the model's documentation. To visualize and analyze high-dimensional embedding space, we employed t-distributed stochastic

neighbor embedding (t-SNE), a machine learning algorithm designed explicitly for dimensionality reduction of high-dimensional data.<sup>62,63</sup> The t-SNE algorithm is applied to the sentence embeddings to project them into a two-dimensional space, facilitating the exploration of their relative distances and clusters. The parameters of the t-SNE are chosen to represent better the distance between the distribution of the reduced vectors (i.e., perplexity of 50, learning rate of 1000, number of components of 2, and random state equivalent to 84). The two-dimensional t-SNE output is analyzed to identify clusters and patterns in the sentence embeddings. We assess similitude between questions to understand semantic similarities or differences by visual inspection only. The entire process was conducted using Python 3.10 (libraries: scikit-learn for t-SNE and matplotlib/seaborn for visualization).

## 2.3 Results

The literature search yielded 1,250 results, and after the removal of duplicates, 678 citations underwent title and abstract screening (Figure 2.1). The exclusion of 660 citations due to lack of relevance to the research question resulted in 18 citations for full-text review. The review of their full texts enabled the exclusion of 4 citations, and finally, 18 studies were included.<sup>9–26</sup> The excluded studies had an ineligible setting ( $n = 2$ ), an ineligible comparator ( $n = 1$ ), and an incorrect study design ( $n = 1$ ).



**Figure 2.1:** Screening and selection of studies. From the 1671 initially identified studies, only 18 met the inclusion/exclusion criteria and were selected for the systematic review.

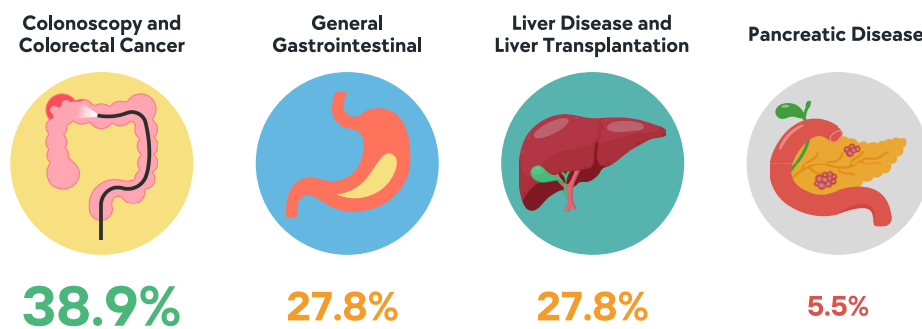
### 2.3.1 Study Characteristics

The main study characteristics are summarized in Table 2.1. As reported in Figure 2.2, out of 18 studies, 5 (27.8%) were performed for Hepatology-related topics<sup>18–22</sup>, 12 (66.7%) were performed for Gastroenterology-related topics<sup>10–17,23–26</sup>, and 1 (5.5%) were performed for topics across

Gastroenterology and Hepatology.<sup>9</sup> All studies used ChatGPT from OpenAI, though some studies (n = 3, 13.7%) investigated multiple LLMs, including Bard from Google<sup>10,18</sup> and YouChat or BingChat.<sup>26</sup> In particular, 12 (66.7%) studies reported the version of the LLM.<sup>10,12–14,16–19,21–23,25</sup> None of the studies used fine-tuning strategies but queried LLM foundational models. Among the included studies, 2 (%) reported using prompt engineering strategies<sup>12,18</sup>, 1 (%) probably used them given the task of assessing disease severity but providing clinical cases without a straightforward task to follow<sup>17</sup>, and 1 study employed prompt-engineering with contextualized guidelines.<sup>23</sup>

### 2.3.2 Question Generation

For studies with available question lists, the average number of questions asked per study was 40.5 ( $\pm 38.4$ ), with a maximum of 164<sup>22</sup> and a minimum of 8<sup>11</sup> questions, with 9 (50%) studies failing to provide clear information on how the questions were generated.<sup>9,10,13,15,17,18,20,21,24</sup> In terms of question availability, 2 (14.3%) studies reported only partial lists<sup>12,13</sup>, and 3 (21.4%) studies reported only case summaries or keywords without a clear prompt to address the aims of the study.<sup>16–18</sup>



**Figure 2.2:** Main topic distribution. Out of 18 studies, 7 (38.9%) involved queries related to colonoscopy and colorectal cancer, 5 (27.8%) involved general gastrointestinal topics (including gastroesophageal reflux disease, inflammatory bowel disease and miscellaneous queries), 5 (27.8%) were performed for hepatology- related topics (including liver radiology, liver transplantation, steatotic liver disease, liver cirrhosis and hepatocellular carcinoma) and 1 (5.5%) focused on pancreatic disease.

Study	Primary Objectives	LLM Type and Version	Number of Questions and Acquisition Strategy	Are questions reported in the manuscript main text or supplementary materials?
Lahat et al., 2023	Evaluate capabilities and limitations of ChatGPT in answering patients' questions in various gastroenterology subjects.	ChatGPT (OpenAI). Version not reported in the manuscript.	NOQ = 110. These real-life questions were gathered from open internet websites providing medical information to diverse groups of patients. The questions were selected to cover a wide range of gastroenterology and hepatology topics in terms of symptoms (n=23), diagnostic tests (n=45), and treatment (n=42). Specific selection criteria or information regarding the websites where the questions were collected is not reported in the manuscript.	Yes, supplementary materials.
Tariq et al., 2023	Define correctness in answering common patient inquiries related to colonoscopy by comparing three different LLMs.	ChatGPT 3.5/4 (OpenAI) and Bard (Google).	NOQ = 47. These questions covered several topics related to colonoscopy preparation, procedure, complications, and expectations. More specific generation/selection criteria are not reported in the manuscript.	Yes, supplementary materials.
Lee et al., 2023	Examine the quality of ChatGPT generated answers to common question about colonoscopy.	ChatGPT (OpenAI). Version not reported in the manuscript.	NOQ = 8. Questions were derived from the top-20 list of the US News & World Report Best Hospitals for gastroenterology and gastrointestinal surgery. The authors randomly selected 3 university-affiliated hospitals (1 from East Coast, 1 Midwest, and 1 West Coast) and retrieved 8 common questions about colonoscopy from publicly available webpages of these three randomly selected hospitals.	Yes, supplementary materials.
Gorelik et al., 2023	Examine ChatGPT ability to process endoscopy and pathology results in various scenarios and to provide guideline-based recommendations and a simple summary for the patient.	ChatGPT 4 (OpenAI)	NOQ = 20. A single gastroenterologist prepared 20 clinical scenarios of colonoscopy and pathology results that covers several recommendations known from society guidelines. Each clinical scenario was articulated in two forms: structured endoscopy report or a free-text clinical note, written with varied structure, terminology, and level of detail.	Yes, supplementary materials. However, partial report (only 6 example questions).
Moazzam, et al., 2023	Evaluate ChatGPT quality of answers to questions pertaining to pancreatic cancer.	ChatGPT (OpenAI) Version not reported in the manuscript.	NOQ = 30. The questions covered general information regarding pancreatic cancer (pre-operative, intra-operative and post-operative). More specific generation/selection criteria are not reported in the manuscript.	Yes, manuscript main text.
Cankurtaran et al., 2023	Assess if ChatGPT can serve as a reliable and useful resources for patients and healthcare professional for questions related to inflammatory bowel disease (IBD).	ChatGPT 4 (OpenAI)	NOQ = 20, half related to Crohn's Disease (CD) and the other half related to Ulcerative Colitis (UC). Patients' questions were identified upon trends in search on Google based on keywords related to CD and UC. Healthcare professionals' questions were generated by a committee of four gastroenterologists, led by an expert gastroenterologist with > 5 years of experience. These questions were related to the classification, diagnosis, activity, poor prognostic indicators, and complications of the disease.	No, the manuscript provides only keywords for each question (e.g., 'diet', 'symptoms', 'treatment') and not complete sentences.
Levartovsky et al., 2023	Assess ChatGPT ability to determine UC severity.	ChatGPT 4 (OpenAI)	NOQ = 20. Distinct acute UC presentation recorded to the local Emergency Department accumulated over 2 years. Case summaries included symptoms, vital signs, and laboratory results. More specific selection criteria are not reported in the manuscript.	No, the manuscript provides only case summaries, however the real prompt input is not present in the manuscript.
Patil et al., 2023	Determine whether AI chatbots recommend imaging consistent with the Appropriateness Criteria (AC) from the American College of Radiology for clinical scenarios related to the liver.	ChatGPT 4 (OpenAI) and Bard (Google).	NOQ = 27. Clinical scenarios directly listed in the AC under the following conditions: "Abnormal Liver Function Tests," "Chronic Liver Disease," "Acute Nonlocalizing Abdominal Pain," "Liver Lesion-Initial Characterization," "Right Upper Quadrant Pain," and "Palpable Abdominal Mass-Suspected Neoplasm".	No, the manuscript provides only the topic for each clinical scenario (e.g., 'abnormal liver function tests', 'chronic liver disease', and 'right upper quadrant pain').
Pugliese et al., 2023	Evaluated the accuracy, completeness, and comprehensiveness of ChatGPT responses to NAFLD-related questions.	ChatGPT 4 (OpenAI)	NOQ = 15. The questions were selected by two expert physicians considering their clinical expertise and guidelines. The questions were grouped into three domains: (1) specialist referral and the ability to detect worsening/improvement of NAFLD; (2) questions focused on diet composition, weight loss, medication, coffee consumption, and alcohol intake; (3) type of physical activity.	Yes, manuscript main text.
Endo et al., 2023	Evaluate the accuracy and completeness of ChatGPT responses to questions related to liver transplantation.	ChatGPT (OpenAI). Version not reported in the manuscript.	NOQ = 29. The questions covered broad general information regarding liver transplantation. In particular, the questions included general questions (n=4), waitlists (n=7), peri-transplant period (n=13), and donor-related information (n=5).	Yes, manuscript main text.
Cao et al., 2023	Evaluate the accuracy of information provided by ChatGPT relating to liver cancer screening, surveillance, and diagnosis.	ChatGPT 3.5 (OpenAI)	NOQ = 20. The questions were created to address fundamental concepts in liver cancer surveillance and diagnosis, with attention to American College of Radiology LI-RADS and American Association for the Study of Liver Disease guide-lines for hepatocellular carcinoma. More specific generation/selection criteria are not reported in the manuscript.	Yes, manuscript main text.
Yeo et al., 2023	Assess the accuracy, completeness, and reproducibility of ChatGPT responses to frequently asked questions about the management and care of patients with cirrhosis and hepatocellular carcinoma.	ChatGPT 3.5 (OpenAI)	NOQ = 164. Frequently asked questions with cirrhosis or hepatocellular carcinoma posted by well-regarded professional societies and institutions and social media such as Facebook. A total of 73 and 91 questions were selected for HCC and cirrhosis, respectively. Questions were divided into the following macro categories: "basic knowledge," "diagnosis," "treatment," "lifestyle," "preventive medicine", and "other".	Yes, supplementary materials.

**Table 2.1:** Summary of Studies Included in the Systematic Review in terms of primary objectives, LLM type and version, use of prompt engineering, number of questions, questions selection strategy, and question availability. LLM = large language model. NOQ = number of questions.

### 2.3.3 Question Clustering

The application of the t-SNE algorithm to sentence embeddings resulted in a two-dimensional representation that reveals distinct clusters and dispersion patterns corresponding to questions from various studies. As depicted in the visualisation, each cluster is colour-coded according to the study it originated from, with a legend provided to identify each study's respective dataset (Figure 2.3A). The t-SNE plot demonstrates that questions from the same study tend to aggregate, suggesting a degree of semantic and syntactic similarity within questions from a single study. This could indicate specific linguistic styles or terminologies employed consistently within individual studies. The distribution and density of the clusters could also reflect the diversity of the questions within each study. Further analysis identified clear clusters when questions were labelled according to the topic (Figure 2.3B) or question

type (Figure 2.3C). When questions were labelled according to their main topic, 57.4% (n = 374) involved gastrointestinal queries, 37.5% (n = 244) involved hepatology-related queries and 33 (5.1%) focused on pancreatic disease. In terms of question type, 83.4% (n = 543) were framed as general questions, whereas 16.6% (n = 108) were framed as clinical cases. On the contrary, clear clusters were not identified when questions were labelled according to patient (n = 121, 18.6%) versus physician perspective (n = 530, 81.4%) as shown in Figure 3D.



**Figure 2.3:** Two-dimensional representation of the high-dimensional embedding space distribution of questions using t-distributed stochastic neighbour embedding (t-SNE). The closer the two points are represented in the two-dimensional space, the more similar the semantic content of the questions. (A) Question distribution according to the study; (B) Question distribution according to the main topic (pancreatic disease, red; liver disease, blue; gastrointestinal disease, green). (C) Question distribution according to the question type (general questions, red; clinical cases, blue). (D) Question distribution according to the question perspective (patient perspective, red; general/physician perspective, blue).

### 2.3.4 Question Evaluators

The number and level of expertise of evaluators who graded LLM responses were reported across the studies<sup>9-26</sup> (Table 2.2). Most studies (n = 16, 88.9%) reported the exact number of evaluators.<sup>9-16,19-26</sup> In addition, only one study reported information related to years of experience and the average number of patients evaluated/procedures performed yearly by each evaluator.

### 2.3.5 Answers Grading Systems

The evaluation systems employed by the selected studies were heterogeneous, with findings summarized in Table 2. Numerical grading systems were employed by a significant portion of the studies (n = 7, 38.8%).<sup>9-11,16,19,21,22</sup> Within this subset, 3 studies reported using a Likert Scale<sup>11,16,19</sup>, with 2 studies employing a 7-point Likert Scale<sup>11,16</sup>, and 1 study utilizing a 6-point Likert Scale.<sup>19</sup> In

contrast, the remaining studies employed diverse evaluation methodologies. These alternative approaches included scoring responses on a scale of 0-2, categorizing them as entirely correct, correct but incomplete, or incorrect.<sup>10</sup> Certain studies focused on compliance with specific guidelines, categorizing responses as appropriate or inappropriate, while others evaluated specificity.<sup>12,14,18</sup> Other studies concentrated on agreement between evaluators and ChatGPT regarding disease severity and hospitalization recommendations.<sup>17</sup> In the context of liver disease, imaging choices were categorized using the Appropriateness Criteria, with points assigned based on recommendation appropriateness.<sup>18,26</sup> Finally, certain studies used qualitative grading categories like "Poor," "Fair," "Good," "Very Good," and "Excellent"<sup>20</sup>, while others categorized responses as accurate, inadequate, or inaccurate, and some employed grading systems based on comprehensive, correct but inadequate, mixed incorrect and incorrect/outdated data or completely incorrect criteria.<sup>22,25</sup>

Study	Evaluators Number	Evaluators Experience	Grading System	Results
Lalhat et al., 2023 <sup>9</sup>	3	The evaluators were gastroenterologists that worked in tertiary medical center and community clinics with more than 20 years of experience. In community clinics, and together cover all sub-specializations of gastroenterology: IBD experts, motility, hepatology, nutrition, and advanced endoscopy.	Each question was evaluated for accuracy, clarity, and efficacy using a scale ranging from 1 (lowest score) to 5 (highest score).	<b>Question Related to Treatments</b> Average accuracy, clarity, and efficacy scores were 3.9 (±0.80), 3.9 (±0.9), and 3.3 (±0.9), respectively. <b>Question Related to Symptoms</b> The average accuracy, clarity, and efficacy scores were 3.4 (±0.8), 3.7 (±0.7), and 3.2 (±0.7), respectively. <b>Question Related to Diagnostic Tests</b> The average accuracy, clarity, and efficacy scores were 3.7 (±1.7), 3.8 (±1.8), and 3.5 (±1.7), respectively.
Tariq et al., 2023 <sup>10</sup>	2	The evaluators were two gastroenterology fellows. No information is reported related to years of experience and/or number of performed procedures. Differences were resolved by third gastroenterologist. No further information provided.	Responses were scored on a scale of 0-2 (completely correct = 2, correct but incomplete = 1, incorrect = 0). Answers considered "unreliable if two responses for same query were considered inconsistent.	Among three models, ChatGPT-4 had the highest rate of completely correct answers (91.4%), with no answers graded as incorrect. 8.6% graded as correct but incomplete. For ChatGPT-3.5, only 3 (6.4%) were graded as completely correct, 40 (85.1%) as correct but incomplete, and 4 (8.5%) as incorrect. No responses for either Chat GPT 3.5 or 4 were considered unreliable. For Bard, only 7 (14.9%) were graded as completely correct, 30 (63.8%) were graded as correct but incomplete, and 10 (21.3%) as incorrect. 4.2% responses were considered unreliable.
Lee et al., 2023 <sup>11</sup>	4	The evaluators were two gastroenterology fellows and two senior gastroenterologists. No information is reported related to years of experience and/or number of performed procedures.	Responses were evaluated on a 7-point Likert Scale (7 = strongly agree, 4 = neutral, 1 = strongly disagree) for ease of understanding, scientific adequacy, and satisfaction with the answer.	The mean value for the task 'The answers are easy to understand' ranged from 5 to 6.4. The mean value for the task 'The answers are scientifically adequate' ranged from 5.4 to 6.5. The mean value for the task 'I am satisfied with the answers' ranged from 4.9 to 6.3. Chat GPT answers compared to non-AI answers were considered equivalent across all parameters.
Gorelik et al., 2023 <sup>12</sup>	2	The evaluators were two senior gastroenterologists, with more than a decade of consultant experience and who performs over 1000 endoscopies annually.	Responses were evaluated as compliant vs. non-compliant to guidelines (i.e., U.S. Multi-Society Task Force on Colorectal Cancer and the American Society for Gastrointestinal Endoscopy).	Of the 20 scenarios, 90% of the responses (18/20) complied with guidelines. In the 2 noncompliant scenarios, ChatGPT recommended a 10-year interval for a polyp with villous histology.
Henson et al., 2023 <sup>13</sup>	3	The evaluators were three board-certified gastroenterologists, including 2 esophagologist. No information is reported related to years of experience and/or number of performed procedures.	Responses were evaluated for appropriateness (completely appropriate, mostly appropriate, and mostly inappropriate) and specificity (only generic information and some specific guidance). Questions were asked 3 times during independent interactions without feedback to assess consistency.	<b>Appropriateness</b> Overall, ChatGPT provided appropriate responses 91.3% of queries, including 29.0% which were considered completely appropriate and 62.3% mostly appropriate. Prompts regarding treatment showed completely appropriate responses in 39.4% of cases, whereas questions involving diagnosis and diagnosis/management showed completely appropriate answers in only 26.7% and 14.3% of cases, respectively. <b>Specificity</b> Overall, ChatGPT provided some specific guidance in 78.3% of queries. Prompts regarding treatment were reported as providing some specific guidance in 75.8% of cases, whereas questions involving diagnosis and diagnosis/management in 93.3% and 71.4% of cases, respectively.
Emile et al., 2023 <sup>14</sup>	3	The evaluators were geographically diverse colorectal surgeons (Italy, Israel, and Brazil) selected based on "their clinical expertise and publications in colon cancer". However, no information is reported related to years of experience and/or number of performed procedures.	Responses were evaluated as follows: appropriate (consistent and factually accurate), inappropriate (consistent and factually inaccurate or contains incorrect information), or inconsistent. Two additional colorectal surgeons assessed ChatGPT answers for their consistency with ASCRS practice parameters.	The percentage of answers rated as appropriate varied among the experts (78.9%, 81.6%, and 122.4%). Overall, at least 2 of 3 experts rated the answers as appropriate for 86.8% of questions. 95% (19 of 20) of questions that applied to ASCRS practice parameters for colon cancer were concordant.
Moazzam, et al., 2023 <sup>15</sup>	20	The evaluators practiced in academic hospitals, with 80% being surgical oncologists, with a mean practice experience of 11.5 years and who performed on average 40 Whipple procedures annually.	Responses were graded as 'poor', 'fair', 'good', 'very good' and 'excellent'.	Across all Chat-GPT answers, 24.5% were evaluated as excellent, 35.2% very good, 21.3% good, 14.2% fair, 4.8% poor.
Cankurtaran et al., 2023 <sup>16</sup>	2	The evaluators were defined as gastroenterology experts. However, no information is reported related to years of experience and/or number of performed procedures.	Responses were graded according to a Likert scale of 1-7 (1 = lowest score, 7 = highest score) for two categories: reliability and usefulness. Inter-rater agreement was evaluated with Cronbach's $\alpha$ .	<b>Reliability</b> UC: Rater 1 reported a mean score of 4.2 (±1.2), whereas Rater 2 reported a mean score of 4.4 (±1.3) – with a Cronbach's $\alpha$ of 0.931. CD: Rater 1 reported a mean score of 4.6 (±1.3), whereas Rater 2 reported a mean score of 4.8 (±1.3) – with a Cronbach's $\alpha$ of 0.936. <b>Usefulness</b> UC: Rater 1 reported a mean score of 4.6 (±1.4), whereas Rater 2 reported a mean score of 4.5 (±1.3) – with a Cronbach's $\alpha$ of 0.986. CD: Rater 1 reported a mean score of 5.0 (±1.2), whereas Rater 2 reported a mean score of 4.9 (±0.9) – with a Cronbach's $\alpha$ of 0.925.
Levartovskiy et al., 2023 <sup>17</sup>	Not reported.	The evaluators were gastroenterologists working at the local hospital. However, no information is reported related to years of experience and/or number of performed procedures. Compared to actual decision to admit vs discharge which was made by the ED physician.	Agreement between evaluators and ChatGPT was calculated in terms of disease severity (mild, moderate, severe) based on the TrueLove and Wits Classification) and recommendation regarding hospitalization (discharge vs. admit).	<b>Disease Severity</b> ChatGPT graded 16/20 (80%) of the patients with the same severity if compared to the evaluators. In the remaining four cases, two severe cases were graded as moderate, and two moderate cases were graded as severe. <b>Hospitalization Recommendation</b> ChatGPT leaned toward hospitalization for 16 of 18 (88.9%) patients. Comparatively, only 12 of the 20 patients were hospitalized in actual clinical practice. In 1 case with moderate UC that was discharged, ChatGPT was in favor of hospitalization based on the patient's age and additional comorbidities. Inconsistencies in 4 cases stemmed

				primarily from inaccurate cut-off values for systemic variables (such as hemoglobin and tachycardia).
Patil et al., 2023 <sup>18</sup>	Not reported.	Not reported.	The Appropriateness Criteria (AC) include specific clinical variants and categorize imaging choices into one of three categories: (1) usually appropriate, (2) may be appropriate, and (3) not appropriate. Chatbot answers were give 1 point if imaging that was usually appropriate was recommended, 0.5 if imaging that may be appropriate was recommended, and 0 points if imaging that was not appropriate was recommended.	ChatGPT-4 was able to correctly identify usually appropriate imaging per the AC more often than Bard (89% versus 63%). ChatGPT-4 identified imaging that may be appropriate per the AC 11% of the time and never (0%) selected imaging that was not appropriate per the AC. Bard identified imaging that may be appropriate per the AC 7% of the time and imaging that was not appropriate per the AC 30% of the time.
Pugliese et al., 2023 <sup>19</sup>	11	Ten were key opinion leaders in NAFLD and one was a non-physician with expertise in patient advocacy in liver disease. However, no information is reported related to years of experience and/or number of patients evaluated each year.	Likert scale for accuracy (scale 1-6, 1 = lowest score, 6 = highest score), completeness (scale 1-3, 1 = lowest score, 3 = highest score), and comprehensiveness (scale 1-3, 1 = lowest score, 3 = highest score).	<b>Accuracy</b> The mean accuracy score was 4.84 ( $\pm 0.74$ ), with the physical activity domain having the highest mean score of 5.56 ( $\pm 0.56$ ), while the specialist referral domain had the lowest mean score of 3.9 ( $\pm 1.44$ ). <b>Completeness</b> The average completeness score was 2.08 ( $\pm 0.3$ ), with the physical activity domain having the highest mean score of 2.46 ( $\pm 0.5$ ), while the specialist referral domain had the lowest score of 1.73 ( $\pm 0.82$ ) <b>Comprehensiveness</b> The overall comprehensiveness rating was 2.87 ( $\pm 0.14$ ).
Endo et al., 2023 <sup>20</sup>	17	Experts in the fields of abdominal transplant surgery, practicing in academic centers, with an average career transplant volume of 396 cases, managing on average 33 cases of liver transplant each year.	Responses were graded as "Poor," "Fair," "Good," "Very Good" and "Excellent".	Most of the 493 total quality grades related to ChatGPT answers were graded as "Very Good" (n = 227, 46.0%) or "Excellent" (n = 149, 30.2%); only 7.5% (n = 37) were graded as "Poor" or "Fair".
Cao et al., 2023 <sup>21</sup>	6	Fellowship trained physicians from three academic liver transplant centers who actively diagnose and treat liver cancer. The evaluators included: 2 abdominal radiologists, 2 interventional radiologists, 1 medical oncologist, and 1 hepatologist.	Responses were graded as accurate (score = 1; all information is true and relevant), inadequate (score = 0; all information is true, but either the information does not fully answer the question or irrelevant information is provided that does not answer the question), or inaccurate (score = -1; any information is false). Each question asked 3 times for 3 answers with mean reviewer evaluations provided for each answer. Accurate response considered mean 0.5 or greater, inaccurate as $\leq 0.5$ . Reliable answers for questions as all 3 individual answers had mean $\geq 0.5$ .	A total of 25% of answers (15 of 60) were evaluated as inaccurate, whereas 48% (29 of 60) were found to be accurate. Five of 20 questions (25%) were considered reliable.
Yeo et al., 2023 <sup>22</sup>	2	Evaluators were board certified/eligible transplant hepatologists. However, no information is reported related to years of experience and/or number of patients evaluated each year. Discrepancies in grading were independently reviewed by 1 senior hepatologist with $\geq 20$ years of experience in transplant hepatology.	Grading system: comprehensive (score = 1), correct but inadequate (score = 2), mixed incorrect and incorrect/outdated data (score = 3), and completely incorrect (score = 4).	The proportion of responses graded as comprehensive or correct but inadequate was 75% or higher for "basic knowledge," "treatment," "lifestyle," and "others." However, this proportion was 66.7% in the "diagnosis" domain and 50% in the "preventive medicine" domain. The proportion of responses that were "mixed with correct and incorrect/outdated data" was 22.2%, 33.3%, 25.0%, 18.1%, and 50.0% in the "basic knowledge," "diagnosis," "treatment," "lifestyle," and "preventive medicine" domains, respectively. 79.1% responses on cirrhosis and 74% on HCC were deemed correct (1 or 2), while 47.3% in cirrhosis & 41.1% in HCC were labelled comprehensive (1 alone).

**Table 2.2:** Summary of Studies Included in the Systematic Review in terms of number of evaluators and their experience, used grading system and reported main results of each study.

### 2.3.6 Accuracy Evaluation

Specific studies reported the performance of LLMs in terms of factual accuracy measured as the ratio of completely correct answers to all answers. It was possible to measure accuracy in seven studies. Among these, only one study evaluated accuracy using ChatGPT-3, showing an accuracy of 28.9% (NOQ = 69).<sup>13</sup> Three studies evaluated accuracy using ChatGPT-3.5, showing accuracies of 6.4% (NOQ = 47)<sup>10</sup>, 25% (NOQ = 20)<sup>21</sup>, and 45.4% (NOQ = 164)<sup>22</sup> accuracy. Three studies evaluated ChatGPT-4, with *Kerbage et al.* showing accuracies for 48% of patient-oriented questions and 40% for physician-oriented questions (NOQ = 65)<sup>25</sup>, and the others showing 90% (NOQ = 20)<sup>12</sup> and 91.4% (NOQ = 47)<sup>10</sup> accuracy respectively. In addition, *Lim et al.* evaluated ChatGPT-4 with contextualized guidelines (NOQ = 62) and showed higher overall accuracy (79% vs. 50.5% for GPT-4 without contextualized guidelines).<sup>23</sup>

Other studies employed diverse metrics to define output accuracy, including numerical scales. For example, *Lahat et al.* reported average efficacy scores (5-point scale) of 3.3, 3.2, and 3.5 for questions related to treatments, symptoms, and diagnostic tests, respectively.<sup>9</sup> Regarding clarity and comprehensiveness, *Lahat et al.* reported average clarity scores (5-point scale) ranging from 3.7 to 3.9<sup>9</sup>, with *Pugliese et al.* showing an average comprehensiveness score of 2.87 (6-point scale).<sup>19</sup> Besides, when compared to non-AI responses, *Lee et al.* found that ChatGPT's answers held equivalent values

across ease of understanding, scientific adequacy, and user satisfaction, with mean values ranging from 4.9 to 6.5 (7-point scale).<sup>11</sup> Interestingly, when ChatGPT was used to make patient recommendations, it demonstrated a tendency to recommend hospitalization more frequently than actual clinical practice (88.9% vs. 60%)<sup>17</sup>, suggested the appropriate imaging in liver disease in 89% of cases<sup>18</sup>, while recommending longer follow-up colonoscopy for patients with high-risk polyps in 10% of cases.<sup>12</sup>

### 2.3.7 Study Bias Evaluation

The quality assessment performed using the JBI critical appraisal tool is reported elsewhere.<sup>56</sup> Four appraisal questions related to sample frame, sample size, and generalizability towards general or local populations do not apply to any studies. In addition, in any of the studies, it was impossible to determine if the findings were robust enough to conduct sensitivity analysis (low number of questions, low number of graders, and general unclear criteria for question generation). This collection of critical issues places all the studies at a high risk of bias. The lack of clarity in the criteria for question generation, combined with the insufficiency in the number of questions addressed and graders involved, significantly undermines the reliability of the findings.

## 2.4 Chapter's Deliverables

This paragraph enumerates the key findings from the systematic review and clarifies their implications for safe clinical use of LLMs in digestive diseases:

- **Corpus & scope:** 18 peer-reviewed studies included after screening. All queried ChatGPT; a minority also assessed Bard/YouChat/Bing. Most studies provided incomplete question lists; none used fine-tuning; few used explicit prompt engineering.
- **Accuracy (completely correct):**
  - ChatGPT-3.5: 6.4%-45.4% (NOQ = 231);
  - ChatGPT-4: 40%-91.4% (NOQ = 132);
  - GPT-4 with contextualized guidelines (RAG + prompts): 79% vs 50.5% without guideline context (NOQ = 62).
- **Safety-relevant behaviors noted in individual studies:** tendency to over-recommend hospitalization (88.9% vs 60% in practice); appropriate imaging suggestions in liver disease 89%; longer-than-guideline colonoscopy follow-up in 10% of high-risk polyp cases.
- **Evaluator & grading heterogeneity:** wide variation in grader numbers/experience and scales (binary, Likert, custom appropriateness), limiting cross-study comparability.
- **Bias assessment:** Overall high risk of bias due to unclear question generation, small samples, limited graders, and inconsistent outcomes; thus no meta-analysis was performed.
- **Clustering insight:** t-SNE shows aggregation by study, topic (GI 57.4%, liver 37.5%, pancreas 5.1%), and question type (general 83.4% vs clinical 16.6%); weaker separation by perspective (patient 18.6% vs physician 81.4%).



## Chapter 3

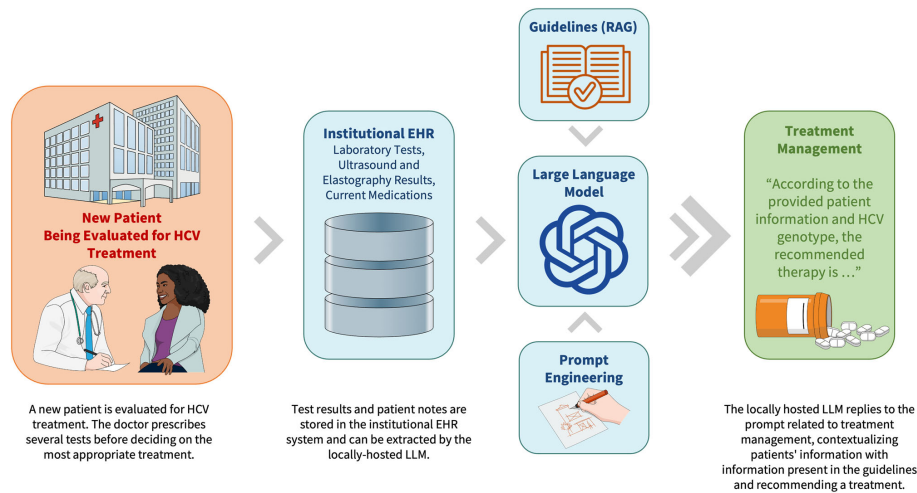
### **Build: Making Hepatitis C Virus Guidelines LLM-Friendly to Align Responses to Evidence-Based Medicine Using Retrieval Augmented Generation**

**“Everything should be made as simple as possible, but not simpler”**

*– Albert Einstein*

### 3.1 Chapter's Overview<sup>3</sup>

In the previous chapter, it emerged that baseline decoder-only LLMs are insufficient for safe use in digestive diseases—showing widely variable accuracy and safety-relevant errors—while performance improves when guideline context is supplied. Building on that evidence, this chapter investigates a guideline-grounded, retrieval-augmented framework—leveraging structured reformatting of recommendations and targeted prompt instructions—to generate clinical decision support systems (CDSS)-ready, guideline-faithful outputs as depicted in Figure 3.1.



**Figure 3.1: Example of a clinical decision support system integrated with large language models.** When a patient is being evaluated for HCV treatment, the doctor prescribes several tests (laboratory and imaging), whose results are stored in the institutional EHR system. The locally hosted LLM has a standardized clinical scenario prompt with laboratory and imaging values that are directly extracted from EHR. Afterward, the standardized prompt is queried to the LLM, which has access to the relevant guidelines to recommend the most appropriate treatment. HCV Hepatitis C virus, EHR electronic health record, RAG retrieval augmented generation, LLM large language model.

In this chapter, using the European Association for the Study of the Liver (EASL) HCV guideline as the reference corpus, the framework (i) transforms heterogeneous guideline content (including tables and flowcharts) into an LLM-friendly text structure; (ii) pairs it with retrieval-augmented generation (RAG) and principled prompting; and (iii) evaluates performance through an ablation program from baseline model behavior to progressively structured, guideline-aware configurations. Clinical decision-support implications are central: the framework is designed to turn patient-specific questions into accurate, transparent, and guideline-faithful recommendations suitable for point-of-care CDSS.

<sup>3</sup>This chapter (text and images) is adapted from the article: *Krešević S., Giuffrè, M., Ajčević, M. et al. Optimization of hepatological clinical guidelines interpretation by large language models: a retrieval augmented generation-based framework. npj Digit. Med. 7, 102 (2024). doi: 10.1038/s41746-024-01091-y*. The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit. Changes: Methods and Results are reproduced verbatim; Introduction/Background, Discussion, and narrative transitions are reformulated for thesis style; figure/table numbering and layout adapted to the thesis format.

The primary aim is to determine whether structured, guideline-integrated RAG can reliably increase the accuracy of GPT-4 Turbo on HCV screening, treatment, and safety questions compared with the baseline model.

Specific objectives include:

1. Build an LLM-friendly guideline corpus: convert the EASL HCV guideline from PDF to text; remove non-informative elements; and convert all non-text sources (notably tables) into text-based lists with consistent section headers (“Paragraph Title,” “Paragraph Text,” “Paragraph Recommendations”).
2. Standardize the question set: create 20 expert-authored prompts spanning text-based items, table-derived items, and realistic clinical scenarios to probe screening, management, drug–drug interactions, and adverse-event handling.
3. Implement an ablation of retrieval and instruction strategies: compare (a) baseline GPT-4 Turbo, (b) raw in-context guideline text, (c) cleaned guideline text with CSV tables, (d) fully reformatted guideline text with table-to-text conversion, and (e) added prompt-engineering; optionally test few-shot augmentation.
4. Quantify expert-graded accuracy and failure modes: use blinded hepatologist review to score complete correctness (primary endpoint) across question types; classify hallucinations using fact-conflicting and input-conflicting categories.
5. Contrast human accuracy with automatic similarity metrics: compute BLEU, ROUGE-L, METEOR, and cosine-similarity scores against expert reference answers to examine how lexical/semantic similarity relates to factual correctness.
6. Derive CDSS design guidance: translate experimental findings into practical recommendations on guideline formatting, retrieval strategy, and prompting needed to deliver reliable, guideline-faithful outputs for hepatology CDSS.

## **3.2 Materials and Methods**

### **3.2.1 Guidelines Selection**

We analyzed the current HCV guidelines from the prominent Northern American and European liver associations. Among these, we selected the European Association for the Study of the Liver (EASL) on the Hepatitis C Virus, entitled “EASL recommendations on treatment of hepatitis C: Final update of the series” published in 2020<sup>64</sup>, to explore our framework. The selected guideline comprised the most complex corpus of text containing broad recommendations on screening and management. In addition, the document contained in-depth information on drug–drug interactions, which was not reported in the Northern American guidelines.<sup>65</sup>

### 3.2.2 Standardized prompts creation

Two expert hepatologists (M.G. and L.S.C.) drafted 20 representative questions (Table 3.1). Fifteen questions addressed screening and management recommendations from each of the major sections, including the guideline main text (10 questions) and graphical tables (5 questions). Tables are a standard feature of clinical guidelines and summarize recommendations in specific ways that may not be reflected in the text. In addition, the two experts drafted five comprehensive clinical cases, each reflecting different HCV-related management strategies, including best treatment selection, drug–drug interaction, and management of treatment severe adverse reactions. All the questions are structured to test reasoning and comprehension from both the main text and tables.

<b>Text-based questions</b>	
1.	A screening blood test before a knee replacement surgery revealed a positive HCV antibody—what test should be performed to confirm HCV infection?
2.	When a patient with HCV can be considered cured after HCV therapy?
3.	Is there any major contraindication to HCV therapy?
4.	Is it possible to apply treatment without determining genotype using grazoprevir/elbasvir?
5.	What are the recommended treatment regimens and duration for a patient with HCV genotype 3 and no cirrhosis?
6.	A patient on the transplantation list for HCC and decompensated liver cirrhosis should be treated before or after transplantation.
7.	Should patients with HCV-positive patients be listed for kidney transplant treated? If yes, why?
8.	Patients with fibrosis F3, according to elastography, should be continuing HCC screening after successful HCV eradication?
9.	Is HCV treatment during pregnancy recommended?
10.	When should children born by an HCV-positive mother be tested for HCV infection?
<b>Table-based Questions</b>	
11.	What test can be used to assess the liver disease severity before treatment?
12.	Is there any interaction between cyclosporine and DAAs?
13.	Is there any interaction between apixaban and DAAs?
14.	Among anticoagulants and antiplatelets which is the one medication with the lowest risk of interactions with DAAs?
15.	What anticonvulsants are at higher risk of inducing drug interactions with DAAs?
<b>Clinical Scenarios</b>	
16.	A 45-year-old male with an unremarkable medical history was scheduled for a routine inguinal hernia surgery. As part of the preoperative evaluation, he was tested for hepatitis C virus (HCV) antibodies, which returned positive. Subsequent HCV RNA testing confirmed active infection, and genotyping identified the virus as HCV genotype 1a. The patient had no prior knowledge of his HCV status and had never been tested or treated for hepatitis C. Before initiating treatment, a liver elastography was performed to assess liver health, yielding a liver stiffness measurement of 5 kPa. What is the recommended treatment for this patient (drugs and duration)?
17.	A 55-year-old patient, previously lost to follow-up, returns to the liver clinic with a history of failed interferon-based therapy for HCV genotype 3. Recent laboratory tests confirm active HCV infection with genotype 3, accompanied by elevated liver enzymes (AST: 100 IU/L, ALT: 150 IU/L). Additional laboratory results include bilirubin at 1.2 mg/dL, creatinine at 0.87 mg/dL, albumin at 3.9 g/dL, and an INR of 1.10. Liver elastography shows a liver stiffness measurement of 15 kPa, without clinical signs of liver decompensation, as observed in the physical examination. What is the recommended therapy for this patient?
18.	A 60-year-old patient with advanced chronic kidney disease (CKD) at stage 4 is diagnosed with Hepatitis C virus (HCV) infection. The patient’s current renal function parameters include a creatinine clearance of 28 mL/min. Additionally, the patient presents with decompensated liver cirrhosis, classified as Child-Pugh Class B8, indicating significant liver dysfunction. What is the recommended therapy for this patient?
19.	A 60-year-old female patient diagnosed with Hepatitis C virus (HCV) genotype 1a, who does not have liver cirrhosis, was recently prescribed a 12-week course of Sofosbuvir (400 mg)/Velpatasvir (100 mg). The patient has a significant medical history of atrial fibrillation, for which she is being treated with amiodarone. During the initial assessment with the hepatologist, the patient inadvertently omitted mentioning their amiodarone treatment. As of now, the patient has not commenced the HCV treatment. Is it advisable for the patient to promptly inform her hepatologist about the amiodarone treatment before starting the HCV therapy?
20.	A 70-year-old female with a recent diagnosis of Hepatitis C Virus (HCV) genotype 1a, confirmed to have no evidence of liver cirrhosis, commenced a treatment regimen consisting of a 12-week course of Sofosbuvir (400 mg) combined with Velpatasvir (100 mg) daily. The patient’s baseline liver function tests were within normal limits, with an Alanine Aminotransferase (ALT) level of 45 IU/L (normal range: 30–45 IU/L). However, upon re-evaluation 4 weeks post-treatment initiation, her ALT levels had markedly elevated to 1123 IU/L. Should the prescribed HCV treatment be discontinued in light of this significant ALT elevation?

**Table 3.1:** Two expert hepatologists drafted 20 questions that specifically refer to information about management recommendations addressing information contained in the guideline main text (10 questions), graphical tables (5 questions), and clinical scenarios (5 questions).

### 3.2.3 Ablation study: customized LLM framework

We used a combination of RAG using EASL HCV guidelines, in different experimental settings with increasing degrees of complexity regarding guideline reformatting, prompt architecture, and few-shot learning to create a customized framework applied to the GPT-4 Turbo model (released by OpenAI, in November 2023 with knowledge updated until April 2023). Experiments with the OpenAI's Application Programming Interface (API) v. 1.17 cannot directly retrieve information from *.pdf* files. Therefore, the original pdf guidelines document was converted to a *.txt* file with UTF-8 encoding using the Python (v. 3.11) library PyPDF2 v3.0.

We carried out an ablation study from the baseline (Experiments 1 through 5) to investigate how different settings in guideline reformatting, prompt architecture, and few-shot learning impact the accuracy and robustness of LLM outputs (Figure 3.2). It is still unknown how non-text sources (e.g., graphical tables and flowcharts) are processed by LLMs and whether the information extracted is accurate. Therefore, we performed preliminary experiments to test the accuracy of the GPT image conversion process and found very low accuracy (16.0%) in extracting pertinent table information, with accuracy ranging from 0% (graphical tables) to 48.0% (only text tables). In light of these findings, we introduced text conversion of tables (non-text sources) into text-based lists and tested their impact on accuracy in Experiments 3, 4, and 5.

**Baseline.** Use of the foundational GPT-4 Turbo without any context. For this experiment, we only provided the questions without any further instruction.

**Experiment 1.** Use of the foundational GPT-4 Turbo with guidelines uploaded in context after pdf-to-text conversion in UTF-8 encoding without any additional text cleaning processes.

**Experiment 2.** Use of the foundational GPT-4 Turbo with guidelines uploaded in context after being manually cleaned with the removal of non-informative data (e.g. page header and bibliography). Tables presented as images in the original text were manually converted into *.csv* files and then provided as context.

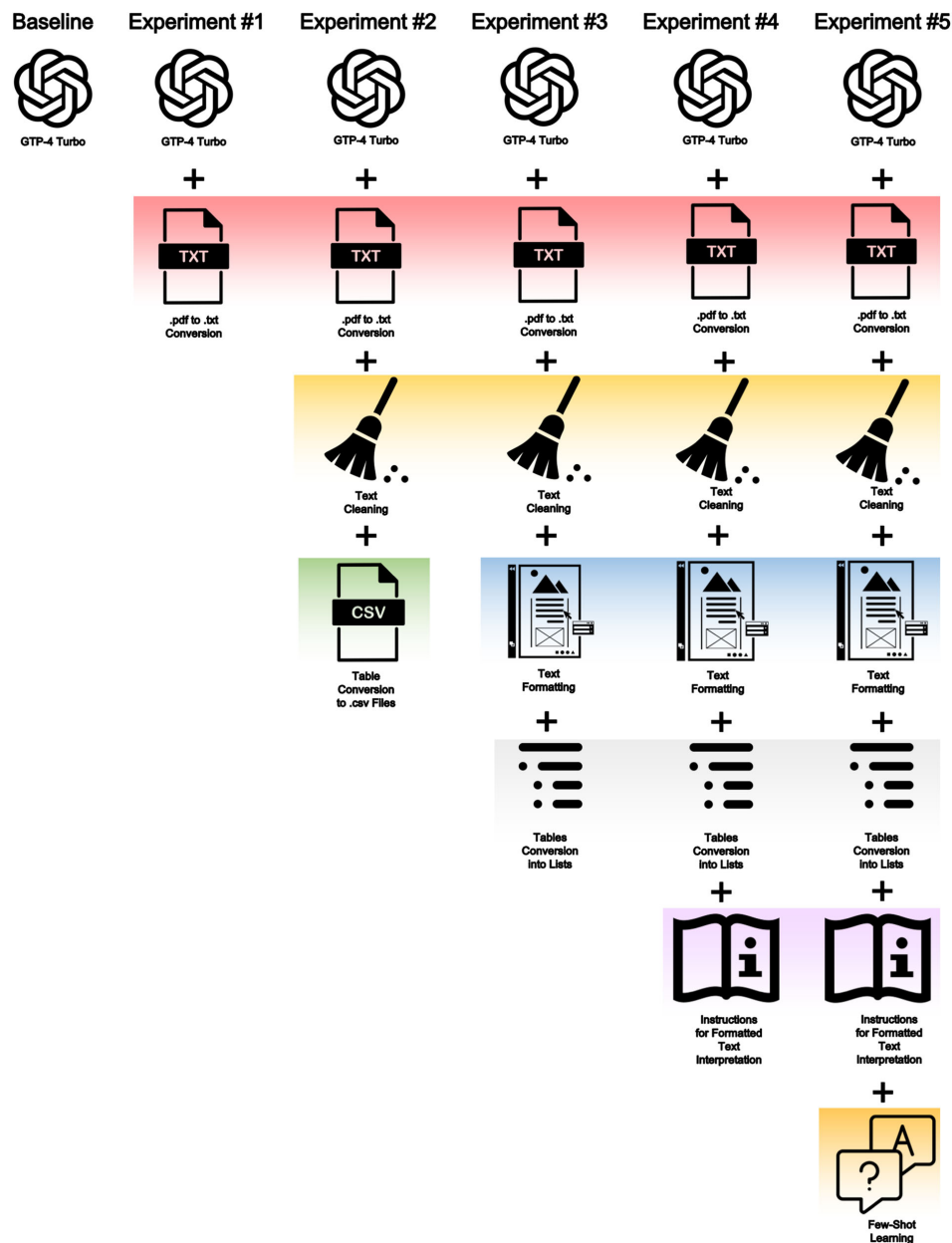
**Experiment 3.** Use of the foundational GPT-4 Turbo with guidelines uploaded as context that were cleaned and formatted to provide a consistent structure alongside the whole document. In addition, we converted all tables from *.csv* files into text-based lists and included them in the main text. Each paragraph title was preceded by "Paragraph Title". All the paragraph recommendations were collected and organized into a list preceded by "Paragraph Recommendations". Evidence reported in the main text was organized and preceded by "Paragraph Text".

**Experiment 4.** Use of the foundational GPT-4 Turbo with guidelines uploaded as context that were cleaned and formatted, with tables converted into text-based lists. We also provided a series of prompts

(i.e., prompt engineering) that instructed the model on how to interpret the structured guidelines elsewhere.<sup>56</sup>

**Experiment 5.** Use of the foundational GPT-4 Turbo with guidelines uploaded as context that were cleaned and formatted, with tables converted into text-based lists. We included the series of prompts (i.e., prompt engineering) and added a series of 54 question-answer pairs (i.e. few-shot learning).<sup>42</sup>

The experiments are summarized in Figure 3.2 and were conducted on a local Python environment with OpenAI API access. Instructions, when provided, are summarized elsewhere. We used foundational model default parameters, selecting a temperature of 0.9, and setting a maximum number of tokens in output equivalent to 800.



**Figure 3.2:** Depiction of Ablation Study experimental settings (Experiment 1 through Experiment 5) to investigate how guideline reformatting, prompt architecture, and few-shot learning impact the accuracy and robustness of LLM outputs.

### 3.2.4 Primary outcome

Our primary outcome was to evaluate qualitative rates of accuracy according to expert grading based on the information reported in EASL guidelines. We repeated the query 5 times each for the 20 questions for each experimental setting and reported the proportion of accurate responses. Each answer was graded with a score of 1 if the text contained completely accurate information or 0 otherwise. Two expert hepatologists (M.G., with four years of experience in treating HCV patients, and L.S.C., with thirty years of experience in treating HCV patients) manually graded each response. The two graders were blind to each other and towards the experimental setting when labeling answers. Disagreements in grading occurred for 5.0% of outputs and were solved by consensus between the two graders.

When outputs are considered inaccurate, the inaccuracy is caused by hallucinations (i.e., the production of plausible sounding but potentially unverified or incorrect information).<sup>32,33</sup> According to the recent definitions of Zhang et al., we defined three types of hallucinations: FCH, ICH, and CCH.<sup>66</sup>

### 3.2.5 Secondary outcome

Our secondary outcome was to evaluate the similarity of LLM-generated responses to the human expert-provided answers used as the gold standard. In particular, an expert hepatologist (M.G.) provided a single answer for each of the 20 questions, which was reviewed and approved by the second expert hepatologist (L.S.C.), and then used as the gold standard expert response to which LLM responses were compared in text-similarity using Recall-Oriented Understudy for Gisting Evaluation (ROUGE)<sup>67</sup>, Bilingual Evaluation Understudy (BLEU)<sup>68,69</sup>, Metric for Evaluation of Translation with Explicit Ordering (METEOR)<sup>69</sup>, and a Custom OpenAI score (based on cosine similarity). The Custom OpenAI score is based on cosine similarity, while the other scores are based on word overlap and semantic coherence between two text sources. We evaluated the similarity by comparing LLM-generated answers to the corresponding ones provided by experts. All these scores are expressed on a scale from 0 to 1, where a score of 1 denotes perfect alignment between two compared text sources. The mean and standard deviation of the similarities were estimated after repeating the query 5 times each for the 20 questions.

### 3.2.6 Statistical analysis

We employed the Chi-Square Test to compare accuracy among experiments qualitatively. We employed the Mann-Whitney U Test to compare differences among continuous scoring for automatic evaluation of answers. We considered statistically significant a two-tailed  $p$ -value  $< 0.05$ . To conduct the analysis, we used Python *v 3.11* and SciPy *v 1.11*.

### 3.3 Results

#### 3.3.1 Accuracy Analysis

The customized LLM framework achieved 99.0% overall accuracy, which was significantly better than the GPT-4 Turbo alone (99.0% vs. 43.0%;  $p < 0.001$ ). Incorporating in-context guidelines improved accuracy (67.0% vs. 43.0%;  $p = 0.001$ ). When the in-context guidelines were cleaned, and tables were converted from images to .csv files, accuracy improved to 78.0% (vs. 43.0%;  $p < 0.001$ ); after the guidelines were formatted with a consistent structure and tables were re-formatted to text-based lists, accuracy further improved to 90.0% (vs. 43.0%;  $p < 0.001$ ). Finally, the addition of custom prompt engineering led to an improvement in accuracy of 99.0% (vs. 43.0%;  $p < 0.001$ ), with no further improvement despite few-shot learning with 54 question-answer pairs (Table 3.2, Figure 3.3).

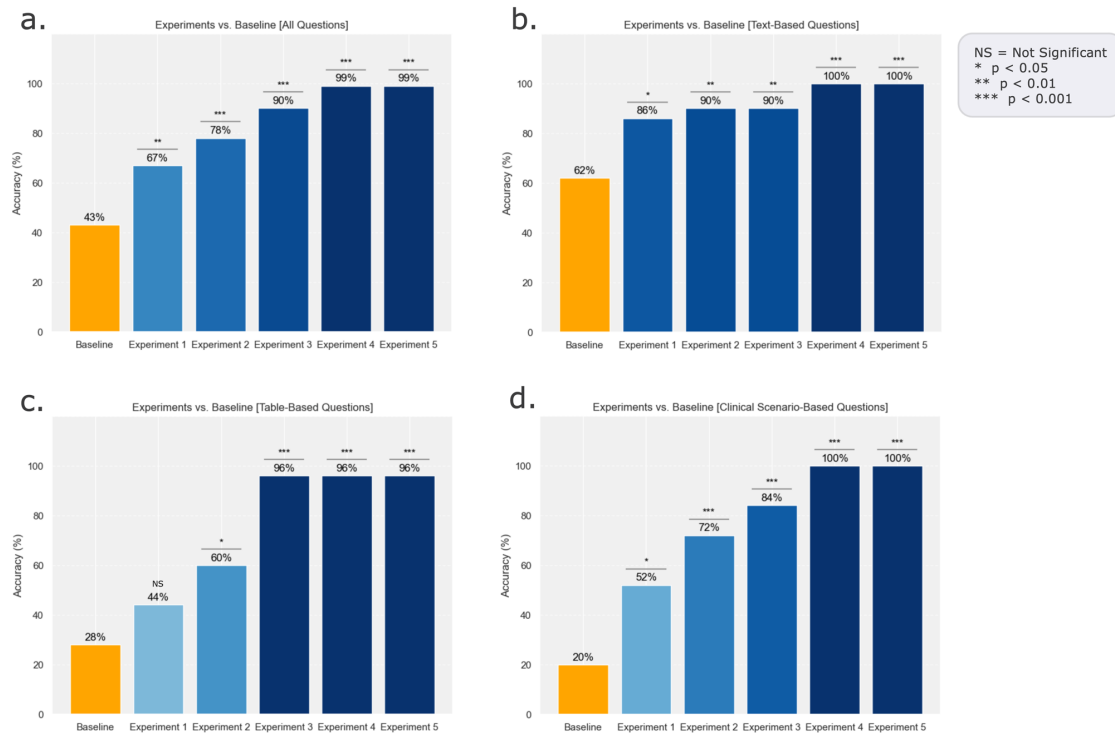
Metrics	Baseline	Experiment 1	Experiment 2	Experiment 3	Experiment 4
All questions:					
Accuracy	43.0%	67.0%	78.0%	90.0%	99.0%
Statistical Significance		$p = 0.001$	$p < 0.001$	$p < 0.001$	$p < 0.001$
Text-based questions:					
Accuracy	62.0%	86.0%	90.0%	90.0%	100.0%
Statistical significance		$p = 0.012$	$p = 0.002$	$p = 0.002$	$p < 0.001$
Table-based questions:					
Accuracy	28.0%	44.0%	60.0%	96.0%	96.0%
Statistical significance		$p = 0.377$	$p = 0.046$	$p < 0.001$	$p < 0.001$
Clinical scenarios questions:					
Accuracy	20.0%	52.0%	72.0%	84.0%	100.0%
Statistical significance		$p = 0.039$	$p < 0.001$	$p < 0.001$	$p < 0.001$

**Table 3.2:** Qualitative evaluation of accuracy based on human expert grading of each answer across all experimental settings. Statistical testing is based on pairwise comparison (Chi-Squared Test) between each experimental setting and the baseline.

For text-based questions, the customized framework achieved 100% overall accuracy, which was better than GPT-4 Turbo alone (100% vs. 62.0%;  $p < 0.001$ ). Incorporating in-context guidelines improved accuracy (86.0% vs. 62.0%;  $p = 0.01$ ); after cleaning the text and conversion of tables from images to .csv, further improvement in accuracy was achieved with no further improvement after formatting the text into a consistent structure and converting tables into text-based lists (90.0% vs. 62.0%;  $p = 0.002$ ). Adding custom prompt engineering resulted in 100% accuracy (100% vs. 62.0%;  $p < 0.001$ ) with equivalent performance after few-shot learning with 54 question-answer pairs (100% vs. 62.0%;  $p < 0.001$ ).

For table-based questions, the customized framework achieved 96.0% overall accuracy, which was better than GPT-4 Turbo alone (96.0% vs. 28.0%;  $p < 0.001$ ). Incorporating in-context guidelines improved accuracy (44.0% vs. 28.0%;  $p = 0.38$ ); after cleaning the text and conversion of tables from images to .csv, accuracy reached 60.0% (vs. 28.0%;  $p = 0.046$ ) with a substantial improvement after

converting tables into text-based lists and formatting the text into a consistent structure (96.0% vs. 28.0%;  $p < 0.001$ ) with similar performance in Experiments 4 and 5 as reported in Table 3.2.



**Figure 3.3.** Qualitative evaluation of accuracy among all experiments from baseline. A. Accuracy for all questions. B. Accuracy only for text-based questions. C. Accuracy for table-based questions. D. Accuracy for clinical scenario-based questions. Statistical testing is based on pairwise comparison (Chi-Squared Test) between each experimental setting and the baseline.

The customized framework achieved 100% overall accuracy for clinical scenarios, which was better than GPT-4 Turbo alone (100% vs. 20.0%;  $p < 0.001$ ). Incorporating in-context guidelines improved accuracy (52.0% vs. 20.0%;  $p = 0.039$ ); after cleaning the text and conversion of tables from images to .csv, accuracy reached 72.0% (vs. 20.0%;  $p < 0.001$ ) with a substantial improvement after converting tables into lists and formatting the text into a consistent structure (84.0% vs. 20.0%;  $p < 0.001$ ). Finally, the addition of custom prompt engineering achieved an accuracy of 100% (vs. 20.0%;  $p < 0.001$ ), with no further improvement despite few-shot learning with 54 question–answer pairs.

When inaccurate outputs were reviewed for hallucinations, we found 112 (90.3%) fact-conflicting hallucinations (FCH) and 12 (9.7%) input-conflicting hallucinations (ICH) across all experiments.

### 3.3.2 Text Similarity Analysis

For the secondary outcomes, we found differences in the customized LLM framework compared to the baseline across similarity scores (BLEU score, ROUGE-LCS F1, METEOR Score F1, and our Custom OpenAI Score) for all questions (Table 3.3).

Metrics	Baseline	Experiment 1	Experiment 2	Experiment 3	Experiment 4
BLEU Score					
Mean ( $\pm$ SD)	0.025 ( $\pm$ 0.023)	0.095 ( $\pm$ 0.088)	0.111 ( $\pm$ 0.143)	0.101 ( $\pm$ 0.094)	0.140 ( $\pm$ 0.119)
Significance		$p < 0.001$	$p < 0.001$	$p < 0.001$	$p < 0.001$
ROUGE-LCS F1					
Mean ( $\pm$ SD)	0.201 ( $\pm$ 0.053)	0.334 ( $\pm$ 0.120)	0.347 ( $\pm$ 0.138)	0.336 ( $\pm$ 0.114)	0.345 ( $\pm$ 0.119)
Significance		$p < 0.001$	$p < 0.001$	$p < 0.001$	$p < 0.001$
METEOR score F1					
Mean ( $\pm$ SD)	0.308 ( $\pm$ 0.059)	0.417 ( $\pm$ 0.104)	0.429 ( $\pm$ 0.126)	0.408 ( $\pm$ 0.101)	0.428 ( $\pm$ 0.115)
Significance		$p < 0.001$	$p < 0.001$	$p < 0.001$	$p < 0.001$
Custom OpenAI Score					
Mean ( $\pm$ SD)	0.939 ( $\pm$ 0.016)	0.954 ( $\pm$ 0.017)	0.956 ( $\pm$ 0.018)	0.956 ( $\pm$ 0.016)	0.957 ( $\pm$ 0.013)
Significance		$p < 0.001$	$p < 0.001$	$p < 0.001$	$p < 0.001$

**Table 3.3:** Evaluation of text-to-text-similarity between LLM-generated outputs and human expert-provided answers used as the gold standard across all questions. Statistical testing is based on pairwise comparison (Mann–Whitney  $U$  Test) between each experimental setting and the baseline.

### 3.4 Chapter’s Deliverables

This section distills the experimental evidence into concrete takeaways for constructing safe, guideline-aligned clinical decision support in hepatology. It translates the ablation results into practical design choices—what to include, how to structure it, and what not to rely on—to achieve reliable, guideline-faithful answers in HCV care. In particular:

- **Task & Dataset setup:**
  - External knowledge: 2020 EASL HCV guideline (selected for breadth and detailed drug–drug interactions).
  - Query set: 20 expert-crafted items spanning knowledge contained into text, graphical elements (e.g., tables), and realistic clinical scenarios.
  - Outcomes: primary—binary expert accuracy; secondary—BLEU, ROUGE-L, METEOR, cosine-based similarity.
- **Ablation finding #1 — RAG + formatting + prompts are decisive:**
  - Overall accuracy rose stepwise from 43% (baseline) → 67% (raw guideline context) → 78% (cleaned text; tables as .csv) → 90% (consistent structure; tables converted to text lists) → 99% (added prompt engineering).
  - Few-shot learning **did not** add further gains beyond RAG + formatting + prompts.
- **Ablation finding #2 — non-text sources are brittle unless converted:**
  - Direct image/table parsing was poor (16% overall; 0–48% depending on table type).
  - Converting tables to text lists and imposing a consistent document schema lifted table-question accuracy from 28% (baseline) to 96%.
- **Error profiling (safety signal):**
  - Among inaccurate outputs, **90.3%** were fact-conflicting hallucinations (FCH) and **9.7%** input-conflicting (ICH), underscoring the value of tight retrieval and structured prompts.
- **Similarity metrics vs. clinical correctness:**

- BLEU/ROUGE-L/METEOR and a cosine-based score **increased** across experiments, yet shifts in these metrics did **not** reliably mirror expert-graded factual accuracy—reinforcing the need for human clinical grading when safety matters.
- **Methodological takeaways for deployment:**
  - Use authoritative guidelines as retrieval corpora.
  - Reformat guidelines into LLM-friendly text (remove boilerplate, convert tables/flowcharts to structured lists, label sections consistently).
  - Add explicit prompt instructions on how to use that structure.
  - Do not expect few-shot examples to replace careful retrieval and formatting.
- **Scope & limitations (as reported)**
  - Focused on HCV (one disease area); limited runs and fixed temperature in the main ablation; did not benchmark non-OpenAI models in this chapter.

**Bottom line:** When guidelines are re-expressed as clean, consistently structured text and paired with targeted prompting, GPT-4 Turbo can deliver **near-perfect, guideline-faithful answers** to HCV questions—including table-derived items and clinical scenarios—whereas baseline, unguided use is insufficient.



## **Chapter 4**

**Validate: Leading Expert-Blinded Validation of Retrieval Augmented Generation and Supervised Fine-Tuning for HCV Management**

**“Our knowledge can only be finite; our ignorance is infinite”**

*– Karl Popper*

## 4.1 Chapter's Overview<sup>4</sup>

In the previous chapter, it emerged that baseline decoder-only LLMs are insufficient for safe use in digestive diseases—showing variable accuracy and safety-relevant errors—whereas performance improves markedly when guideline context is integrated through RAG and structured formatting. These findings demonstrated that LLMs can approximate expert reasoning when supplied with high-quality, guideline-grounded inputs. However, to determine whether such systems can achieve clinically trustworthy performance, a formal validation against domain experts is required.

This chapter therefore focuses on the next step: testing the reproducibility, robustness, and clinical validity of guideline-aligned LLM frameworks when evaluated by international hepatology specialists. The study leverages the participation of EASL HCV guideline authors and European hepatologists, who served as blinded graders of model-generated outputs. Their involvement provides an authoritative benchmark to assess whether retrieval optimization, fine-tuning, and decoding control can elevate performance to a level consistent with expert consensus.

The primary aim is to determine whether structured RAG and SFT can significantly improve the accuracy and clarity of GPT-4-based outputs for HCV management when validated by international hepatology experts.

Specific objectives include:

1. Optimize retrieval architecture: Test multiple text-chunking strategies (sentence-based, fixed-length, and paragraph-based) and retrieval depths (Top-1 vs Top-10) to identify the configuration that maximizes contextual relevance and efficiency.
2. Tune model hyperparameters: Systematically vary *temperature* and *top-p* to minimize hallucinations and stabilize factual accuracy across open-ended medical queries.
3. Evaluate fine-tuning approaches: Implement domain-specific supervised fine-tuning on guideline-derived question–answer pairs to assess whether internalized expert reasoning enhances accuracy beyond retrieval alone.
4. Conduct blinded expert validation: Compare baseline, RAG, and SFT configurations through dual-panel grading by hepatologists (guideline authors and independent clinicians), assessing accuracy and clarity on a standardized Likert scale and binary metrics.

---

<sup>4</sup>This chapter (text and images) is adapted from the article: *Giuffrè M., Pugliese N., Kresevic S., et al. From Guidelines to Real-Time Conversation: Expert-Validated Retrieval-Augmented and Fine-Tuned GPT-4 for Hepatitis C Management. Liver Int. 2025 Oct;45(10):e70349. doi: 10.1111/liv.70349.* The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit. Changes: Methods and Results are reproduced verbatim; Introduction/Background, Discussion, and narrative transitions are reformulated for thesis style; figure/table numbering and layout adapted to the thesis format.

5. Assess clinical applicability: Extend evaluation to simulated HCV clinical scenarios to determine whether optimized models can recommend appropriate antiviral regimens aligned with expert consensus.

## **4.2 Materials and Methods**

### **4.2.1 Question and Clinical Cases Generation**

Two junior hepatologists, M.G. (four years of experience, 100 patients treated/yearly) and N.P. (six years of experience, 100 HCV patients treated/yearly), selected four principal domains in the guideline (i.e., general indications to treatment, pre-therapeutic assessment, drug-drug interactions, and treatment of HCV patients with renal impairment) and five principal topics within the four domain and drafted representative questions to reflect general, patient-oriented, and physician-oriented perspectives for each of the topic (Table 4.1). Each question was generated to specifically target the retrieval of a set of information contained in single or multiple paragraphs of the guideline. Each question was reviewed and approved by an expert hepatologist who contributed as one of the main panel members of European Association for the Study of the Liver (EASL) HCV guidelines, A.A. (twenty years of experience, 100 HCV patients treated/yearly). In addition, the two junior hepatologists (M.G. and N.P.) who drafted the questions also provided free-text answers, which were reviewed and approved by the senior hepatologist (A.A.). These approved answers were then used as the gold standard text to which LLM responses were compared using various text-similarity metrics to assess the best chunking strategy and hyperparameters tuning.

In addition to the question set, M.G. and N.P. developed 25 simulated clinical cases (fully reported elsewhere)<sup>70</sup> to evaluate the models' performance in recommending appropriate HCV treatments. These cases reviewed and validated by the senior hepatologist (A.A.) to ensure clinical relevance and accuracy. The clinical cases were designed to represent the diverse spectrum of patients encountered in HCV clinical practice. Each case included specific details about: patient demographics (age, sex), HCV genotype and viral load, previous treatment history (treatment-naïve or treatment-experienced), presence and degree of liver fibrosis/cirrhosis, renal function status, relevant comorbidities, and concomitant medications with potential for drug-drug interactions.

<b>Domain #1: General indications to treatment: who should be treated?</b>	
General	Is therapy with direct-acting antivirals recommended for all patients diagnosed with chronic hepatitis C?
Patient	I have been diagnosed with hepatitis C. Do I need to take antiviral therapy even if blood tests and abdominal ultrasound show no signs of liver disease?
Physician	Should I treat a young patient recently diagnosed with chronic hepatitis C with normal transaminases and no fibrosis assessed by non-invasive methods with direct-acting antivirals?
<b>Domain #2: Pretherapeutic assessment: HCV genotype determination</b>	
General	Is it necessary to know the HCV genotype before starting therapy with direct-acting antivirals?
Patient	I have been diagnosed with hepatitis C and am about to start antiviral therapy. The doctor recommended viral genotyping, but I did not understand its usefulness. Do I need to have the result before starting antiviral therapy?
Physician	I have recently seen a newly diagnosed HCV patient from a geographical area where HCV subtypes that are intrinsically resistant to NS5A inhibitors are common. Do I need to test him for HCV genotype before starting first-line antiviral therapy?
<b>Domain #3: Drug-Drug Interaction Assessment</b>	
General	Is it possible to start therapy with direct-acting antivirals in a patient with chronic HCV infection and on anticoagulants therapy?
General	Is it possible to start therapy with direct-acting antivirals in a patient with chronic HCV infection and on statins?
Patient	This morning the hepatologist prescribed me Sofosbuvir/Velpatasvir for a chronic HCV infection. I forgot to tell him that I'm taking dabigatran for atrial fibrillation. Can I take the drugs together without any problems?
Patient	This morning the hepatologist prescribed me Glecaprevir/Pibrentasvir for a chronic HCV infection. I forgot to tell him that I'm taking atorvastatin for dyslipidemia. Can I take the drugs together without any problems?
Physician	I need to start therapy with direct-acting antivirals in a 75-year-old patient recently diagnosed with HCV. I do not know his genotype. He has no evidence of advanced chronic liver disease. He is taking dabigatran for a recent pulmonary thromboembolism. What is the best treatment option considering potential drug-drug interactions?
Physician	A 55-year-old patient has just been diagnosed with chronic HCV infection. She needs to start antiviral therapy with direct-acting antivirals as soon as possible. She has been taking atorvastatin for dyslipidemia which is currently under excellent control. What is the best treatment option?
<b>Domain #4: Treatment in patients with renal impairment</b>	
General	Can renal function influence the choice of treatment regimen to be used in a newly diagnosed patient with chronic HCV infection who is candidate for therapy with direct-acting antivirals?
Patient	I have been chronically infected with HCV for about 20 years but have been advised against starting antiviral therapy because I am on haemodialysis. Should I go to a liver disease referral centre?
Physician	I should start antiviral therapy for a chronic HCV infection, genotype 1 and naive to antiviral therapy, in a non-cirrhotic patient with severe renal impairment (eGFR < 30 ml/min/1.73m <sup>2</sup> ). Can I use Sofosbuvir/Velpatasvir or would it be better to use Glecaprevir/Pibrentasvir?

**Table 4.1:** List of questions across four domains and three perspectives (general, patient, and physician-oriented) generated by two hepatologists and reviewed by one of the panel members of Hepatitis C Virus Management Guidelines published by the European Association for the Study of the Liver.

#### 4.2.2 Qualitative Answer Evaluation of Open-Ended Questions

The primary aim of the study was to have HCV experts to evaluate qualitatively accuracy and clarity of LLM-generated responses. To ensure comprehensive assessment, we recruited two distinct groups of evaluators. The first group consisted of four expert hepatologists selected from the main authors and chairs of the EASL HCV guidelines<sup>64</sup> (F.N., M.P., J.M.P., and X.F.), which were used as the dataset for both the RAG external knowledge and for training in SFT. The second group comprised four expert hepatologists from a tertiary hepatology referral center (Humanitas University Hospital, Milan) who had no direct role in HCV guideline development, providing an independent perspective from clinical practitioners. None of the evaluators from either group were involved in any phase of question creation, answer generation, or dataset creation for fine-tuning purposes. Evaluators, who were blinded to model configuration, reviewed the same set of responses from all configurations for each question before moving on to the next group of responses. Each expert was free to review and score the response sets at their own pace, taking as much time as needed to ensure a thorough assessment. This self-paced approach helped mitigate potential fatigue or rush-related biases.

*Accuracy.* Given the lack of a widely accepted definition of accuracy and the absence of standardized quantification methods, as well as the use of multiple accuracy definitions in current literature, we opted to evaluate accuracy using both a binary grading system and a 10-point Likert scale. We chose a 10-

point scale, despite other studies using 5- or 7-point scales<sup>56</sup>, because scales with < 10 points have been proven to offer greater sensitivity and the ability to detect more nuanced differences in responses, while increasing reliability, validity, and data consistency.<sup>71,72</sup> This dual approach enables us to distinguish not only completely correct answers from incorrect ones but also to gauge the level of inaccuracy among the incorrect responses. While the binary classification helps identify totally accurate answers, the Likert scale provides insights into the relative accuracy of all answers, which is especially important for identifying and addressing partial inaccuracies. The complete classification for the 10-point Likert scale is reported in Table 4.2. For the binary scale, graders assigned a score of one for answers with entirely correct and presenting informative text without any hallucinations, and a score of zero for answers that did not meet these criteria.

*Clarity.* Clarity was defined as the presence of relevant information that could be easily understood, straight to the point, explicit, and free from ambiguity. In evaluating clarity, the graders provided a binary score of one for clear answers and a score of zero for unclear answers.

Grades	Definition
Grade: 0	Completely Inaccurate: The output is entirely irrelevant, incorrect, or nonsensical, showing no understanding of the query.
Grade: 1	Extremely Inaccurate: The response contains major errors and misunderstandings, with very little relevant or accurate information.
Grade: 2	Highly Inaccurate: Significant inaccuracies dominate the response, though there might be a minor element of relevance or accuracy.
Grade: 3	Very Inaccurate: While mostly inaccurate, the response shows some basic understanding of the topic, but with major errors.
Grade: 4	Inaccurate: The response has a mix of correct and incorrect information, but the inaccuracies are more prominent.
Grade: 5	Moderately Accurate: The response is a balance of accurate and inaccurate information, showing an equal mix of correct insights and errors.
Grade: 6	Somewhat Accurate: The response is more accurate than not, with some notable inaccuracies but a general understanding of the topic.
Grade: 7	Mostly Accurate: The response contains mostly correct information, with minor errors or inaccuracies.
Grade: 8	Very Accurate: The response is highly accurate, with only very slight inaccuracies or areas of uncertainty.
Grade: 9	Highly Accurate: The response is extremely accurate, showing a deep understanding of the topic with almost no inaccuracies.
Grade: 10	Completely Accurate: The response is entirely accurate, with no discernible inaccuracies or errors, perfectly addressing the query.

**Table 4.2:** Grades and definitions of the 10-point scale to evaluate the accuracy of answers generated by each model.

#### 4.2.3 Quantitative Answer Evaluation of Open-Ended Questions

Two junior hepatologists (M.G. and N.P.) who drafted the questions also provided free-text answers, which were reviewed and approved by the senior hepatologist (A.A.). The expert responses were used as the reference text to develop and evaluate embedding-based text similarity metrics (i.e., cosine similarity) to determine the best chunking strategy in RAG framework, and in the process of hyperparameters tuning. For the embedding-based metrics, text was first tokenized internally by OpenAI's cl100k\_base tokenizer, which splits text into subword tokens using byte-pair encoding (BPE), and then converted into its embedding representation. An embedding is a high-dimensional vector representation of data, typically text, which captures the semantic and syntactic nuances of the input.<sup>73</sup> We used one of the currently available text-embedding models provided by OpenAI (i.e., text-embedding-3-large)<sup>61</sup>, which generates embeddings with 3072 dimensions for each text input. The embedding vectors were then used to compute similarity scores between expert responses and model-generated answers. For better visualization of the relative gap between the cosine similarity score from

different models, we provide the transformation of first normalizing the similarity raw score with its maximum attainable score and then applying the logit function.

#### **4.2.4 Quantitative Evaluation of Clinical Scenario Treatment Recommendation**

For each case, a gold standard for appropriate DAA regimen selection was established through a consensus process involving four independent hepatologists from tertiary referral centers, who were not involved in case creation. Their recommendations were determined by majority vote, with subsequent expert review for cases lacking initial consensus. We used these consensus-based gold standard recommendations to evaluate the performance of each model configuration across all 25 clinical scenarios reported in details elsewhere.<sup>70</sup> Model performance was assessed by comparing model recommendations to the expert majority consensus, using complete accuracy as the metric: ‘all recommended treatment regimens without any incorrect recommendations’. This metric assessed whether the model could identify all appropriate treatment options according to the expert consensus, without suggesting any treatments that diverged from the consensus. This represents an optimal standard for clinical decision support, as it would provide clinicians with the full range of appropriate treatment options.

#### **4.2.5 Model Configurations and Answer Generation**

Experiments were conducted on a local Python (*version* 3.11) environment, using the OpenAI Action Programming Interface (API) (*version* 1.17) to interact with GPT models using the model “gpt-4-turbo-2024-04-09-preview” according to OpenAI’s nomenclature. We deliberately chose not to evaluate advanced prompting strategies in this study, as our previous research has demonstrated that performance gains achieved through sophisticated prompting techniques are substantially less significant compared to those obtained through retrieval augmentation or fine-tuning approaches.<sup>23,42</sup> Additionally, limiting the number of experimental conditions allowed us to reduce the substantial labor burden associated with expert human grading while focusing on the most promising enhancement methods for clinical applications.

However, we developed and integrated tailored heuristic prompts to direct answer generation for each model configuration, with specific prompting strategies designed separately for the 15 open-ended questions and the 25 clinical scenarios. These prompts were crafted to elicit relevant medical information and treatment recommendations in a clinically appropriate format, as detailed in elsewhere.<sup>70</sup>

Prior to the main evaluation, we conducted a comprehensive hyperparameter optimization process, which is essential to determine the correct hyperparameter combination that reduces the model’s tendency to explore less likely tokens, thus reducing the risk of hallucinations.<sup>74-76</sup> Using the open-ended questions as a test bed, we systematically explored the impact of Temperature and Top\_p settings

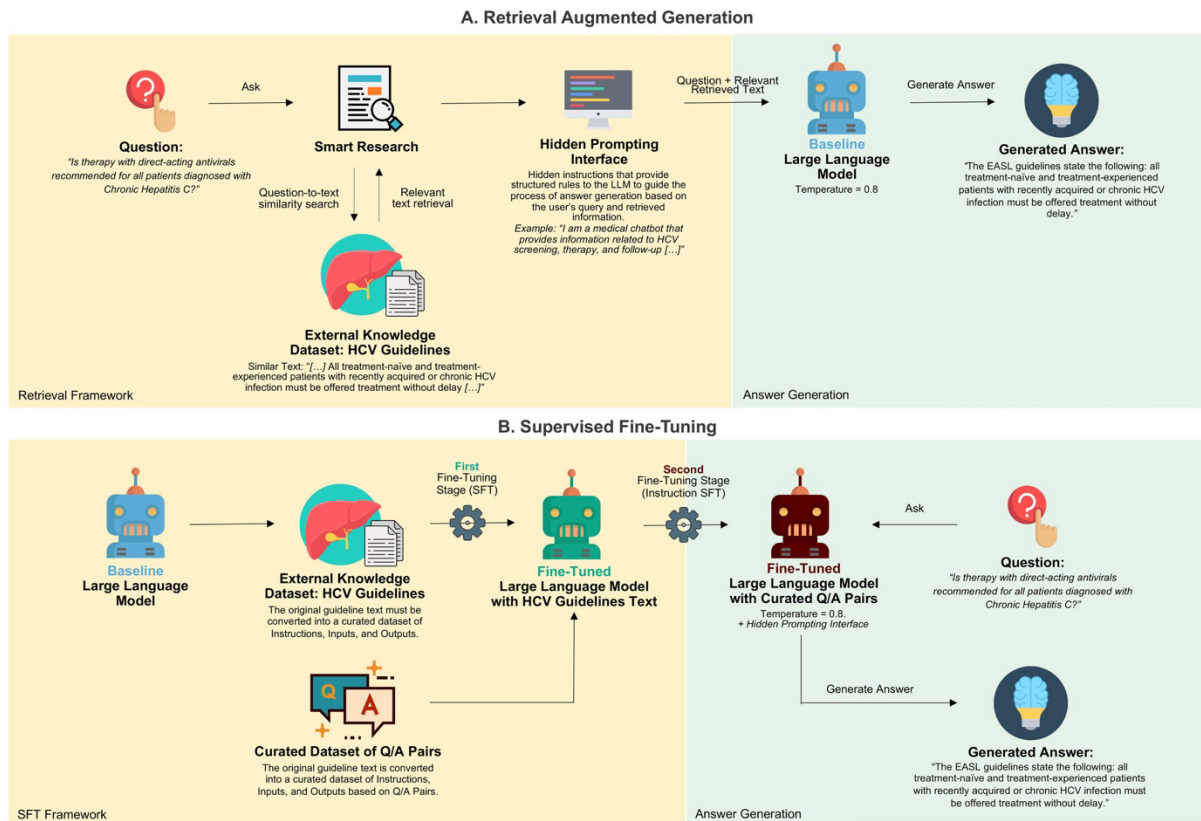
on response quality. For Temperature, we generated 10 responses at each threshold from 0 to 2.0 in increments of 0.2, while holding Top\_p constant at its default value. Similarly, for Top\_p, we generated 10 responses at thresholds from 0 to 1.0 in increments of 0.1, while maintaining temperature at its default setting. Response quality was evaluated by measuring similarity to expert-provided answers using embedding-based metrics. We simultaneously optimized the RAG framework by testing three different chunking strategies across varying numbers of retrieved chunks, generating 10 responses for chunk number. This systematic evaluation allowed us to identify the optimal chunking strategy and retrieval depth. Based on these optimization experiments, we selected the highest-performing hyperparameter values and RAG configuration to generate the responses for expert evaluation. For qualitative grading, only the first generated response was considered, simulating real-world application conditions. The optimized parameters were subsequently applied to generate responses for the 25 clinical scenarios, where we again conducted 10 interaction rounds per case to assess recommendation accuracy and consistency.

*Baseline Model.* As the baseline, we used the model employing only the heuristic prompts detailed elsewhere<sup>70</sup> without any enhancement techniques.

*Retrieval Augmented Generation.* For the RAG<sup>77</sup> framework, we chose the EASL guidelines on HCV<sup>64</sup> as the reference text based on our group of expert regions of origin and medical practice. The reference text was provided in an LLM-friendly version as previously validated by manually removing non-informative data (e.g., headers or reference numbers), and converting all non-textual sources (e.g., graphical tables containing drug-drug interactions) to text-based lists.<sup>42</sup> We explored the impact of the number of retrieved chunks across three chunking strategies<sup>78</sup> to optimize retrieval within the RAG framework: sentence-based chunking, in which a special delimiter was inserted at each period to isolate individual sentences; fixed-length chunking, which segmented the text into 512-token units with a 100-token overlap between adjacent chunks; and paragraph-based chunking, in which paragraphs were manually delineated to preserve semantic coherence. Each chunk was then encoded into a 3072-dimensional embedding vector by the embedding model using OpenAI's text-embedding-3-large.<sup>61</sup> To optimize retrieval efficiency, we utilized the Facebook AI Similarity Search (FAISS) library for vector indexing and similarity search. The embeddings were normalized and stored as float32 numpy arrays to ensure computational efficiency. For retrieval, user queries underwent the same tokenization and embedding process; cosine similarity scores were computed between the query embedding and each chunk embedding using FAISS, and chunks with highest score were ranked accordingly. To generate answers, we employed the paragraph-based chunking strategy (which performed best across all tests) with two distinct retrieval configurations: a 'RAG-Top1' configuration that retrieved only the single most relevant chunk, and a 'RAG-Top10' configuration that retrieved the top 10 most relevant chunks based on cosine similarity scores. While the similarity improvement from 1 to 10 chunks was not

dramatically different (with both configurations providing good performance), our analysis (Figure 4.2) showed that performance gains reached a plateau after approximately 10 chunks, offering an optimal balance between answer quality and computational efficiency. Finally, the complete RAG-based system was integrated into a publicly accessible chatbot interface (<https://github.com/liver-mainds/hcv-gpt>).

*Supervised Fine-Tuning.* The fine-tuning framework was conducted on the OpenAI's platform through a structured process involving two stages and four curated datasets, all available in the GitHub repository (<https://github.com/liveraidlab/hcv-gpt>), and developed according to OpenAI's instructions<sup>79</sup> using the *jsonl* file format. Each entry in the dataset contained a json object with fields for role, content, and metadata. The training hyperparameters of batch size, learning rate multiplier, and the number of epochs were set to "auto". OpenAI internally ran extensive experiments on billions of tokens and hundreds of tasks during the development of InstructGPT. The default hyperparameter with the "auto" set-up are those that, on average, offer the best trade-off between convergence speed, robustness, and the risk of overfitting on heterogeneous datasets.<sup>80</sup> To avoid the loss of relevant information by splitting the EASL guideline dataset, we used the EASL guideline to build the training dataset and the HCV guidelines developed by the American Association of the Study of The Liver to construct the validation dataset.<sup>81</sup> As reported on the fine-tuning command panel, this initial fine-tuning phase used 723837 tokens and showed a 0.45 training loss and 0.97 validation loss (these parameters are determined automatically after successful training and validation and reported on OpenAI's fine-tuning platform). We further enhanced the fine-tuned model using a curated question-and-answer dataset based on expert knowledge to align the model's outputs with the expertise. For this stage of Instruction Fine-Tuning a curated dataset was created combining the original training data with question-answer (Q/A) pairs derived from the EASL guidelines to capture all the content contained in the guideline text in a Q/A format. A separate validation set from the American guidelines was used to evaluate model performance. As reported on the fine-tuning command panel, the second stage used 172887 tokens and showed a 0.01 training loss and a 1.34 validation loss (these parameters are determined automatically after successful training and validation and reported on OpenAI's fine-tuning platform). The interpretation of these metrics remains proprietary of OpenAI and, therefore, is reported only for reproduction purposes. The Q/A pairs used for each round of instruction fine-tuning differ from the questions used for performance evaluation, and they are available for download on our GitHub repository (<https://github.com/liver-mainds/hcv-gpt>). A summary of the SFT framework is depicted graphically in Figure 4.1.



**Figure 4.1:** Overview of retrieval-augmented generation and supervised fine-tuning. In the retrieval-augmented generation approach (Section A), a user question initiates a question-to-text similarity search that retrieves relevant information from external knowledge datasets, such as HCV guidelines. A hidden prompting interface then provides structured instructions to guide the large language model in generating an accurate answer based on the query and retrieved text. The baseline large language model, such as GPT-4 Turbo, processes the question and relevant text to generate a comprehensive answer. In the supervised fine-tuning (Section B), the process begins with a baseline large language model, which undergoes a first fine-tuning stage using a dataset that includes HCV guidelines text. This is followed by a second fine-tuning stage (instruction fine-tuning) incorporating a curated dataset of Q/A pairs. The final fine-tuned model responses are guided by hidden prompting instructions and can generate responses without searching external datasets. HCV, hepatitis C virus; SFT, supervised fine-tuning; Q/A, question/answer.

#### 4.2.6 Statistical Analysis

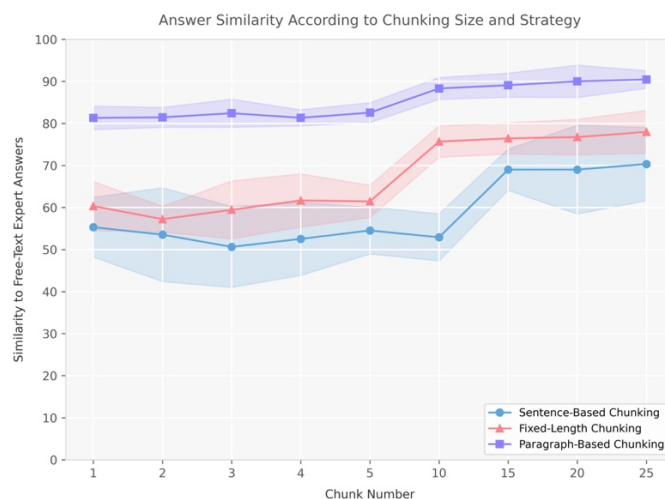
For the continuous 10-point Likert Scale, agreement among graders was evaluated using the One-Way Random Intraclass Correlation Coefficient<sup>82</sup>, with values reported in an interval between 0 and 1. According to Cicchetti et al.<sup>83</sup>, we considered the following interpretations: ICC < 0.40: poor agreement; ICC ≥ 0.40 and < 0.60: fair agreement; ICC ≥ 0.60 and < 0.75: good agreement; ICC ≥ 0.75: excellent agreement. For the categorical evaluations, agreement among graders was evaluated using Fleiss' Kappa, whose values were interpreted as follows: Kappa < 0.20: slight agreement; Kappa ≥ 0.20 and < 0.40: fair agreement; Kappa ≥ 0.40 and < 0.60: moderate agreement; Kappa ≥ 0.60 and < 0.80: substantial agreement; Kappa ≥ 0.80: almost perfect or excellent agreement<sup>84</sup>. Data were reported as the mean and standard deviation for the 10-point accuracy scale. We employed the Chi-Square Test to compare (binary) accuracy and clarity between all configurations of experimental settings (baseline, RAG, and SFT) and the paired sample Student's t-Test to compare differences among the 10-point accuracy scale grading between all configurations (baseline, RAG, and SFT). Due to the number of

multiple comparisons ( $n = 4$ ), we applied the Bonferroni correction and considered significant a two-tailed p-value  $< 0.012$ . To conduct the analysis, we used Python v3.11 and SciPy v1.11.

## 4.3 Results

### 4.3.1 Optimal Number and Strategy of Text Chunking

To build the optimal RAG framework, we evaluated three chunking strategies—sentence-based, fixed-length, and paragraph-based—across nine different chunk counts (1, 2, 3, 4, 5, 10, 15, 20, 25) and measured similarity to free-text expert answers (Figure 4.2). Paragraph-based chunking consistently outperformed both the sentence-based baseline (55–70% similarity) and fixed-length chunking (60–78%), achieving 81–90% similarity as the number of chunks increased. Notably, most of the gain in paragraph-based performance occurs by retrieving the top ten chunks, reaching a plateau thereafter. To strike the best balance between answer quality and computational cost, we therefore focused subsequent comparisons on two retrieval scenarios: Top-1 versus Top-10 chunk numbers using the paragraph-based strategy. By doing so, we capture the steep initial improvement afforded by paragraph chunks while containing the overhead of larger retrieval sets.

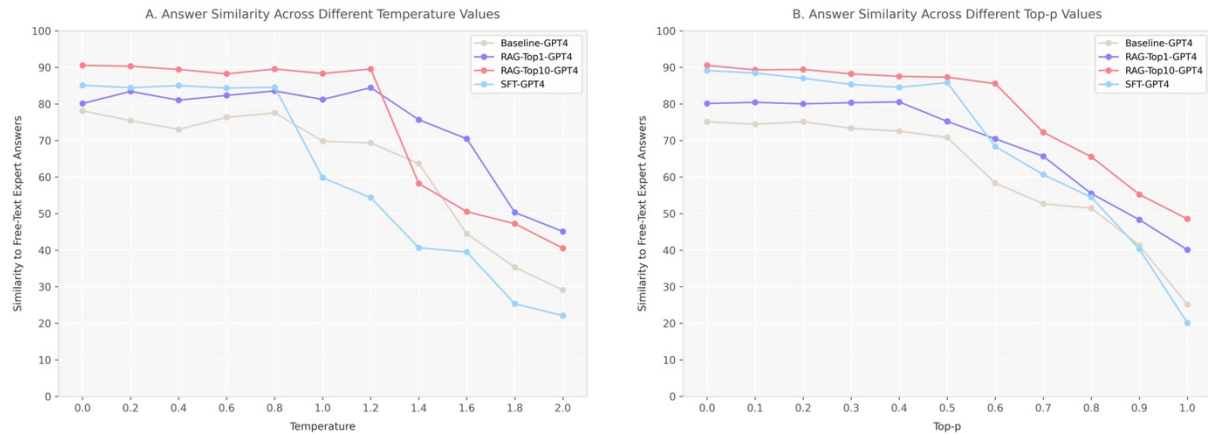


**Figure 4.2:** Comparison of answer similarity to expert responses across different chunking strategies and chunk numbers in the retrieval framework. Three chunking approaches were evaluated: Sentence-based (blue), fixed-length (red) and paragraph-based (purple). The shaded areas around each line represent the distribution range (minimum to maximum values) of similarity scores across multiple test questions for each chunking strategy, illustrating the consistency and variability of performance at different chunk numbers.

### 4.3.2 Optimal Hyperparameters Tuning

Prior to generating responses for expert evaluation, we conducted a comprehensive hyperparameter analysis to determine the optimal settings for answer generation using cosine similarity with the free-text answers provided by one of the guideline authors using the question sample reported in Table 4.1. Figure 4.3 illustrate the impact of temperature and Top-p hyperparameters on model response similarity to expert answers. For all models (Baseline-GPT4, RAG-Top1-GPT4, RAG-Top10-GPT4, and SFT-GPT4), optimal performance is achieved with low values of both temperature (0-0.8) and top-p (0-0.5). Notably, the RAG-Top10-GPT4 model maintains high accuracy ( $>85\%$ ) up to temperature 1.2 and top-p 0.6, while other models deteriorate more rapidly. All models exhibit a dramatic performance decline with temperature  $\geq 1.4$  and Top-p  $\geq 0.7$ , confirming that higher values introduce excessive variability in responses. These findings suggest that to optimize accuracy in specialized question-answering tasks,

it is preferable to maintain temperature between 0-0.8 and Top-p between 0-0.5, with the RAG-Top10-GPT4 model offering the greatest robustness to variations in these parameters. Based on these results, we selected the final values of temperature = 0.8 and Top-p = 0.5 for response generation, as these represent the highest parameter values before the rapid performance decline, balancing response quality with appropriate diversity.



**Figure 4.3:** Impact of hyperparameter settings on answer similarity across different model configurations. Panel A shows the effect of temperature values (0–2.0) on similarity to expert answers, while Panel B illustrates the impact of Top- $p$  values (0–1.0).

### 4.3.3 Inter-grader Agreement

The inter-grader agreement between experts was evaluated using the Intraclass Correlation Coefficient (ICC) for the 10-point scale and Fleiss' kappa for binary evaluations, as summarized in Table 4.3. Among guideline authors, the highest agreement was observed for RAG-Top1 with an ICC of 0.87 (95% C.I. 0.76-0.94,  $p < 0.001$ ) for 10-point accuracy and a kappa of 0.68 (95% C.I. 0.45-0.88,  $p < 0.001$ ) for binary accuracy, while RAG-Top10 showed the highest clarity agreement (kappa = 0.49, 95% C.I. 0.60-0.82,  $p = 0.031$ ). Among tertiary referral center experts, RAG-Top10 demonstrated the highest agreement across all metrics with an ICC of 0.72 (95% C.I. 0.56-0.82,  $p < 0.001$ ) for 10-point accuracy and kappa values of 0.51 (95% C.I. 0.19-0.78,  $p < 0.001$ ) and 0.59 (95% C.I. 0.10-0.92,  $p < 0.001$ ) for binary accuracy and clarity, respectively. For complete results across all configurations, see Table 3.

### 4.3.4 Evaluation of Accuracy and Clarity by Human Graders

The comparative performance in accuracy and clarity grading between the baseline GPT-4 Turbo, RAG-Top1, RAG-Top10, and SFT configurations is summarized in Table 4.4. Among guideline authors, RAG-Top10 showed the highest mean accuracy of 9.45 ( $\pm 1.42$ ), significantly higher ( $p < 0.001$ ) than the baseline GPT-4 Turbo's 6.4 ( $\pm 2.30$ ). RAG-Top1 achieved 8.1 ( $\pm 3.13$ ), while the SFT model performed well with 8.3 ( $\pm 2.03$ ), both significantly outperforming the baseline ( $p = 0.001$  and  $p < 0.001$  respectively). In binary accuracy grading, RAG-Top10 showed the highest performance (91.7% accuracy), followed by RAG-Top1 (81.7%) and SFT (71.7%), all significantly outperforming the baseline (36.6%,  $p < 0.001$ ). Regarding clarity, RAG-Top10 (91.7%), RAG-Top1 (86.6%), and SFT

(88.3%) all significantly outperformed the baseline (46.6%,  $p < 0.001$ ). Among tertiary referral center experts, RAG-Top10 demonstrated the best performance with a mean accuracy of 8.40 ( $\pm 0.99$ ), significantly higher than baseline's 7.23 ( $\pm 1.66$ ,  $p < 0.001$ ) and RAG-Top1's 7.38 ( $\pm 1.42$ ,  $p < 0.001$ ), with no significant difference from SFT's 8.25 ( $\pm 1.43$ ). For binary accuracy, RAG-Top10 achieved the highest rating (93.3%), significantly outperforming baseline (50%,  $p < 0.001$ ), RAG-Top1 (60%,  $p < 0.001$ ), and SFT (76.7%,  $p = 0.010$ ). In clarity evaluations, RAG-Top10 again led with 96.7%, significantly better than baseline (65%,  $p < 0.001$ ) and RAG-Top1 (75%,  $p < 0.001$ ), with no significant difference from SFT (83.3%).

Configuration	10-Point Scale Accuracy Grading		Binary Accuracy Grading		Binary Clarity Grading	
	ICC (95% C.I.)	Significance	Kappa (95% C.I.)	Significance	Kappa (95% C.I.)	Significance
<b>Expert Hepatologists (Guideline Authors)</b>						
Baseline	0.76 (0.49-0.91)	$p < 0.001$	0.30 (0.09-0.50)	$p = 0.049$	0.10 (-0.13-0.45)	NS
RAG-Top1	0.87 (0.76-0.94)	$p < 0.001$	0.68 (0.45-0.88)	$p < 0.001$	0.42 (0.16-0.67)	$p = 0.002$
RAG-Top10	0.83 (0.61-0.93)	$p < 0.001$	0.41 (0.20-0.65)	$p = 0.01$	0.49 (0.60-0.82)	$p = 0.031$
SFT	0.80 (0.70-0.90)	$p < 0.001$	0.39 (0.21-0.50)	$p = 0.042$	0.41 (0.20-0.59)	$p = 0.032$
<b>Expert Hepatologists (Tertiary Referral Center)</b>						
Baseline	0.24 (0.10-0.47)	$p = 0.004$	0.24 (-0.04-0.48)	NS	0.12 (-0.16-0.33)	NS
RAG-Top1	0.41 (0.10-0.71)	$p < 0.001$	0.43 (0.10-0.67)	$p = 0.033$	0.45 (0.30-0.60)	$p = 0.026$
RAG-Top10	0.72 (0.56-0.82)	$p < 0.001$	0.51 (0.19-0.78)	$p < 0.001$	0.59 (0.10-0.92)	$p < 0.001$
SFT	0.47 (0.22-0.65)	$p < 0.001$	0.41 (0.22-0.69)	$p = 0.037$	0.45 (0.10-0.60)	$p = 0.031$

**Table 4.3:** Inter-grader agreement was calculated using the Intraclass Correlation Coefficient (ICC) for the 10-Point Likert Scale and using the Fleiss' Kappa for the binary evaluation of accuracy and clarity. Abbreviations: 95% C.I.: 95% Confidence Interval; NS: not significant.

#### 4.3.5 Model Performance on Simulated Clinical Scenarios

To establish a gold standard for direct-acting antivirals (DAAs) regimen selection, recommendations were determined by majority vote among four hepatologists from tertiary referral centers across 25 simulated clinical scenarios. Expert consensus was achieved for 24 (96%) sofosbuvir-alone recommendations, 20 (80%) sofosbuvir-velpatasvir recommendations, 23 (92%) glecaprevir-pibrentasvir recommendations, 20 (80%) grazoprevir-elbasvir recommendations, and all 25 (100%) decisions regarding ribavirin addition. For cases lacking initial majority consensus, subsequent expert review yielded final gold standard prescribing patterns: sofosbuvir-alone in 0 (0%) cases, sofosbuvir-velpatasvir in 23 (92%) cases, sofosbuvir-velpatasvir-voxilaprevir in 7 (28%) cases, glecaprevir-pibrentasvir in 18 (72%) cases, grazoprevir-elbasvir in 3 (12%) cases, and ribavirin in 0 (0%) cases. We evaluated various model configurations by conducting ten independent interactive sessions per configuration across all 25 clinical scenarios, using identical prompts throughout. Model performance was assessed using two complementary metrics, with median accuracies reported. When evaluating the primary metric of "at least one recommended treatment without any incorrect recommendations"

(Figure 4A), performance varied significantly by configuration. The baseline model achieved a median accuracy of 30%, while the RAG-Top1 configuration demonstrated substantially improved performance at 70% (vs. baseline,  $p < 0.001$ ). The RAG-Top10 configuration exhibited the highest performance with 87% accuracy (vs. baseline,  $p < 0.001$ ), while the SFT model achieved 68% accuracy (vs. baseline,  $p < 0.001$ ). When evaluating the more stringent metric of "all recommended treatments without any incorrect recommendations" (Figure 4B), we observed significant performance differences across model configurations. The baseline model demonstrated a median accuracy of 24.0% the RAG-Top1 configuration showed marked improvement with a median accuracy of 58.0% (vs. baseline,  $p < 0.001$ ), the RAG-Top10 configuration emerged as the superior performer, achieving a median accuracy of 76.0% (vs. baseline,  $p < 0.001$ ), and the SFT model demonstrated a median accuracy of 66.0% (vs. baseline,  $p < 0.001$ ).

Configuration	10-Point Scale Accuracy Grading		Binary Accuracy Grading		Binary Clarity Grading	
	Mean ( $\pm$ SD)	Significance	Number (%)	Significance	Number (%)	Significance
<b>Expert Hepatologists (Guideline Authors)</b>						
Baseline	6.4 ( $\pm$ 2.30)	vs. RAG-Short: $p = 0.001$ vs. RAG-Long: $p < 0.001$ vs. SFT: $p < 0.001$	22 (36.6%)	vs. RAG-Short: $p < 0.001$ vs. RAG-Long: $p < 0.001$ vs. SFT: $p < 0.001$	28 (46.6%)	vs. RAG-Short: $p < 0.001$ vs. RAG-Long: $p < 0.001$ vs. SFT: $p < 0.001$
RAG-Top1	8.1 ( $\pm$ 3.13)	vs. Baseline: $p = 0.001$ vs. RAG-Long: $p = 0.005$ vs. SFT: NS	49 (81.7%)	vs. Baseline: $p < 0.001$ vs. RAG-Long: NS vs. SFT: NS	52 (86.6%)	vs. Baseline: $p < 0.001$ vs. RAG-Long: NS vs. SFT: NS
RAG-Top10	9.45 ( $\pm$ 1.42)	vs. Baseline: $p < 0.001$ vs. RAG-Short: $p = 0.005$ vs. SFT: $p < 0.001$	55 (91.7%)	vs. Baseline: $p < 0.001$ vs. RAG-Short: NS vs. SFT: NS	55 (91.7%)	vs. Baseline: $p < 0.001$ vs. RAG-Short: NS vs. SFT: NS
SFT	8.3 ( $\pm$ 2.03)	vs. Baseline: $p < 0.001$ vs. RAG-Short: NS vs. RAG-Long: $p < 0.001$	43 (71.7%)	vs. Baseline: $p < 0.001$ vs. RAG-Short: NS vs. RAG-Long: NS	53 (88.3%)	vs. Baseline: $p < 0.001$ vs. RAG-Short: NS vs. RAG-Long: NS
<b>Expert Hepatologists (Tertiary Referral Center)</b>						
Baseline	7.23 ( $\pm$ 1.66)	vs. RAG-Top1: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: $p < 0.001$	30 (50%)	vs. RAG-Top1: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: $p = 0.002$	39 (65%)	vs. RAG-Top1: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: NS
RAG-Top1	7.38 ( $\pm$ 1.42)	vs. Baseline: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: $p = 0.001$	36 (60%)	vs. Baseline: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: NS	45 (75%)	vs. Baseline: NS vs. RAG-Top10: $p < 0.001$ vs. SFT: NS
RAG-Top10	8.40 ( $\pm$ 0.99)	vs. Baseline: $p < 0.001$ vs. RAG-Top1: $p < 0.001$ vs. SFT: NS	56 (93.3%)	vs. Baseline: $p < 0.001$ vs. RAG-Top1: $p < 0.001$ vs. SFT: $p = 0.010$	58 (96.7%)	vs. Baseline: $p < 0.001$ vs. RAG-Top1: $p < 0.001$ vs. SFT: NS
SFT	8.25 ( $\pm$ 1.43)	vs. Baseline: $p < 0.001$ vs. RAG-Top1: $p = 0.001$ vs. RAG-Top10: NS	46 (76.7%)	vs. Baseline: $p = 0.002$ vs. RAG-Top1: NS vs. RAG-Top10: $p = 0.010$	50 (83.3%)	vs. Baseline: NS vs. RAG-Top1: NS vs. RAG-Top10: NS

**Table 4.4:** Comparisons of accuracy and clarity grading between the baseline model and the RAG and SFT model pipelines. Statistical testing is based on pairwise comparison (paired Student's T Test for continuous variables and Chi-Square Test for binary variables). According to the Bonferroni correction a  $p$ -value  $< 0.012$  can be considered statistically significant. Abbreviations: RAG: Retrieval Augmented Generation; SFT: Supervised Fine-Tuning.

## 4.4 Chapter’s Deliverables

This section distills the experimental program into concrete choices for building reliable, guideline-aligned LLM for digestive diseases. In particular:

- **Task and Evaluation Setup:**
  - **External knowledge:** 2020 EASL HCV guideline (cleaned; non-text tables converted to structured text lists).
  - **Question set:** 15 open-ended items across four domains and three perspectives, plus 25 simulated clinical scenarios for DAA selection.
  - **Grading:** expert hepatologists (guideline authors and tertiary-center clinicians), blinded; accuracy (binary + 10-point), clarity (binary); agreement via ICC and Fleiss’  $\kappa$ .
  - **Optimization signals:** embedding-based cosine similarity to expert answers for chunking and hyperparameter selection.
- **Design finding #1 — Retrieval format and depth matter:**
  - Paragraph-based chunking outperforms sentence- and fixed-length approaches; performance plateaus around Top-10 retrieved chunks.
  - Recommended decoding ranges for accuracy-critical use: temperature 0–0.8, top-p 0–0.5 (RAG-Top10 most robust to variation).
- **Design finding #2 — RAG and SFT both help; RAG-Top10 leads:**
  - Against baseline GPT-4 Turbo, RAG-Top10 attains the highest expert-graded accuracy and clarity; RAG-Top1 and SFT also significantly improve over baseline.
  - Graded (10-point) accuracy shows good-to-excellent ICC, while binary metrics yield only fair/moderate  $\kappa$ —favoring scale-based evaluation for nuanced clinical judgments.
- **Clinical decision support impact (25 cases):**
  - On “ $\geq 1$  correct regimen without incorrects,” median accuracy rises from 30% (baseline)  $\rightarrow$  70% (RAG-Top1)  $\rightarrow$  87% (RAG-Top10)  $\rightarrow$  68% (SFT).
  - On “all correct regimens without incorrects,” median accuracy rises from 24%  $\rightarrow$  58%  $\rightarrow$  76%  $\rightarrow$  66%, respectively.
- **Method guidance for deployment:**
  - Ground responses in authoritative guidelines via RAG; prefer paragraph chunks with  $\sim$  Top-10 retrieval.
  - Keep sampling low-variance decoding (temperature/top-p as above).
  - Use heuristic prompting that enforces structured, guideline-aligned outputs.
  - Validate with blinded expert grading before clinical use; complement with similarity metrics for tuning but do not substitute them for factual accuracy checks.
- **Scope & limitations (as reported):** Focus on HCV; GPT-family models only; no RLHF; one-guideline context per query; persistent accuracy gaps remain in edge cases.

**Bottom line:** With cleaned, consistently structured guideline text and paragraph-based **Top-10** retrieval under low-variance decoding, GPT-4 Turbo transitions from suboptimal baseline behavior to **clinically usable, guideline-faithful** recommendations—supporting safe integration into hepatology CDSS.



## **Chapter 5**

### **Verify: The Expert-of-Experts Verification and Alignment (EVAL) Framework for Safe, Aligned and Automatically Verified Outputs**

**“Precaution is better than cure”**

*– Edward Coke*

## 5.1 Chapter’s Overview<sup>5</sup>

In Chapter 2, it emerged that baseline decoder-only LLMs are unreliable for safe use in digestive diseases, while accuracy improves when guideline context is supplied. Chapter 3 then showed that transforming guidelines into LLM-friendly text and coupling them with RAG and principled prompting can yield near-perfect, guideline-faithful answers for HCV. Chapter 4 extended this by benchmarking retrieval design, hyperparameters, and SFT against blinded hepatology experts, confirming that expert-validated, guideline-grounded systems outperform unguided models.

Building on these precursors, this chapter introduces an Expert-of-Experts Verification and Alignment (EVAL) framework for AI safety in gastroenterology. The framework anchors “ground truth” to free-text answers authored by guideline leaders (“golden labels”) and evaluates safety at two levels: model-level ranking via unsupervised embedding similarity to expert answers, and answer-level screening via a trained reward model that filters inaccurate responses across temperature regimes. As a distinct test bed beyond HCV, the use case is upper gastrointestinal bleeding (UGIB)—a high-stakes condition with an incidence up to 116 per 100,000<sup>85</sup> and mortality up to 11%<sup>86</sup>, governed by robust international guidelines yet marked by variable adherence ( $\approx 14.3\text{--}95.7\%$ ) and low uptake of guideline risk scores ( $\sim 30\%$  of practitioners).<sup>87–89</sup>

The primary aim is to establish an expert-anchored, scalable framework that (i) ranks LLM configurations by alignment to guideline-author answers and (ii) automatically screens individual outputs for accuracy, thereby improving the safety of provider-facing LLMs in gastroenterology.

Specific objectives include:

1. Define expert ground truth: Collect “golden-label” free-text answers from UGIB guideline senior authors across pre-, intra-, and post-endoscopic care.
2. Assemble multi-tier benchmarks: Evaluate across three datasets—expert-generated UGIB questions (n=13), ACG multiple-choice questions (n=40), and real-world trainee questions from simulation (n=117).

---

<sup>5</sup>This chapter (text and images) is adapted from the article: **Giuffrè M., You K., Pang Z., et al.** *Expert of Experts Verification and Alignment (EVAL) Framework for Large Language Models Safety in Gastroenterology*. *NPJ Digit Med*. 2025 May 3;8(1):242. doi: 10.1038/s41746-025-01589-z. The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit. Changes: Methods and Results are reproduced verbatim; Introduction/Background, Discussion, and narrative transitions are reformulated for thesis style; figure/table numbering and layout adapted to the thesis format.

3. Rank models with unsupervised similarity: Compare TF-IDF, Sentence Transformers, and fine-tuned ColBERT to identify which metric best tracks human-graded accuracy and MCQ performance.
4. Train a reward model for answer-level safety: Use expert-graded outputs to build a grader that detects accurate vs. inaccurate responses across temperature ranges and enables automated rejection sampling.
5. Compare configurations: Benchmark baseline, RAG, SFT, and RAG+SFT across GPT, Claude, Llama, and Mistral families, using guideline-reformatted corpora.
6. Quantify real-world applicability: Test whether model-level ranking plus answer-level filtering improves accuracy on free-form, clinician-posed questions relevant to UGIB management.

## 5.2 Materials and Methods

### 5.2.1 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>37, 90-94</sup> Following our previously published protocol<sup>42,95</sup>, we reformatted the original documents from raw PDF formats to ones suitable for LLMs, as described elsewhere.<sup>42</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.

#### 5.2.1.1 Retrieval Augmented Generation

For retrieval augmented generation (RAG)<sup>34</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 128000 tokens, whereas Anthropic’s Claude-3-Opus has a context window of up to 200000 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.

### 5.2.1.2 Supervised Fine-Tuning

Supervised fine-tuning was performed using low-rank adaptation (LoRA)<sup>96,97</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.

In the fine-tuning process, we start with a dataset  $D$  consisting of relevant documents and an existing parametric language model  $p_\theta$ . The goal is to update the model’s parameters from  $\theta$  to  $\theta'$  using the documents contained in dataset  $D$ . This update aims to produce a new language model that can more effectively answer questions requiring knowledge from the dataset  $D$ . To illustrate the fine-tuning process, consider a single document  $d_i$ . This document can be represented as a sequence of tokens  $[\langle s \rangle, t_1, t_2, \dots, t_n]$ , where  $\langle s \rangle$  is the start token. The current language model’s performance on this sequence can be evaluated by predicting each token using the preceding tokens:

$$Perplexity = \exp\left(-\frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i | x_{<i})\right),$$

which is known as the perplexity metric in language modeling. Perplexity measures the model’s uncertainty or surprise in predicting token  $t_i$  after observing a sequence of tokens  $t_1, \dots, t_{i-1}$ . Generally, a well-aligned model  $p_\theta$  with the corpus  $D$  is expected to have lower perplexity values. The differentiability of the perplexity objective allows it to be used as a loss function for the language model  $p_\theta$ . Minimizing the perplexity on the documents  $D$  during training leads to a model  $p_\theta$  that is more likely to produce responses aligned with the text in the corpus. Large language models are typically very large, with many open-source models ranging from 7 to 100 billion parameters. Training such models requires substantial GPU memory, typically between 48GB and 500GB, depending on the model size and architecture. Consequently, training necessitates large GPU clusters, making the process resource-

intensive and costly. To address these challenges, two common techniques are often employed: (1) quantization and (2) low-rank adaptation. Quantization is a process of reducing the model’s weights from 16-bit floating point to 8-bit or 4-bit floating point precision, which has been shown to maintain performance with minimal degradation. Low-rank adaptation is based on the concept that only a small subset of the model’s parameters  $\theta$  needs to be fine-tuned. Tuning as little as 0.01% of the model’s parameters can significantly reduce perplexity on the document corpus. Further details of this training procedure are provided below.

For Low-Rank Adaptation (LoRA) for Efficient Fine-Tuning, The model  $p_\theta$  consists of multiple attention blocks, each employing query, key, and value projections, represented by the projection matrices  $W_Q$ ,  $W_K$ , and  $W_V$  of dimension  $d \times k$ . The large dimensionality of these matrices can lead to high computational costs during fine-tuning. To mitigate this, LoRA is utilized to reduce the dimensionality of the training space. In LoRA, the projection matrices are frozen, and we introduce randomly initialized trainable projection matrices  $A$  of size  $r \times k$  and  $B$  of size  $d \times r$ , both of which are filled with zeros in the beginning. Typically,  $d$  and  $k$  are large values, while  $r$  is chosen to be much smaller than  $d$  and  $k$  to reduce the number of trainable parameters. The original projection matrices  $W_Q$ ,  $W_K$ , and  $W_V$  remain fixed, while the product  $BA$  is added to these matrices. Hence, for any text embedding vector  $x$  passed through a projection matrix  $W$ , it can be decomposed into the sum of the frozen and trainable portions:

$$Wx = W_0x + BAx = (W_0 + BA)x$$

LoRA has advantages for fine-tuning because  $W_0$  can be pre-trained and kept fixed, while  $BA$ , with its significantly smaller number of parameters, can be easily trained and integrated with  $W_0$ , thereby necessitating fewer computational resources.

### 5.2.2 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 5.1). These cases were specifically designed to test the ability to integrate multiple guideline recommendations in realistic clinical contexts.

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 hours? Should she undergo red blood cell transfusion or upper endoscopy within 24 hours?
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?

**Table 5.1: List of Expert-Generated Questions for Upper Gastrointestinal Bleeding Management.** 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).

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 elsewhere.<sup>98</sup>

Using these expert-curated questions, we also generated responses using all LLM configurations at a temperature setting of 0.8<sup>99</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 elsewhere.<sup>98</sup> 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>100</sup> Each clinical case-question pair is reported elsewhere.<sup>98</sup>

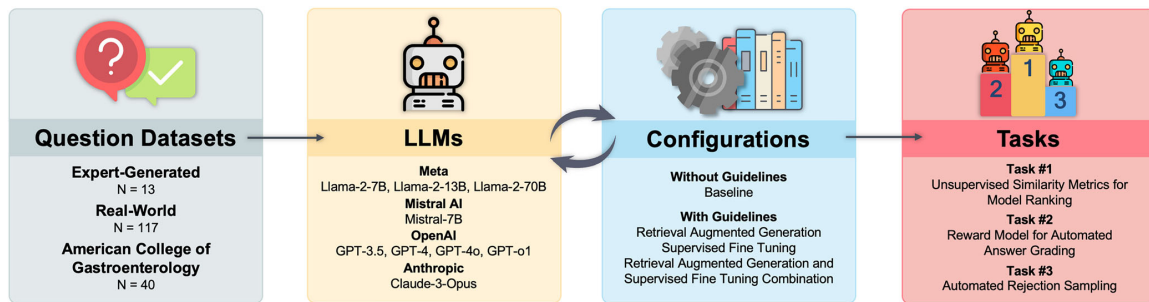
Each LLM configuration was tested using a heuristic prompting approach, necessary due to the unpredictable nature of trainee questions. The prompts were structured to include complete clinical case analysis, providing all relevant context (including patient demographics, laboratory findings, and clinical 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.

### 5.2.3 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 Figure 5.1.



**Figure 5.1:** 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/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.

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 Figure 5.2.

We used ColBERT<sup>101</sup> to quantify the alignment between responses generated by LLMs and responses by experts (Figure 5.2). 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 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>102</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 Table 5.2). 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. 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>103</sup>, a common existing LLM-based method for textual similarity, and TF-IDF<sup>104</sup>, which is a classical method based on word and document statistics.

Question	Triplet Example ColBERT		
	Expert Answer	Guidelines Text (Related to Query)	Guidelines Text (Unrelated to Query)
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?	“Glasgow-Blatchford Score. Most recommend score 0-1 as threshold”.	“We suggest that patients presenting to the emergency department with UGIB who are classified as very low risk, defined as a risk assessment score with $\leq 1\%$ false negative rate for the outcome of hospital-based intervention or death (e.g., Glasgow-Blatchford score = 0–1), be discharged with outpatient follow-up rather than admitted to hospital (conditional recommendation, very-low-quality evidence)”.	“We suggest a restrictive policy of red blood cell (RBC) transfusion with a threshold for transfusion at a hemoglobin of 7 g/dL for patients with UGIB (conditional recommendation, low-quality evidence). Summary of evidence: For the general population of patients with anemia, not restricted to UGIB, current US guidelines make a strong recommendation for a restrictive RBC transfusion threshold of 7 g/dL in hospitalized hemodynamically stable patients, including critical care patients, and a threshold of 8 g/dL in those undergoing orthopedic or cardiac surgery and those with existing cardiovascular disease”.

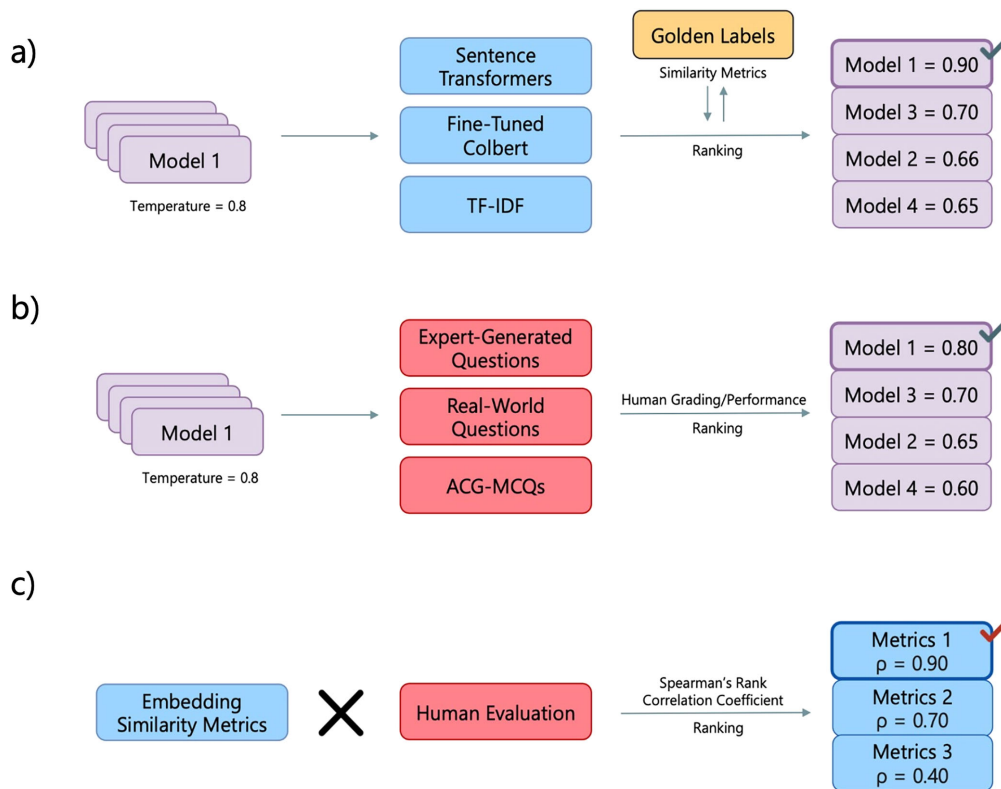
**Table 5.2:** Triplet Example for Fine-tuning Colbert with Expert Answers. This table illustrates the structure of triplets used for fine-tuning the ColBERT model. Each triplet consists of an expert-generated answer (gold standard), a related guideline text (providing context aligned with the query), and an unrelated guideline text (serving as a negative example). These triplets enable the ColBERT model to differentiate between relevant and irrelevant textual information, enhancing its alignment with expert guidance.

For the Sentence Transformers-based similarity metrics, we use the publicly available pre-trained embedding model, all-MiniLM-L6-v2, from Sentence Transformer<sup>105</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>106</sup> Lastly, we calculate the cosine similarity between the calculated TF-IDF score to serve as the final similarity score.

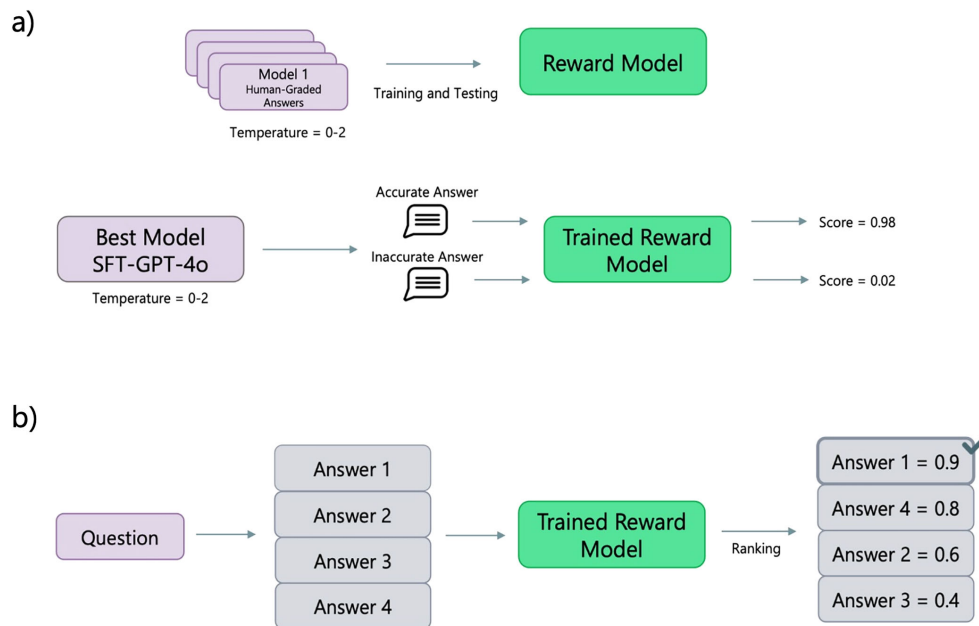
For each similarity method, we performed pairwise t-tests comparing the highest-scoring model configuration against all other configurations individually. Similarly, we conducted pairwise t-tests for human-graded accuracies across the three evaluation sets (expert-generated questions, real-world questions, 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.



**Figure 5.2:** 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.

### 5.2.4 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>107</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 Figure 5.3.



**Figure 5.3:** 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.

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 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>108,109</sup> For larger or more nuanced topics, a dataset exceeding 50000 pairs may be necessary.<sup>110</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>111–</sup>  
<sup>113</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 less than 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 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.

### 5.2.5 Automated Rejection Sampling

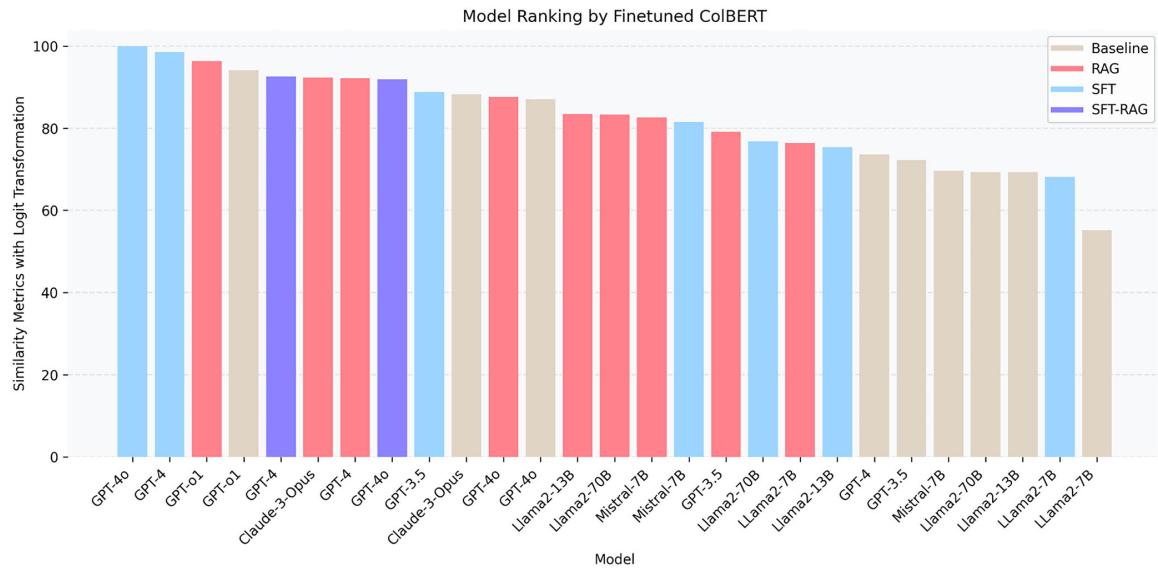
Extending the reward model pipeline, we can incorporate the reward function directly into the answer pipeline by using a rejection sampling approach. For each question, the LLM agent generates  $K$  candidate answers. These  $K$  answers are evaluated by the reward model, and only the top-scoring answer 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 elsewhere.<sup>98</sup> 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 elsewhere<sup>98</sup> to illustrate the trade-offs and performance improvements.

## 5.3 Results

### 5.3.1 Model Ranking by Similarity Metrics

In terms of model ranking by similarity metrics (Table 5.3), Claude-3-Opus in the baseline configuration achieved the best performance in both 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 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 under RAG configuration showing a less stringent statistical significance ( $p < 0.05$ ). For 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 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 ( $0.679 \pm 0.006$ ), GPT-o1 RAG ( $0.687 \pm 0.004$ ), 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 ( $0.681 \pm 0.015$ ).

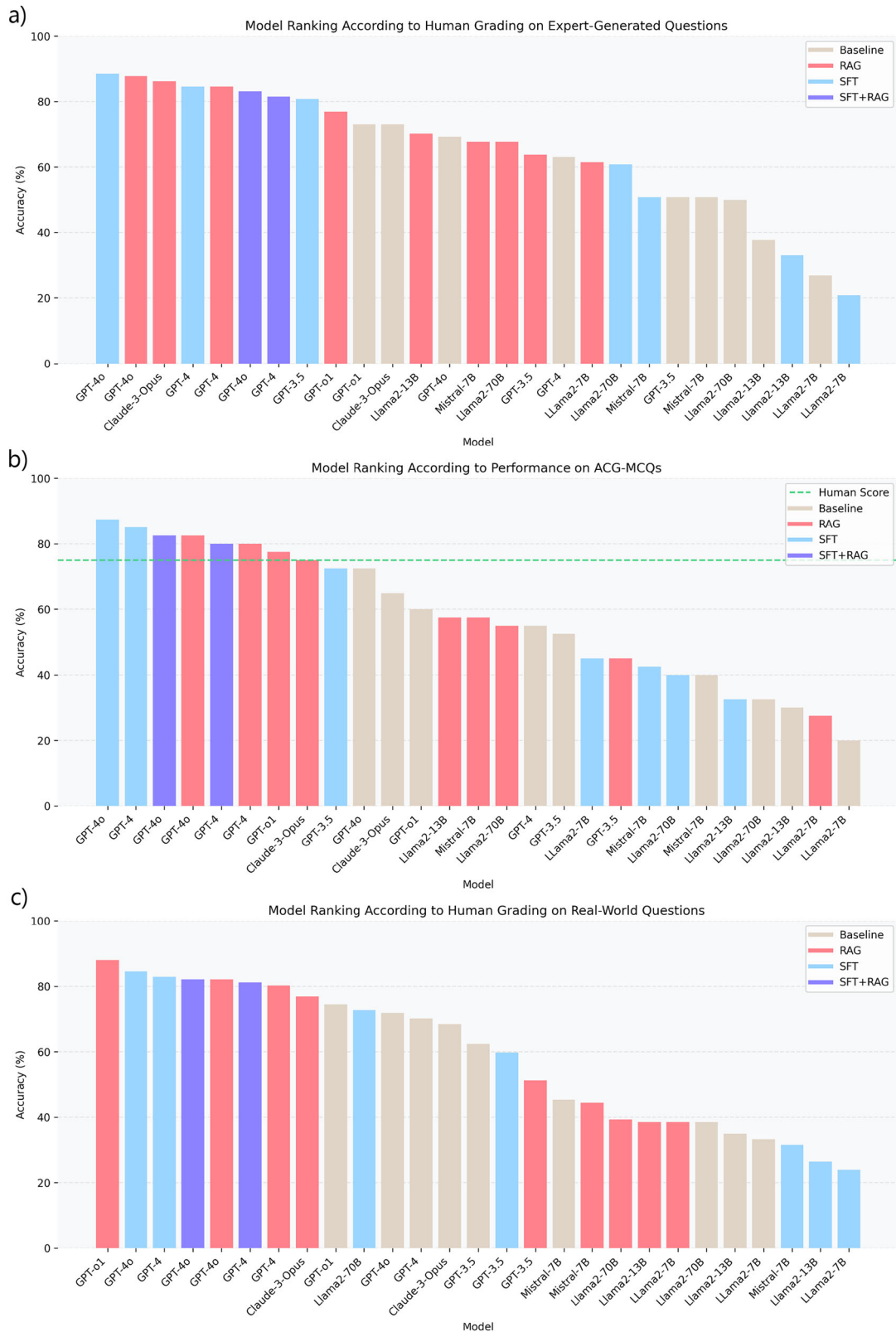


**Figure 5.4:** Model Performance Ranking Based on Fine-tuned ColBERT Similarity Scores. The figure shows the ranking of different LLM configurations based on their similarity to expert-generated responses, as measured by Fine-tuned ColBERT scores after logit transformation. Models are grouped by configuration 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.

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 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 ColBERT scores in Figure 5.4, as this transformation maintains the monotonic relationship between scores while providing better visual differentiation between high-performing models.

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 N (%)	ACG-MCQs Performance N (%)	Real-World Questions N (%)
<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>

**Table 5.3: Model Ranking Comparison across similarity-based metrics, human grading, and performance of multiple-choice questions (MCQs) dataset.** 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.



**Figure 5.5:** 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.

### 5.3.2 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 5.3 and depicted in Figure 5.5. 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 all 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 AC3G-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 all 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%), RAG-Claude-3-Opus (76.9%), 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 baseline RAG-Claude-3-Opus (76.9%,  $p < 0.05$ ), with all other model configurations demonstrated statistically significant differences ( $p < 0.01$ ).

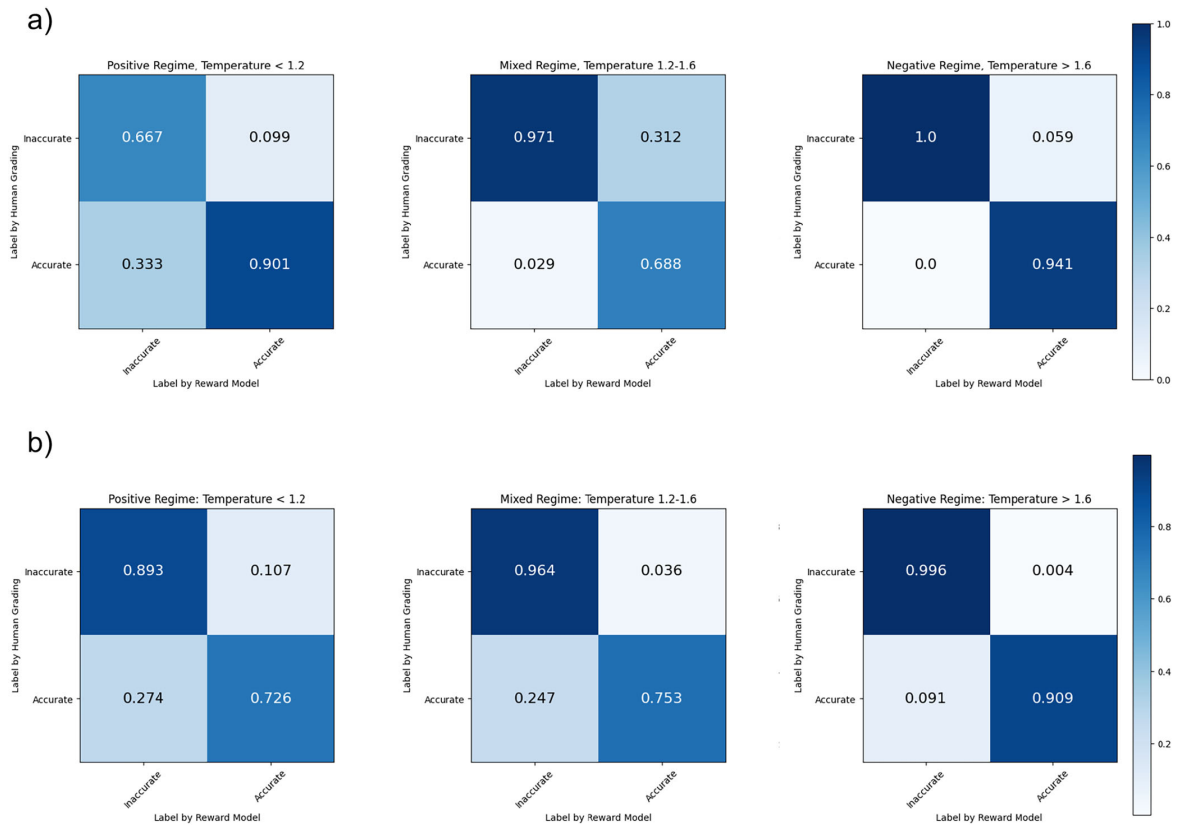
### 5.3.3 Alignment between Similarity Metrics and Human Performance

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.

### 5.3.4 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 elsewhere.<sup>98</sup> 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 (Figure 5.6). 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 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.



**Figure 5.6:** 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.

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 Table 5.4 and Table 5.5 respectively.

Temperature Value	True Labels RAG-GPT-4	True Labels SFT-GPT-4o
T = 0.0	0.938	0.754
T = 0.2	0.900	0.769
T = 0.4	0.884	0.731
T = 0.6	0.900	0.731
T = 0.8	0.892	0.731
T = 1.0	0.892	0.746
T = 1.2	0.646	0.769
T = 1.4	0.723	0.823
T = 1.6	0.915	0.892
T = 1.8	0.992	0.985
T = 2.0	0.992	1.000

**Table 5.4:** Alignment, expressed as the number of true labels, among the Reward Model and RAG-GPT-4 and SFT-GPT-4o for expert-generated questions. Abbreviations: RAG: Retrieval Augmented Generation; SFT: Supervised Fine-Tuning.

Model	Baseline	RAG	SFT	RAG-SFT
Llama-2-7B	0.512	0.513	0.462	N/A
Llama-2-13B	0.538	0.521	0.265	N/A
Llama-2-70B	0.547	0.555	0.453	N/A
Mistral-7B	0.598	0.547	0.521	N/A
Claude-3-Opus	0.761	0.769	N/A	N/A
GPT-3.5	0.743	0.769	0.547	N/A
GPT-4	0.769	0.855	0.744	0.667
GPT-4o	0.777	0.812	0.821	0.761
GPT-o1	0.803	0.889	N/A	N/A

**Table 5.5:** Alignment, expressed as the number of true labels, among the Reward Model and all model configurations in real-world question-answer pairs at temperature = 0.8. Abbreviations: RAG: Retrieval Augmented Generation; SFT: Supervised Fine-Tuning; N/A: Not Available.

### 5.3.5 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 5.6). 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% to 54.9%) and of 23.60% for SFT-GPT-4o (increasing accuracy from 64.4% to 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.

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	<b>0.559</b>	<b>0.890</b>	<b>0.549</b>	<b>0.240</b>
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%

**Table 5.6:** Rejection Sampling for automated grading. 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.

## 5.4 Chapter’s Deliverables

In light of Chapters 2–4—where baseline decoder-only LLMs proved insufficient, and guideline-grounded retrieval plus expert validation emerged as prerequisites for safety—this chapter formalizes a scalable, expert-anchored evaluation stack. The EVAL framework operationalize “ground truth” via golden-label answers from guideline leaders and test whether automated proxies (embedding similarity, reward modeling, and rejection sampling) can reliably triage models and suppress unsafe outputs across decoding regimes. What follows distills the empirical signals into deployment-ready guidance for building, ranking, and safeguarding LLM-based CDSS in gastroenterology. In particular:

- **Task and Benchmark Setup:**
  - Expert ground truth: Free-text “golden labels” from UGIB guideline senior authors; used to anchor accuracy.
  - Datasets: 13 expert-generated questions; 40 ACG UGIB MCQs; 117 real-world trainee questions from simulation.
  - Model families & configs: GPT-3.5/4/4o/o1, Claude-3-Opus, Llama-2 (7B/13B/70B), Mistral-7B; evaluated as baseline, RAG, SFT, and RAG+SFT (subject to vendor constraints).
- **Finding #1 — Fine-tuned ColBERT best reflects human accuracy:** Fine-tuned ColBERT showed the strongest alignment with human evaluation: Spearman  $\rho = 0.91$  (expert-generated), 0.86 (ACG-MCQs), 0.81 (real-world); Sentence Transformers: moderate; TF-IDF: weakest. Fine-Tuned ColBERT can be used for model-level ranking, enabling scalable identification of top configurations without manual grading.
- **Finding #2 — Best-performing configurations differ by task:**
  - SFT-GPT-4o ranked best on expert-generated (88.5%) and ACG-MCQs (87.5%).
  - RAG-GPT-o1 ranked best on real-world questions (88.0%).

- Several advanced GPT/Claude configurations were not significantly different from the top performer on expert/MCQ tasks, indicating a competitive cohort at the high end.
- **Finding #3 — Reward model enables answer-level safety controls:**
  - Reward model (OPT-350M) reproduced human labels in:
    - RAG-GPT-4 internal validation: 87.9% overall; 90.0% at  $T < 1.2$ ; 99.2% at  $T > 1.6$ .
    - SFT-GPT-4o external validation: 81.8% overall; robust detection of inaccurate answers across regimes.
  - Sensitivity across temperature showed strongest filtering in high-temperature (“negative”) regimes, where hallucination risk is highest.
- **Finding #4 — Rejection sampling improves accuracy, especially when harder:**
  - With  $K=5$  candidates and reward-model filtering:
    - RAG-GPT-4: +9.39% overall (to 0.559); +98.35% in high-temperature regime.
    - SFT-GPT-4o: +8.36% overall (to 0.598); +23.60% in mixed regime; +121.43% in high-temperature regime.
  - Practical implication: automated filtering can partially rescue accuracy under exploratory decoding.
- **Operational takeaways for deployment:**
  - Use expert-authored golden labels to stand up a scalable model-ranking pipeline (fine-tuned ColBERT recommended).
  - Add an answer-level reward model and rejection sampling to screen outputs—particularly valuable when temperature  $\geq 1.2$  or when prompts are less structured.
  - Expect task-dependent winners (expert/MCQ vs. free-form clinical questions); maintain configuration portfolios rather than a single “best” model.
  - Continue to ground CDSS in reformatted, guideline-centric corpora (as established in Chapters 3–4).
- **Scope & limitations:**
  - UGIB-focused use case; external validity to other conditions requires new expert labels and guideline corpora.
  - Training/test split aligned with geography (US vs. EU/AP) may introduce subtle biases.
  - Reward-model alignment, while strong, is not a substitute for human oversight in high-stakes care.

**Bottom line:** An expert-anchored, two-level safety stack—**model-level ranking via fine-tuned ColBERT plus answer-level filtering via a reward model with rejection sampling**—can identify reliable configurations and suppress inaccurate outputs. Coupled with LLM-friendly, guideline-

grounded retrieval (Chapter 3) and expert validation (Chapter 4), EVAL offers a practical path to safer, evidence-aligned LLM use in gastroenterology CDSS.



## **Chapter 6**

**Scale: From Single-Disease Prototyping to Society-Level Benchmarking: Guideline-Grounded Testing on European Association for the Study of the Liver Multiple-Choice Benchmark**

**“More is different”**

*– Philip W. Anderson*

## 6.1 Chapter's Overview<sup>6</sup>

In the preceding chapters, baseline decoder-only LLMs proved insufficient for safe clinical use; accuracy rose materially only when guideline knowledge was injected (Chapter 3), structured and optimized via chunking/hyperparameters with expert validation (Chapter 4), and safeguarded with expert-anchored ranking and reward-model filtering (Chapter 5).

This chapter extends that trajectory to a comprehensive, guideline-referenced benchmark for hepatology using the publicly available EASL Campus MCQs. The aim is to move beyond single-disease probes and test whether domain-knowledge injection (RAG, SFT, and RAG+SFT) reliably lifts performance across liver medicine, including items with images, under standardized prompting and controlled decoding.

The primary aim is to quantify how much guideline grounding (RAG/SFT) improves LLM accuracy on hepatology MCQs versus pooled physician performance, and identify configuration(s) that generalize across subdomains.

Specific objectives include:

1. Build a clean, EASL-aligned MCQ benchmark (n=110 after exclusions) and a human baseline (mean physician accuracy 56.9%).
2. Compare leading model families (Gemini-1.5-Pro, Claude-3-Opus, GPT-4o) across four configurations: baseline, RAG, SFT (GPT-4o), and RAG+SFT with fixed hyperparameters (temperature 0.8; top-p 0.5).
3. Standardize ingestion of text-only and image-containing questions (base64 workflow) under a uniform “answer-letter-only” prompt.
4. Estimate absolute and relative accuracy gains versus physicians and test statistical significance (Fisher’s exact).
5. Run subdomain sensitivity (tumors; viral hepatitis; cirrhosis/complications; metabolism/alcohol/toxicity; immune/cholestatic; general hepatology) to surface areas of strength/weakness.

---

<sup>6</sup>This chapter (text and images) is adapted from the article: **Giuffrè M., Distefano A., Kresevic S., et al. Guideline-enhanced large language models outperform physician-test takers on EASL Campus quizzes multiple choice questions. *JHEP Rep.* 2025 Jul 14;7(10):101523. doi: 10.1016/j.jhepr.2025.101523.** The article is Open Access under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits use, sharing and reproduction with appropriate credit. Changes: Methods and Results are reproduced verbatim; Introduction/Background, Discussion, and narrative transitions are reformulated for thesis style; figure/table numbering and layout adapted to the thesis format.

6. Distill deployment guidance on when to prefer RAG, SFT, or their combination for guideline-centric knowledge assessments.

## 6.2 Materials and Methods

### 6.2.1 Dataset Creation and Model Comparison

We compiled a dataset of 110 MCQs from all available quizzes up to November 30<sup>th</sup>, excluding questions on basic science ( $n = 2$ ) and rare liver diseases ( $n = 3$ ) due to insufficient knowledge outlined in guideline documents. To establish a human performance benchmark, we used data from previous online physician test-takers. Specifically, for each question, we determined the percentage of participants who answered correctly, then averaged these percentages across all questions to obtain an overall physician accuracy of 56.9%. We evaluated state-of-the-art LLMs (reported in the following section). We incorporated specific domain knowledge from current EASL guidelines (published up to October 15<sup>th</sup>, 2024), using retrieval-augmented generation (RAG), supervised fine-tuning (SFT), only available for GPT-4o, or a combined approach. The accuracy of each LLM configuration was measured as the percentage of correctly answered questions and compared against human performance by physician test-takers, using Fisher's exact test. We defined statistical significance as a two-tailed  $p$  value less than 0.05.

### 6.2.2 Large Language Model Configuration

Experiments were conducted on a local Python (version 3.11) environment, using the respective Action Programming Interface (API) to interact with Google AI's Gemini 1.5 Pro, Anthropic's Claude 3 Opus, and OpenAI's GPT-4o. For all experiments, we used the following hyperparameter: temperature = 0.8 and top\_p = 0.5, max\_token = 150. This hyperparameter combination was previously tested and provided the highest rate of accurate answers. In particular, setting the temperature  $\leq 0.8$  and top\_p = 0.5 helps in generating responses that are more deterministic and consistent, as a lower temperature reduces the model's tendency to explore less likely tokens, thus increasing the fidelity to the external data source. By restricting the token selection to only the most probable option at each step, we prioritize the model's confidence in its outputs, potentially increasing the likelihood of generating responses that closely align with the retrieved information and the given context.

For the retrieval-augmented generation (RAG) framework, we used different embedding models optimized for each LLM according to their respective documentations: "textembedding-3-large" for GPT-4o, "voyage-3-large" for Claude-3-Opus, and "text-embedding004" for Gemini-1.5-Pro. We used the complete corpus of European Association for the Study of the Liver (EASL) guidelines (published up to October 15<sup>th</sup>, 2024) as our reference text. The guidelines were preprocessed into an LLM-friendly version by removing non-informative content (e.g., headers, reference numbers) and converting non-textual sources (e.g., tables, figures) into text format, following previously validated methods. The

preprocessed guidelines were then segmented at the paragraph level to create discrete chunks suitable for embedding. Each chunk was encoded into high-dimensional embeddings using the respective embedding models for each LLM. For all three models, we implemented a RAG configuration that retrieved the top 15 most relevant chunks based on cosine similarity between the query and chunk embeddings. Query embeddings were generated using the same model-specific embedding approaches for similarity matching. This approach allowed us to evaluate the performance of each LLM when augmented with relevant guideline knowledge across all major areas of hepatology.

OpenAI has not allowed fine-tuning GPT-4 to most research groups but only to a selected group of high-volume developers for experimental purposes. Therefore, for the SFT, we used the model “gpt-4-turbo”, according to OpenAI’s nomenclature. The fine-tuning framework was conducted on the OpenAI’s platform through a structured process involving two stages and four curated datasets, all available elsewhere,<sup>98</sup> and developed according to OpenAI’s instructions using the jsonl file format. Each entry in the dataset contained a json object with fields for role, content, and metadata. The training hyperparameters of batch size, learning rate multiplier, and the number of epochs were set to “auto”. Prompting Structure For all experiments, we employed a standardized basic prompt engineering approach. The system prompt was designed to instruct the LLM to act as an expert hepatologist specializing in MCQ evaluation based on EASL guidelines.

The prompt structure was as follows: "You are an expert hepatologist tasked with answering Multiple Choice Questions based on EASL guidelines. Your role is to: Carefully read each question, including any associated image [image input in base64 format] Analyze all provided answer options (A-D or A-E) Provide ONLY the letter corresponding to the correct answer without any additional explanation Answer format: For 4-option questions: select from A, B, C, or D For 5-option questions: select from A, B, C, D, or E IMPORTANT: Respond ONLY with the letter of the correct answer."

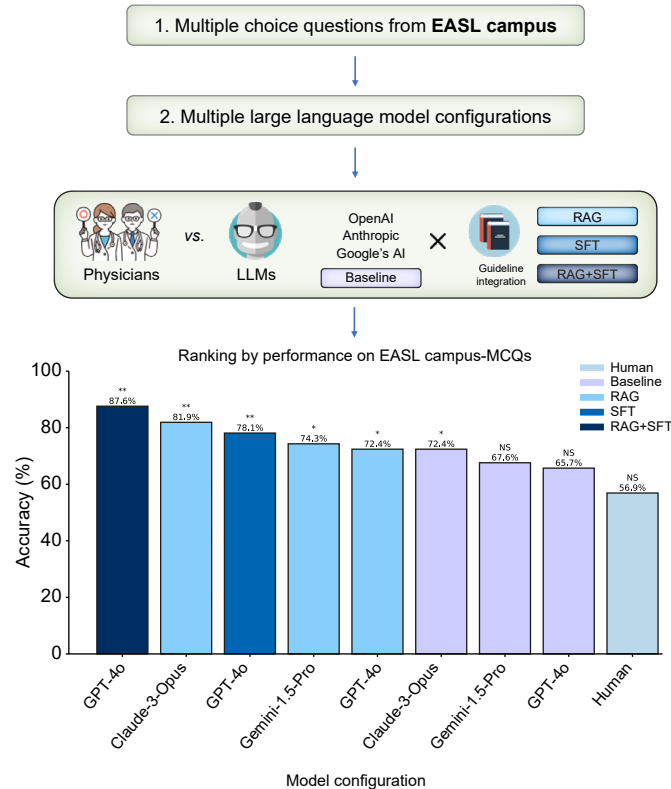
This prompt was consistently used across all LLMs to ensure standardization of the experimental conditions. When images were present in the MCQs, they were encoded in base64 format and included in the prompt structure following each model's specific image handling requirements. This approach was designed to maintain consistency while allowing for the evaluation of both text-only and image-containing questions within the same framework.

For MCQs containing images, these were preprocessed and integrated into the prompts using standardized protocols. Images were converted to base64 format using Python's base64 library. Images were accepted in JPEG format with a maximum file size of 20MB. For optimal processing, images were preprocessed to maintain resolution while reducing file size when necessary.

### 6.3 Results

As shown in Figure 6.3, the baseline performance of each LLM showed varying degrees of improvement over human performance, with GPT-4o and Gemini-1.5-Pro achieving non-statistically significant higher accuracies of 65.7% and 67.6%, respectively, while baseline Claude-3-Opus demonstrated statistically significant higher accuracy than physicians (72.4%,  $p = 0.026$ ). The implementation of RAG consistently enhanced model performance across all models, with RAG-enhanced versions showing significant improvements over physician performance: GPT-4o (72.4%,  $p = 0.026$ ), Gemini-1.5-Pro (74.3%,  $p = 0.017$ ), and Claude-3-Opus (81.9%,  $p < 0.001$ ). SFT-GPT-4o also yielded substantial improvements (78.1%,  $p = 0.002$ ). The combined SFT-RAG-GPT-4o configuration achieved the highest overall accuracy of 87.6% ( $p < 0.001$ ).

We have conducted also a sensitivity analysis of accuracy among hepatology subdomains as categorized by the EALS society. In liver tumors, it achieved 95.0% accuracy compared to 50.7% human performance ( $p < 0.001$ ). For viral hepatitis, the configuration reached 80.0% accuracy versus 55.1% for physicians ( $p < 0.001$ ), while in cirrhosis and complications, it reached 80.0% accuracy compared to 58.2 ( $p = 0.001$ ). In the metabolism, alcohol and toxicity domain, it showed 53.9% accuracy versus 46.3% ( $p = 0.3$ ). For immune-mediated and cholestatic disease, the configuration achieved 100% accuracy compared to 69.6% ( $p < 0.001$ ), and in general hepatology, it reached 100% accuracy versus 61.9% ( $p < 0.001$ ).



**Figure 6.3:** Evaluation framework of liver disease knowledge using EASL Campus multiple choice questions. A total of 105 multiple choice questions were collected and used to assess both human performance and various LLM configurations. The LLMs were tested in three configurations: (1) Baseline performance using standard models, (2) RAG-enhanced models incorporating relevant EASL guideline content, and (3) models improved through SFT or combined RAG-SFT approaches. Performance is shown as percentage accuracy for each configuration compared to human baseline performance. n.s., not significant; \* $p < 0.05$ ; \*\* $p < 0.01$ . LLM, large language model; RAG, retrieval-augmented generation; SFT, supervised fine-tuning.

## 6.4 Chapter's Deliverables

This section translates the EASL-Campus benchmark into actionable guidance for building and evaluating guideline-grounded LLMs in hepatology. Beyond reporting raw accuracies, it clarifies *where* domain knowledge (RAG, SFT, and their combination) delivers the largest, statistically significant gains over pooled physician performance, *which* configurations are most reliable across subdomains and image-based items, and *how* to operationalize stable decoding for safer outputs. Read it as a “deployment brief”: what corpus to use, how to structure retrieval, when to fine-tune, and where to prioritize future curation when performance remains uneven. In particular:

- **Task and Benchmark Setup:**
  - Expert reference: EASL Campus MCQs curated to guideline scope (n=110 after exclusions).
  - Human baseline: pooled online physician accuracy 56.9%.
  - Models/configs: Gemini-1.5-Pro, Claude-3-Opus, GPT-4o; tested as Baseline, RAG(Top-15 paragraph chunks), SFT-GPT-4o, and RAG+SFT-GPT-4o; uniform

prompt; temperature 0.8, top-p 0.5; image items handled via standardized base64 pipeline.

- **Finding #1 — Domain knowledge reliably beats physicians:**
  - Baseline models were comparable to or modestly better than physicians (e.g., GPT-4o 65.7%, Gemini 67.6%; Claude-3-Opus 72.4%,  $p=0.026$ ).
  - RAG lifted accuracy across models: GPT-4o 72.4% ( $p=0.026$ ), Gemini 74.3% ( $p=0.017$ ), Claude-3-Opus 81.9% ( $p<0.001$ ).
  - SFT-GPT-4o: 78.1% ( $p=0.002$ ).
  - RAG+SFT-GPT-4o: 87.6% ( $p<0.001$ ), the best overall.
- **Finding #2 — Gains are uneven by subdomain (where to trust):**
  - Liver tumors: 95.0% vs 50.7% human ( $p<0.001$ ).
  - Viral hepatitis: 80.0% vs 55.1% ( $p<0.001$ ).
  - Cirrhosis & complications: 80.0% vs 58.2% ( $p=0.001$ ).
  - Immune/cholestatic and general hepatology: 100% vs 69.6% / 61.9% ( $p<0.001$ ).
  - Metabolism, alcohol, toxicity: 53.9% vs 46.3% (ns) → an area needing richer context or targeted tuning.
- **Operational guidance:**
  - For knowledge checks/board-style content, prefer RAG+SFT where feasible; RAG alone is a strong default when fine-tuning is not available.
  - Maintain subdomain dashboards: expect heterogeneous payoffs; allocate curation/tuning to weak areas (e.g., metabolism/toxicity).
  - Keep prompting minimal and standardized (letter-only answers) to reduce spurious variance, especially with image items.
- **Scope and limitations:**
  - MCQs are a necessary but insufficient proxy for clinical deployment; real-world CDS requires free-text reasoning, uncertainty handling, and multimodal EHR integration.
  - SFT access was limited to GPT-4o; vendor constraints may alter relative rankings.

**Bottom line:** On a public, guideline-referenced hepatology benchmark, **injecting guideline knowledge (RAG), and especially combining it with SFT**, yields large, statistically significant accuracy gains over pooled physician performance, with the strongest configuration reaching **87.6%**. These results extend prior chapters from single-disease HCV to **broad hepatology**, offering a reproducible yardstick and concrete settings for safer, guideline-aligned LLM use.



## **Chapter 7**

**Reflect: Discussion and Conclusions: What Holds, What Does Not,  
and Where to Go Next**

**“All models are wrong, but some are useful”**

*– George E. P. Box*

## 7.1 Discussion

The present doctoral work set out to examine the readiness, safety, and optimization of LLMs for clinical decision support across the spectrum of digestive diseases. While AI has already transformed image recognition and structured-data analytics in medicine, the introduction of generative models capable of natural-language reasoning represents a paradigm shift—one that demands rigorous validation before integration into high-stakes clinical environments. The thesis followed a sequential logic: first establishing a baseline of LLM performance and methodological heterogeneity in published studies; then developing a structured, guideline-grounded retrieval and formatting framework for HCV management; next testing fine-tuning and expert validation strategies; subsequently introducing scalable verification and alignment methods through expert embeddings and reward models; and finally benchmarking multi-model performance across hepatology using the EASL Campus question bank as an openly accessible, guideline-referenced benchmark. Together, these chapters delineate a roadmap for transitioning LLMs from uncontrolled general-purpose chatbots into evidence-based, expert-aligned, and verifiable systems for clinical decision support in hepatology.

The work began with a systematic review of all published studies employing decoder-only LLMs in gastroenterology and hepatology. This effort, detailed in Chapter 2, revealed profound variability in accuracy—ranging from 6.4% to 91.4%—depending on model version, prompting strategies, and question type. Across the 18 included studies, ChatGPT-3.5 and 4 dominated the literature, yet even the latter failed to demonstrate consistency across domains. Most concerningly, safety-relevant behaviors such as over-recommendation of hospital admission and misclassification of follow-up intervals were recurrently observed. Equally striking was the lack of methodological transparency: few studies reported the exact model version, hyperparameters, or evaluator expertise, and nearly all used bespoke grading schemes that precluded meta-analysis. This heterogeneity underscored the absence of standardized evaluation metrics or reporting criteria, an omission that severely hampers reproducibility and comparability in this nascent field. The review therefore established two critical premises for the subsequent experimental work. First, that baseline LLMs—when used in zero-shot or general contexts—are not sufficiently accurate or safe for clinical deployment. Second, that reliable evaluation frameworks and benchmarks are essential prerequisites for the responsible advancement of LLMs in medicine. These conclusions mirror findings from broader evaluations of medical-chatbot performance, where LLMs have failed to meet passing thresholds on specialty examinations and displayed highly variable accuracy across digestive-disease tasks.<sup>114,115</sup>

Building on this foundation, Chapter 3 developed and validated a methodological framework for transforming clinical guidelines into machine-interpretable corpora suitable for RAG. Focusing on the EASL guideline for chronic HCV infection, this work demonstrated how structured reformatting and principled prompt engineering can dramatically improve the accuracy and factual faithfulness of LLM

outputs. The framework systematically converted heterogeneous guideline materials—including tables and flowcharts—into consistent, hierarchical text segments, coupled with explicit prompting templates. Through a five-stage ablation program, accuracy rose stepwise from 43% (baseline GPT-4 Turbo) to 99% when cleanly reformatted guideline text, structured retrieval, and optimized prompts were combined. Notably, performance gains were largest when non-textual elements such as tables were converted into plain-text lists, confirming that LLMs—even in their multimodal versions—struggle to interpret embedded or graphical data without explicit textual representation. These findings have both technical and epistemological implications: technically, they confirm that input structure and retrieval fidelity are decisive for performance; conceptually, they suggest that “model intelligence” alone cannot substitute for transparent knowledge encoding. Consistent with studies showing weak correspondence between lexical similarity metrics and clinical correctness, the chapter emphasized human expert grading for safety-clinical evaluation.<sup>116-119</sup> In practical terms, the chapter established a template for building guideline-grounded, clinically verifiable RAG frameworks that could underpin future decision-support systems in hepatology.

The next phase of the thesis, presented in Chapter 4, moved beyond single-disease evaluation to test generalization, optimization, and expert alignment. Here, SFT was introduced to complement RAG, using curated datasets derived from expert-authored questions and model responses graded by world leaders in HCV guideline development. The combination of RAG and SFT was assessed across multiple metrics—accuracy, clarity, and inter-rater reliability—under different chunking strategies and hyperparameter regimes. Results revealed several consistent trends. Paragraph-level chunking significantly outperformed sentence-based or fixed-length segmentation, preserving semantic coherence while maximizing retrieval relevance. Optimal temperature and top-p values were between 0.0–0.8 and 0.0–0.5 respectively, yielding outputs that were both accurate and consistent. These optimal configurations are concordant with recent multi-configuration studies in digestive medicine that reported similar temperature and top-p ranges and diminishing returns beyond ~Top10 retrieval.<sup>34,120</sup> Importantly, RAG and SFT each improved accuracy and clarity compared with the baseline, but their effects were context-dependent: RAG excelled in tasks requiring retrieval from heterogeneous text, whereas SFT enhanced stylistic and conceptual precision in open-ended explanations. This pattern matches prior work where grounding in guidelines reduced hallucinations and training bias,<sup>121</sup> baseline GPT-4 showed only mixed performance across GI subdomains, and RAG improved post-polypectomy follow-up recommendations from ~50.5% to ~79% versus foundational models.<sup>23</sup> Together, these findings validated the dual role of retrieval and fine-tuning as complementary mechanisms—one augmenting external factual grounding, the other strengthening internal linguistic alignment. They also contextualize our earlier HCV-focused pipeline that reached ~99% when guidelines were fully reformatted and tightly prompted.

Another major insight from this chapter concerned the evaluation process itself. When expert graders used binary metrics (correct/incorrect), agreement was only slight to fair, while a continuous 10-point Likert scale yielded good to excellent interclass correlation. This divergence highlights a central challenge in the assessment of generative AI for clinical use: accuracy is often multidimensional, encompassing factual correctness, interpretive reasoning, and communicative clarity. Binary labels risk oversimplifying these nuances, potentially underestimating clinical utility in borderline responses. The inclusion of multiple expert panels—both guideline authors and tertiary hepatologists—further demonstrated that domain familiarity influences grading, reinforcing the necessity of transparent reporting and consensus-based evaluation frameworks. Overall, Chapter 4 advanced the thesis from a methodological proof-of-concept toward an expert-aligned validation paradigm, showing that LLMs, when fine-tuned and evaluated within the framework of guideline expertise, can approach human-level performance on complex hepatology questions.

Chapter 5 extended these principles into a scalable safety framework—the Expert-of-Experts Verification and Alignment (EVAL) model. Recognizing that continuous expert supervision is infeasible for large-scale deployment, EVAL introduced a two-tier automated evaluation pipeline. At the model level, fine-tuned ColBERT embeddings were used to rank LLM configurations by similarity to golden-label answers authored by international guideline leaders. At the answer level, a lightweight reward model (OPT-350M) was trained to classify outputs as accurate or inaccurate across temperature settings. This dual system enables automated identification of both reliable model configurations and unsafe outputs, effectively bridging the gap between manual expert validation and autonomous LLM operation. EVAL is situated within a growing literature that has assessed LLMs mainly on diagnostic case vignettes<sup>122</sup> or multiple-choice testing<sup>123</sup>, and that has tended to use simple retrieval to inject guidelines<sup>124</sup>; our contribution advances this line by coupling expert-anchored embeddings with an explicit reward-model filter.

Empirically, EVAL demonstrated strong concordance with human judgment. Fine-tuned ColBERT achieved Spearman correlations up to 0.91 with expert rankings, while the reward model reproduced human labels with over 85% accuracy across multiple datasets. When applied to rejection sampling, the reward model improved factual accuracy by up to 98% in high-temperature regimes, mitigating the risk of hallucinations during exploratory generation. This temperature-sensitive rescue echoes reports that higher temperatures may aid reasoning while increasing hallucination risk.<sup>125,126</sup> Notably, different configurations excelled in different contexts—SFT-GPT-4o in expert-authored and MCQ datasets, RAG-GPT-o1 in real-world simulation questions—suggesting that LLM safety optimization must remain task-specific rather than model-specific. The framework also revealed the diminishing marginal gains of combining RAG and SFT in overlapping knowledge domains, a phenomenon akin to interference or redundancy between internal fine-tuned representations and external retrieval sources.

This observation is consistent with literature on catastrophic forgetting and interference during continued adaptation.<sup>127</sup> By systematically quantifying such interactions, EVAL represents one of the first practical pipelines for model-agnostic safety evaluation and tuning of LLMs in evidence-based medicine.

Finally, Chapter 6 applied the developed methodology to a publicly accessible benchmark, the EASL Campus MCQ repository, thereby extending validation across the full hepatology domain and enabling direct comparison with human physicians. This benchmark represents a unique resource: a guideline-referenced, peer-curated question set spanning major liver diseases. Using this dataset, the study compared baseline, RAG, SFT, and combined RAG-SFT configurations of three state-of-the-art proprietary models—GPT-4o, Claude-3-Opus, and Gemini-1.5-Pro—under standardized prompting and hyperparameter conditions. The results were striking. While baseline models achieved accuracies similar to or slightly above human test-takers (approximately 65–72%), RAG and SFT each produced statistically significant improvements, and their combination achieved up to 87.6% accuracy overall. Across subdomains, the highest accuracies were observed for liver tumors (95%) and immune-mediated diseases (100%), underscoring that well-structured guideline content strongly benefits retrieval-based systems. Conversely, performance remained modest in metabolic and toxic liver diseases, where guidelines themselves are less standardized—a reminder that AI performance is ultimately bounded by the quality and granularity of human knowledge sources. These findings both complement and extend earlier reports showing that ChatGPT-3 and GPT-4 failed to pass ACG self-assessments<sup>115</sup> and that, within UGIB-focused testing, SFT-GPT-4o could approach ~90% on expert questions whereas baseline models lagged physician averages (~75%).<sup>98</sup>

Crucially, this chapter marked the first demonstration that guideline-integrated LLMs can not only match but surpass average physician accuracy on a validated, society-endorsed benchmark. Yet it also highlighted the persistent gap between controlled test settings and clinical reality. Unlike MCQs, real clinical scenarios require synthesis of multimodal, incomplete, and longitudinal data, as well as probabilistic reasoning under uncertainty. Thus, while domain knowledge injection through RAG and SFT yields impressive results in structured settings, translation to practice will demand hybrid architectures that integrate reasoning modules, multimodal encoders, and continual expert oversight.

Taken together, these six chapters trace a clear developmental trajectory from descriptive assessment to actionable design principles for safe and effective clinical LLM deployment. The systematic review exposed foundational gaps—variability, opacity, and the absence of benchmarks—that mirror the early stages of evidence-based medicine itself, before standardized trial reporting transformed clinical research. The subsequent chapters progressively closed these gaps by introducing structured retrieval, fine-tuning, expert grading, automated verification, and benchmarking. Methodologically, the thesis

advances several innovations. First, it formalizes “LLM-friendly guideline engineering,” a process of re-expressing clinical knowledge into structured text optimized for machine retrieval without altering its clinical semantics. Second, it operationalizes expert-anchored evaluation through EVAL, creating a scalable path to align model output with domain expertise. Third, it establishes EASL Campus as a public benchmark for hepatology, enabling reproducible cross-model comparisons and tracking of progress over time. Conceptually, these contributions collectively bridge the divide between technical performance metrics and clinical epistemology, reframing LLM evaluation as an extension of guideline-based reasoning rather than mere language generation.

Despite these advances, several limitations warrant consideration. The experimental studies primarily focused on hepatology, which, although rich in structured guidelines, may not generalize to fields with less codified clinical pathways. The use of proprietary models limits reproducibility and external validation, as API-based implementations can change over time. Additionally, fine-tuning and RAG configurations were constrained by vendor access and cost, preventing exhaustive exploration of open-source models such as LLaMA or Mistral in equivalent settings. From an evaluation standpoint, expert grading—while rigorous—remains subject to cognitive bias and inter-rater variability, emphasizing the continued need for hybrid human-AI evaluation frameworks. Finally, the reliance on synthetic or test-bank questions, though methodologically necessary, does not fully capture the complexity of real-world diagnostic reasoning, where patient data are noisy, incomplete, and multimodal.

Nonetheless, the findings of this thesis carry profound implications for the future of AI-assisted hepatology. They suggest that, when appropriately structured, retrieved, and validated, LLMs can serve not merely as passive chatbots but as dynamic, guideline-conscious assistants capable of supporting clinician reasoning. In practice, such systems could be embedded within clinical decision-support dashboards, automatically contextualizing patient data against evidence-based recommendations, flagging inconsistencies, and providing just-in-time educational feedback. Moreover, the demonstrated capacity of RAG- and SFT-enhanced LLMs to outperform average human physicians on standardized tasks implies potential for real-world utility in training, triage, and quality assurance. Yet realizing this potential will depend on robust governance, transparency, and ongoing alignment with human expertise.

Future research should pursue several directions. First, prospective validation in simulated and real clinical workflows is essential to test safety, trustworthiness, and impact on decision-making. Second, guideline organizations such as EASL and AASLD should consider developing dual-purpose guideline versions: one optimized for human readability and one for machine interpretability, ensuring that future updates remain synchronized across modalities. Third, integrating reward models and uncertainty quantification into LLM outputs could facilitate risk-aware decision support, where the model explicitly communicates confidence and defers uncertain cases to human experts. Finally, interdisciplinary

collaborations between clinicians, computer scientists, and regulatory agencies will be necessary to define evaluation standards, liability frameworks, and data-protection mechanisms for clinical LLMs.

In summary, this thesis establishes a rigorous, stepwise framework for evaluating and optimizing LLMs in hepatology. Beginning from the recognition of unreliability in baseline models, it proceeds through the design of structured retrieval and alignment mechanisms, culminating in expert-anchored validation and benchmarking. The collective findings demonstrate that the safety and efficacy of LLMs in clinical medicine do not depend solely on model scale or training data volume, but on the deliberate design of knowledge interfaces, evaluation pipelines, and expert feedback loops. As medicine enters the era of generative AI, these principles offer a blueprint for developing trustworthy, interpretable, and clinically aligned decision-support systems that augment rather than replace the physician's judgment—anchoring technological innovation within the enduring ethos of evidence-based care.

## 7.2 Conclusions

This thesis demonstrates that the path toward safe, evidence-based integration of large language models in medicine depends not on scale or novelty alone, but on rigor, structure, and alignment. Through six sequential studies—progressing from systematic review to experimental validation and benchmarking—it has shown that the trustworthiness of generative AI in hepatology arises from disciplined engineering of both knowledge and evaluation. When clinical guidelines are reformatted into machine-readable structures, when retrieval and fine-tuning are judiciously combined, and when expert oversight is formalized through scalable verification frameworks such as EVAL, LLMs can approach and even surpass human-level accuracy on structured diagnostic and management tasks.

Yet, these achievements do not diminish the need for humility: models remain tools, not arbiters, whose strength lies in amplification of clinical reasoning rather than its replacement. The evidence assembled here offers a blueprint for responsible progress—one that treats transparency, reproducibility, and expert anchoring as non-negotiable design principles. As hepatology, and medicine more broadly, transitions into an AI-augmented future, the central insight of this work is that generative systems become safe only when they are made *interpretive rather than intuitive, guided rather than autonomous, and aligned with the ethical and epistemic foundations of clinical care.*



# Bibliography

1. Naveed, H. et al. A Comprehensive Overview of Large Language Models. (2023).
2. Nazi, Z. et al. Large language models in healthcare and medical domain: A review. (2023).
3. Wachter, R. et al. Will Generative Artificial Intelligence Deliver on Its Promise in Health Care? *JAMA* 331, 65 (2024).
4. Webster, P. Six ways large language models are changing healthcare. *Nat Med* 29, 2969–2971 (2023).
5. Soroush, A. et al. Generative Artificial Intelligence in Clinical Medicine and Impact on Gastroenterology. *Gastroenterology* <https://doi.org/10.1053/j.gastro.2025.03.038> (2025).
6. Vaswani, A. et al. Attention Is All You Need. (2023).
7. Bahdanau, D. et al. Neural Machine Translation by Jointly Learning to Align and Translate. (2016).
8. Sutskever, I. et al. Sequence to Sequence Learning with Neural Networks. (2014).
9. Lahat, A. et al. Evaluating the Utility of a Large Language Model in Answering Common Patients' Gastrointestinal Health-Related Questions: Are We There Yet? *Diagnostics* 13, 1950 (2023).
10. Tariq, R. et al. Evolving Landscape of Large Language Models: An Evaluation of ChatGPT and Bard in Answering Patient Queries on Colonoscopy. *Gastroenterology* <https://doi.org/10.1053/j.gastro.2023.08.033> (2023).
11. Lee, T.-C. et al. ChatGPT Answers Common Patient Questions About Colonoscopy. *Gastroenterology* 165, 509-511.e7 (2023).
12. Gorelik, Y. et al. Harnessing language models for streamlined postcolonoscopy patient management: a novel approach. *Gastrointest Endosc* 98, 639-641.e4 (2023).
13. Henson, J. B. et al. Evaluation of the Potential Utility of an Artificial Intelligence Chatbot in Gastroesophageal Reflux Disease Management. *American Journal of Gastroenterology* <https://doi.org/10.14309/ajg.0000000000002397> (2023).
14. Emile, S. H. et al. How appropriate are answers of online chat-based artificial intelligence (ChatGPT) to common questions on colon cancer? *Surgery* 174, 1273–1275 (2023).
15. Moazzam, Z. et al. Quality of ChatGPT Responses to Questions Related to Pancreatic Cancer and its Surgical Care. *Ann Surg Oncol* 30, 6284–6286 (2023).
16. Cankurtaran, R. E. et al. Reliability and Usefulness of ChatGPT for Inflammatory Bowel Diseases: An Analysis for Patients and Healthcare Professionals. *Cureus* <https://doi.org/10.7759/cureus.46736> (2023) doi:10.7759/cureus.46736.
17. Levartovsky, A. et al. Towards AI-Augmented Clinical Decision-Making: An Examination of ChatGPT's Utility in Acute Ulcerative Colitis Presentations. *American Journal of Gastroenterology* <https://doi.org/10.14309/ajg.0000000000002483> (2023)
18. Patil, N. S. et al. Using Artificial Intelligence Chatbots as a Radiologic Decision-Making Tool for Liver Imaging: Do ChatGPT and Bard Communicate Information Consistent With the ACR

Appropriateness Criteria? *Journal of the American College of Radiology* <https://doi.org/10.1016/j.jacr.2023.07.010> (2023).

19. Pugliese, N. et al. Accuracy, Reliability, and Comprehensiveness of ChatGPT-Generated Medical Responses for Patients With Nonalcoholic Fatty Liver Disease. *Clinical Gastroenterology and Hepatology* <https://doi.org/10.1016/j.cgh.2023.08.033> (2023).

20. Endo, Y. et al. Quality of ChatGPT Responses to Questions Related To Liver Transplantation. *Journal of Gastrointestinal Surgery* 27, 1716–1719 (2023).

21. Cao, J. J. et al. Accuracy of Information Provided by ChatGPT Regarding Liver Cancer Surveillance and Diagnosis. *American Journal of Roentgenology* 221, 556–559 (2023).

22. Yeo, Y. H. et al. Assessing the performance of ChatGPT in answering questions regarding cirrhosis and hepatocellular carcinoma. *Clin Mol Hepatol* 29, 721–732 (2023).

23. Lim, D. Y. Z. et al. ChatGPT on guidelines: Providing contextual knowledge to GPT allows it to provide advice on appropriate colonoscopy intervals. *J Gastroenterol Hepatol* 39, 81–106 (2024).

24. Mukherjee, S. et al. Assessing ChatGPT's Ability to Reply to Queries Regarding Colon Cancer Screening Based on Multisociety Guidelines. *Gastro Hep Advances* 2, 1040–1043 (2023).

25. Kerbage, A. et al. Accuracy of ChatGPT in Common Gastrointestinal Diseases: Impact for Patients and Providers. *Clinical Gastroenterology and Hepatology* <https://doi.org/10.1016/j.cgh.2023.11.008> (2023).

26. Atarere, J. et al. Applicability of Online Chat-Based Artificial Intelligence Models to Colorectal Cancer Screening. *Dig Dis Sci* <https://doi.org/10.1007/s10620-024-08274-3> (2024).

27. Cacciamani, G. E. et al. ChatGPT: standard reporting guidelines for responsible use. *Nature* 618, 238–238 (2023).

28. Suchman, K. et al. Chat Generative Pretrained Transformer Fails the Multiple-Choice American College of Gastroenterology Self-Assessment Test. *American Journal of Gastroenterology* 118, 2280–2282 (2023).

29. Kung, T. H. et al. Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models. *PLOS Digital Health* 2, e0000198 (2023).

30. Lee, P. et al. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine. *New England Journal of Medicine* 388, 1233–1239 (2023).

31. Eriksen, A. V. et al. Use of GPT-4 to Diagnose Complex Clinical Cases. *NEJM AI* 1, (2023).

32. Giuffrè, M. et al. Evaluating ChatGPT in Medical Contexts: The Imperative to Guard Against Hallucinations and Partial Accuracies. *Clinical Gastroenterology and Hepatology* <https://doi.org/10.1016/j.cgh.2023.09.035> (2023).

33. Giuffrè, M. et al. Scrutinizing ChatGPT Applications in Gastroenterology: A Call for Methodological Rigor to Define Accuracy and Preserve Privacy. *Clinical Gastroenterology and Hepatology* <https://doi.org/10.1016/j.cgh.2024.01.024> (2024).

34. Lewis, P. et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. (2020).

35. Chen, B. et al. Unleashing the potential of prompt engineering in Large Language Models: a comprehensive review. (2023).
36. Sivarajkumar, S. et al. An Empirical Evaluation of Prompting Strategies for Large Language Models in Zero-Shot Clinical Natural Language Processing. (2023).
37. Laine, L. et al. ACG Clinical Guideline: Upper Gastrointestinal and Ulcer Bleeding. *American Journal of Gastroenterology* 116, 899–917 (2021).
38. Chang, P. W. et al. ChatGPT4 Outperforms Endoscopists for Determination of Postcolonoscopy Rescreening and Surveillance Recommendations. *Clinical Gastroenterology and Hepatology* 22, 1917-1925.e17 (2024).
39. Simhi, A. et al. Interpreting Embedding Spaces by Conceptualization. (2022).
40. Ratner, N. et al. Parallel Context Windows for Large Language Models. (2022).
41. Chang, Y. et al. BoookScore: A systematic exploration of book-length summarization in the era of LLMs. (2023).
42. 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).
43. Raffel, C. et al. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. (2019).
44. Xue, L. et al. mT5: A massively multilingual pre-trained text-to-text transformer. (2020).
45. Askell, A. et al. A General Language Assistant as a Laboratory for Alignment. (2021).
46. Han, T. et al. MedAlpaca -- An Open-Source Collection of Medical Conversational AI Models and Training Data. (2023).
47. Van Veen, D. et al. Adapted large language models can outperform medical experts in clinical text summarization. *Nat Med* <https://doi.org/10.1038/s41591-024-02855-5> (2024).
48. Wolf, Y. et al. Fundamental Limitations of Alignment in Large Language Models. (2023).
49. Wang, Y. et al. Aligning Large Language Models with Human: A Survey. (2023).
50. Ziegler, D. M. et al. Fine-Tuning Language Models from Human Preferences. (2019).
51. Xie, Q. et al. Me LLaMA: Foundation Large Language Models for Medical Applications. (2024).
52. Li, J. Security Implications of AI Chatbots in Health Care. *J Med Internet Res* 25, e47551 (2023).
53. Dwivedi, Y. K. et al. Opinion Paper: "So what if ChatGPT wrote it?" Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy. *Int J Inf Manage* 71, 102642 (2023).
54. Page, M. J. et al. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *BMJ* n71 (2021) doi:10.1136/bmj.n71.
55. Campbell, M. et al. Synthesis without meta-analysis (SWiM) in systematic reviews: reporting guideline. *BMJ* l6890 (2020) doi:10.1136/bmj.l6890.
56. McGowan, J. et al. PRESS Peer Review of Electronic Search Strategies: 2015 Guideline Statement. *J Clin Epidemiol* 75, 40–46 (2016).

57. Haddaway, N. R. et al. Citationchaser: A tool for transparent and efficient forward and backward citation chasing in systematic searching. *Res Synth Methods* 13, 533–545 (2022).
58. Yale University Harvey Cushing/John Hay Whitney Medical Library. Reference Deduplicator [Internet]. 2021. <https://library.medicine.yale.edu/referencededuplicator>.
59. Munn, Z. et al. Methodological guidance for systematic reviews of observational epidemiological studies reporting prevalence and cumulative incidence data. *Int J Evid Based Healthc* 13, 147–153 (2015).
60. Embeddings. Accessed March 10, 2024. <https://platform.openai.com/docs/guides/embeddings>.
61. Cai, T. T. & Ma, R. Theoretical Foundations of t-SNE for Visualizing High-Dimensional Clustered Data. arXiv: 2105.07536 (2021).
62. Van Der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9, 2579–2605 (2008).
63. Pawlotsky, J.-M. et al. EASL recommendations on treatment of hepatitis C: Final update of the series☆. *J Hepatol* 73, 1170–1218 (2020).
64. Chung, R. T. et al. Hepatitis C guidance: AASLD-IDSAs recommendations for testing, managing, and treating adults infected with hepatitis C virus. *Hepatology* 62, 932–954 (2015).
65. Zhang, Y. et al. Siren's Song in the AI Ocean: A Survey on Hallucination in Large Language Models. (2025).
66. Ganesan, K. ROUGE 2.0: Updated and Improved Measures for Evaluation of Summarization Tasks. (2018).
67. Papineni, K. et al. BLEU. in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02* 311 (Association for Computational Linguistics, Morristown, NJ, USA, 2001). doi:10.3115/1073083.1073135.
68. Lavie, A. & Denkowski, M. J. The Meteor metric for automatic evaluation of machine translation. *Machine Translation* 23, 105–115 (2009).
69. Giuffrè, M. et al. Systematic review: The use of large language models as medical chatbots in digestive diseases. *Aliment Pharmacol Ther* <https://doi.org/10.1111/apt.18058> (2024) doi:10.1111/apt.18058.
70. Preston, C. C. et al. Optimal number of response categories in rating scales: reliability, validity, discriminating power, and respondent preferences. *Acta Psychol (Amst)* 104, 1–15 (2000).
71. Alwin, D. F. Feeling Thermometers Versus 7-Point Scales. *Sociol Methods Res* 25, 318–340 (1997).
72. Zhang, B. et al. A Computational Approach to Interpreting the Embedding Space of Dimension Reduction (2024). doi:10.1101/2024.06.23.600292.
73. 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).
74. Wang, S. et al. LLM can Achieve Self-Regulation via Hyperparameter Aware Generation. (2024).

75. Atil, B. et al. LLM Stability: A detailed analysis with some surprises. (2024).
76. Radford, A. et al. Language Models Are Unsupervised Multitask Learners. <https://github.com/codelucas/newspaper>.
77. Giuffrè, M. Distinguishing Retrieval Augmented Generation From Prompt Engineering: Implications for Reproducibility in Large Language Model Research and Applications. *American Journal of Gastroenterology* <https://doi.org/10.14309/ajg.0000000000003439> (2025).
78. Yepes, A. J. et al. Financial Report Chunking for Effective Retrieval Augmented Generation. (2024).
79. OpenAI - Accessed May 03 2024. Fine-Tuning. <https://platform.openai.com/docs/guides/fine-tuning>.
80. Ouyang, L. et al. Training language models to follow instructions with human feedback. (2022).
81. Bhattacharya, D. et al. Hepatitis C Guidance 2023 Update: American Association for the Study of Liver Diseases– Infectious Diseases Society of America Recommendations for Testing, Managing, and Treating Hepatitis C Virus Infection. *Clinical Infectious Diseases* <https://doi.org/10.1093/cid/ciad319> (2023) doi:10.1093/cid/ciad319.
82. Shrout, P. E. et al. Intraclass correlations: Uses in assessing rater reliability. *Psychol Bull* 86, 420–428 (1979).
83. Cicchetti, D. V. Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. *Psychol Assess* 6, 284–290 (1994).
84. Fleiss, J. L. Measuring nominal scale agreement among many raters. *Psychol Bull* 76, 378–382 (1971).
85. Zheng, N. S. et al. Trends in characteristics, management, and outcomes of patients presenting with gastrointestinal bleeding to emergency departments in the United States from 2006 to 2019. *Aliment Pharmacol Ther* 56, 1543–1555 (2022).
86. 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. *American Journal of Gastroenterology* 108, 1449–1457 (2013).
87. 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).
88. Lu, Y. et al. 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).
89. Liang, P. S. & Saltzman, J. R. A National Survey on the Initial Management of Upper Gastrointestinal Bleeding. *J Clin Gastroenterol* 48, e93–e98 (2014).
90. 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).

91. 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 Periendoscopic Period. *American Journal of Gastroenterology* 117, 542–558 (2022).
92. de Franchis, R. et al. Baveno VII – Renewing consensus in portal hypertension. *J Hepatol* 76, 959–974 (2022).
93. 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).
94. Sung, J. J. et al. Asia-Pacific working group consensus on non-variceal upper gastrointestinal bleeding: an update 2018. *Gut* 67, 1757–1768 (2018).
95. Giuffrè, M. et al. Optimizing large language models in digestive disease: strategies and challenges to improve clinical outcomes. *Liver International* <https://doi.org/10.1111/liv.15974> (2024) doi:10.1111/liv.15974.
96. Dettmers, T. et al. QLoRA: Efficient Finetuning of Quantized LLMs. (2023).
97. Hu, E. J. et al. LoRA: Low-Rank Adaptation of Large Language Models. (2021).
98. 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 Recommended Management Of Patients With Upper Gastrointestinal Bleeding. *Gastroenterology* 166, S-889-S-890 (2024).
99. Rajashekar, N. C. et al. Human-Algorithmic Interaction Using a Large Language Model-Augmented Artificial Intelligence Clinical Decision Support System. in *Proceedings of the CHI Conference on Human Factors in Computing Systems* 1–20 (ACM, New York, NY, USA, 2024). doi:10.1145/3613904.3642024.
100. Khattab, O. et al. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. (2020).
101. Devlin, J. et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. (2018).
102. Reimers, N. et al. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. (2019).
103. Sparck Jones, K. A Statistical Interpretation Of Term Specificity And Its Application In Retrieval. *Journal of Documentation* 28, 11–21 (1972).
104. Feng, F. et al. Language-agnostic BERT Sentence Embedding (2022).
105. Pedregosa, F. et al. Scikit-learn: Machine Learning in Python. (2012).
106. Dhuliawala, S. et al. Chain-of-Verification Reduces Hallucination in Large Language Models. (2023).
107. Nath, S. et al. Leveraging Domain Knowledge for Efficient Reward Modelling in RLHF: A Case-Study in E-Commerce Opinion Summarization. (2024).

108. Wang, Z. et al. HelpSteer2: Open-source dataset for training top-performing reward models. (2024).
109. Zhang, S. et al. Instruction Tuning for Large Language Models: A Survey (2024).
110. Rafailov, R. et al. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. (2023).
111. Bai, Y. et al. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback. (2022).
112. Stiennon, N. et al. Learning to summarize from human feedback. (2020).
113. Liu, M. et al. Performance of ChatGPT Across Different Versions in Medical Licensing Examinations Worldwide: Systematic Review and Meta-Analysis. *J Med Internet Res* 26, e60807 (2024).
114. Weng, T. L. et al. ChatGPT failed Taiwan's Family Medicine Board Exam. *J Chin Med Assoc* 86, 762–766 (2023).
115. Chen, A. et al. Evaluating Question Answering Evaluation. in *Proceedings of the 2nd Workshop on Machine Reading for Question Answering* 119–124 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2019). doi:10.18653/v1/D19-5817.
116. Tang, L. et al. Evaluating large language models on medical evidence summarization. *NPJ Digit Med* 6, 158 (2023).
117. Fabbri, A. R. et al. SummEval: Re-evaluating Summarization Evaluation. *Trans Assoc Comput Linguist* 9, 391–409 (2021).
118. Blagec, K. et al. A global analysis of metrics used for measuring performance in natural language processing. in *Proceedings of NLP Power! The First Workshop on Efficient Benchmarking in NLP* 52–63 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2022). doi:10.18653/v1/2022.nlppower-1.6.
119. Izacard, G. & Grave, E. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. (2020).
120. Ayoub, N. F. et al. Inherent Bias in Large Language Models: A Random Sampling Analysis. *Mayo Clinic Proceedings: Digital Health* 2, 186–191 (2024).
121. Knoedler, L. et al. In-depth analysis of ChatGPT's performance based on specific signaling words and phrases in the question stem of 2377 USMLE step 1 style questions. *Sci Rep* 14, 13553 (2024).
122. Bicknell, B. T. et al. ChatGPT-4 Omni Performance in USMLE Disciplines and Clinical Skills: Comparative Analysis. *JMIR Med Educ* 10, e63430 (2024) doi:10.2196/63430.
123. Unlu, O. et al. Retrieval-Augmented Generation–Enabled GPT-4 for Clinical Trial Screening. *NEJM AI* 1, (2024).
124. Renze, M. & Guven, E. The Effect of Sampling Temperature on Problem Solving in Large Language Models. <https://doi.org/10.18653/v1/2024.findings-emnlp.432> (2024).

125. Windisch, P. et al. The Impact of Temperature on Extracting Information From Clinical Trial Publications Using Large Language Models. *Cureus* <https://doi.org/10.7759/cureus.75748> (2024) doi:10.7759/cureus.75748.
126. Luo, Y. et al. An Empirical Study of Catastrophic Forgetting in Large Language Models During Continual Fine-tuning. (2023).
127. Giuffrè, M. et al. Expert of Experts Verification and Alignment (EVAL) Framework for Large Language Models Safety in Gastroenterology. *NPJ Digit Med* 8, 242 (2025).