

# LLM-Based Multi-Agent System with Retrieval-Augmented Generation for Medical Care Planning Generation in Sickle Cell Disease

Luana Bringel Leite<sup>1</sup>, David Eduardo Pereira<sup>1</sup>, Eyshila Buriti de Araujo Azevedo<sup>1</sup>,  
Leonardo Mota Meira Filho<sup>1</sup>, Eliane Cristina Araújo<sup>1</sup>, Claudio E. C. Campelo<sup>1</sup>,  
Taciana R. O. C. Marques<sup>2</sup>, Leticia B. de Almeida<sup>2</sup>, Herman Martins Gomes<sup>1</sup>

<sup>1</sup>Federal University of Campina Grande (UFCG) - Academic Unity of Systems and Computing

<sup>2</sup>Federal University of Campina Grande (UFCG) - Center for Biological and Health Sciences

Correspondence: [luana.leite@ccc.ufcg.edu.br](mailto:luana.leite@ccc.ufcg.edu.br)

## Abstract

Ensuring safety in clinical applications of large language models (LLMs) remains an unresolved challenge, particularly for high-risk and underrepresented conditions such as Sickle Cell Disease (SCD). Consequently, these models may exhibit limited reliability for SCD, including hallucinations and clinical non-adherence. This paper proposes an LLM-based Multi-Agent System (MAS) enhanced by Retrieval-Augmented Generation (RAG) to automate the generation of medical care plans for SCD. The MAS decomposes clinical reasoning into specialized agents responsible for diagnosis, investigation, and treatment planning. Retrieval is framed not as a performance optimization, but as a safety control mechanism. Three RAG strategies, namely LLM-Guided Tree Retrieval, Metadata-Filtered Retrieval, and Semantic Similarity Retrieval, are evaluated alongside a baseline. Our experiments considered LLM-as-a-Judge evaluations and independent assessments by physicians. The results demonstrate high clinical quality, with safety scores exceeding 4 on a 5-point scale. While average performance was similar between RAG and baseline conditions, the Tree Retrieval strategy reduced the frequency of clinically unsafe outputs compared to conventional Semantic Retrieval, indicating fewer clinically unsafe outputs. These findings provide evidence that average performance is insufficient to evaluate clinical AI systems, particularly in high-risk scenarios where retrieval serves as a safety control layer.

## 1 Introduction

Creating a medical care plan for a patient is a critical task for healthcare professionals. A care plan has a direct impact on a patient's life quality (Abdeldafie and Alaajmi, 2022), as it guides clinical procedures, medication management, diagnostic examinations, monitoring protocols, among other interventions. This process becomes more challenging when designing care plans for patients with

rare diseases or conditions, a scenario that physicians rarely encounter during medical education or routine practice. (Walkowiak and Domaradzki, 2021).

Sickle Cell Disease (SCD) is a rare genetic disorder that, according to existing studies, will affect around 400,000 individuals worldwide until 2050, including a significant number in Brazil (estimate of 30,000) (Kato et al., 2018). Given its prevalence in specific populations, SCD is often underrepresented in medical education and insufficiently understood by physicians, leading to inadequate management in clinical practice (Druye et al., 2024). This lack of familiarity is further influenced by broader social, cultural, and historical factors, since SCD disproportionately affects individuals of African descent due to genetic inheritance.

Aligned with this social and cultural marginalization, significant knowledge gaps persist among physicians regarding SCD (Druye et al., 2024). Structural inequalities have contributed to its underrepresentation in medical education, research and clinical practice (Reich et al., 2022). This issue extends to technological development within the field of computer science and healthcare informatics. Studies indicate that SCD, along with other rare and African diseases, is underrepresented in healthcare evaluation datasets (Mutisya et al., 2025). As a result, AI models may exhibit limited performance and reduced reliability when addressing these conditions, reinforcing gaps in access to specialized care.

This research proposes an LLM-based Multi-Agent System (MAS) enhanced by Retrieval-Augmented Generation (RAG) for automated medical care plan generation in SCD as a module that is part of *HemaChat*<sup>1</sup> software, a multi-agent clinical reasoning and decision support system designed to expand access to safe medical guidance for SCD

<sup>1</sup><https://bit.ly/hemachat>

patients. The proposed MAS explicitly decomposes clinical reasoning into sequential specialized agents. By integrating structured retrieval at each reasoning stage, the system constrains generation within validated clinical protocols, improving reliability, interpretability, and clinical safety. Previous studies indicate that RAG can significantly improve the quality and reliability of LLM-based applications in healthcare (Amugongo et al., 2025). In safety-critical domains, retrieval architecture should be understood as a safety control layer rather than a performance optimization, as it directly constrains generation and reduces the likelihood of harmful outputs.

This research evaluates different types of RAG techniques in the context of medical care plan generation for SCD on patient-reported symptoms. Since existing models can demonstrate limited capabilities when dealing with rare diseases, the application of RAG can enhance reliability and performance in medical LLM-based systems. Therefore, the present study evaluates three distinct RAG techniques (Semantic Similarity Retrieval, Metadata-Filtered Retrieval, and LLM-Guided Tree Retrieval), incorporating both **LLM-as-a-judge** (Li et al., 2025) and **human** evaluations conducted by physicians.

This study offers three key contributions: (1) a structured multi-agent architecture aligned with clinical reasoning workflows; (2) an evaluation of retrieval mechanisms for safety-critical clinical applications; and (3) empirical evidence supporting retrieval as a safety control layer in generative AI systems, demonstrating its role in reducing unsafe outputs in high-risk clinical scenarios.

The remainder of this paper is structured as follows. Section 2 reviews related work on LLMs and Retrieval-Augmented Generation in healthcare. Section 3 presents the proposed multi-agent architecture, the RAG strategies, and the clinical evaluation protocol. Section 4 reports the experimental results. Finally, Section 5 discusses the implications, limitations, and future research directions.

## 2 Related Work

The use of LLMs in healthcare is a rapidly expanding research field (Wang et al., 2024b). AI systems have demonstrated the ability to support a wide range of healthcare-related tasks. Existing review studies show that LLM applications in healthcare are broad, encompassing the summa-

rization of complex clinical information, medical knowledge retrieval to support question answering and examinations, improved public access to medical information, predictive tasks such as diagnosis, treatment support, and drug interaction analysis, as well as administrative activities, including clinical documentation and public health data collection (Wang et al., 2024b).

Despite these broad applications, the use of LLMs in healthcare raises several ethical concerns that must be addressed to prevent harm (Wang et al., 2023). These concerns span legal, humanistic, algorithmic, and informational dimensions, including unclear liability in cases of patient harm, risks to patient privacy, potential disruption of the physician–patient relationship, erosion of trust due to over-reliance on AI, and challenges related to transparency, bias, and explainability (Wang et al., 2023).

Recent advances in LLMs demonstrate improved reasoning capabilities, reduced latency, and multimodal functionality, which help mitigate some challenges in healthcare applications (Neha et al., 2025). Notably, many LLMs have been fine-tuned on biomedical corpora to enhance domain-specific comprehension (Neha et al., 2025), including recent ChatGPT variants specifically designed for healthcare-related queries<sup>2</sup>. However, these models cannot continuously incorporate evolving clinical knowledge, which limits their adaptability in dynamic healthcare environments and are susceptible to hallucinations.

To address these limitations, RAG techniques have emerged as a promising approach to improve the reliability of LLMs. RAG helps mitigate hallucinations and reduces over-reliance on static model training data (Arslan et al., 2024). Nevertheless, RAG is not a silver bullet. Studies indicate that hallucinations can still occur and that factual inconsistencies remain a persistent issue even in RAG-based systems applied to healthcare scenarios (Amugongo et al., 2025).

The study proposed by Neha et al. (2025) emphasizes the use of RAG in healthcare domains such as diagnostic assistance, electronic health record and discharge note summarization, medical question answering, patient education and conversational agents, clinical trial matching, and biomedical literature synthesis. Beyond these areas, research

<sup>2</sup><https://openai.com/pt-BR/index/introducing-chatgpt-health/>

has explored the application of RAG in more specialized healthcare domains. These include mental health-related solutions (Kermani et al., 2025), patient simulation for educational purposes (Yu et al., 2025), health problem identification in home healthcare settings (Zhang et al., 2025), inclusive urban public healthcare services (Sun et al., 2025), among other context-specific applications.

It is important to highlight that RAG can be a valuable tool for mitigating hallucinations and other LLM-related issues in healthcare scenarios. However, studies also reveal significant limitations. Most research focuses on English and Chinese, leaving many other languages underrepresented (Amugongo et al., 2025). Moreover, bias can persist or even be reproduced through RAG pipelines, as biased source data can propagate biased outputs, potentially leading to harmful or misleading information. Another limitation is that current evaluation metrics are often insufficient to assess RAG performance in healthcare contexts, as they may not adequately capture clinical relevance or safety considerations (Neha et al., 2025).

For these reasons, this research investigates the effectiveness of LLM-based MAS for medical care plan generation in SCD, enhanced by RAG techniques. It combines quantitative metrics with assessments generated by specialized human and LLM evaluations. Furthermore, the focus on SCD provides an important contribution, given the limitations of LLM in rare disease contexts that are often underrepresented in datasets and LLM models (Mutisya et al., 2025). This work also takes into account the Brazilian Portuguese language setting, addressing another critical gap in current healthcare-focused LLM research, which is predominantly centered on the English language.

It is important to note that research on rare diseases, such as SCD, remains limited and often receives insufficient attention, mainly due to the relatively small number of affected individuals (Visibelli et al., 2023). When surveying the literature on SCD and AI, only a small number of relevant studies can be identified. These include research on the use of LLMs in ambulatory devices for home health diagnostics, a case study focused on sickle cell anemia management (Ogundare and Sofolahan, 2023). Additionally, a question-answering study addresses rare diseases more broadly rather than exclusively focusing on SCD (Wang et al., 2024a).

Furthermore, when examining research on medical care plan generation, a similar scarcity of stud-

ies is observed. One notable example is MED-Plan (Hsu et al., 2025), an approach for medical care plan generation that leverages LLMs and RAG within the Subjective, Objective, Assessment, Plan (SOAP) framework, however, it does not directly address the limitations of the SCD context. To the best of our knowledge, there is currently no research investigating the use of LLM-based MAS enhanced by RAG in the context of the generation of medical care plans specifically tailored to SCD.

### 3 Methodology

This study is a prospective, blinded clinical validation to evaluate the adequacy, completeness, and safety of medical care plans generated by an LLM-based MAS enhanced with RAG within a broader multi-agent clinical reasoning and decision support system. All clinical cases, system outputs, and evaluation procedures were predefined prior to assessment, and physician evaluators were blinded to the retrieval strategy, ensuring unbiased assessment.

The system itself is intentionally engineered to replicate the sequential reasoning process employed by physicians in emergency care settings as a decision support tool (Croskerry, 2009). In real clinical workflows, physicians do not generate diagnoses, investigations, and treatments simultaneously; rather, they follow a *structured reasoning sequence* in which initial clinical observations inform diagnostic hypotheses, which in turn guide selection of confirmatory investigations and ultimately determine therapeutic interventions. This process is iterative and uncertainty-aware, and is mirrored by the proposed MAS system (see Section 3.3 for details).

The experimental protocol compared three RAG strategies: LLM-Guided Tree Retrieval (proposed method), Metadata-Filtered Retrieval, and Semantic Similarity Retrieval, each evaluated against a baseline without retrieval. Each generated medical care plan was independently evaluated by three blinded physicians and, in parallel, by an automated LLM-as-a-Judge framework (Gu et al., 2024). This dual evaluation design allowed for a direct comparison between human clinical judgement and automated evaluation methods.

#### 3.1 Medical Care Plan Data

To ensure validity and faithful representation of real-world clinical complexity, 10 clinical case vi-

gnettes were prospectively developed by a senior pediatric and SCD specialist. These cases were specifically designed to represent the spectrum of acute SCD emergencies most frequently encountered in emergency departments, including vaso-occlusive crisis, acute chest syndrome, splenic sequestration, ischemic priapism, and acute neurological complications such as ischemic stroke and seizure presentations, which together represent the most common and clinically significant acute manifestations of SCD (Rees et al., 2010; Piel et al., 2017).

These scenarios represent high-risk emergency conditions that require timely and appropriate management, making them suitable for the evaluation of critical safety systems. Although the number of clinical cases included in this study is limited, this study prioritizes depth of clinical validation over scale. Each vignette represents a high-risk emergency scenario and was carefully designed by a specialist to reflect real-world complexity and clinical decision-making requirements and cover most recurrent acute SCD complications.

In addition, each vignette was constructed using a realistic emergency department triage structure. A specialist physician generated a comprehensive gold-standard medical care plan representing the clinically optimal management approach for each patient. The standardized clinical structure comprised patient identification and clinical context, clinical history and presenting complaint, physical examination findings, differential diagnosis, diagnostic investigations, therapeutic management, and follow-up monitoring. A detailed description and a complete example of the vignettes are provided in Figure 4 in the Appendix A.6 for illustration.

These reference medical care plans were constructed using current institutional protocols and established clinical guidelines and served as the *clinical reference benchmark* against which AI-generated medical care plans were evaluated. Importantly, the Gold Standard was not directly provided to the AI system during generation, ensuring that evaluation reflected true generalization rather than memorization.

### 3.2 RAG Knowledge Base

To ensure clinical safety and prevent the generation of recommendations based on unreliable or unverified sources, the knowledge base of the system was constructed exclusively from validated institutional clinical protocols and professional society

guidelines.

The knowledge corpus consists of:

1. **The Pediatric Emergency Care Protocol** of the *Hospital de Clínicas de Porto Alegre* (HCPA) (Hospital de Clínicas de Porto Alegre, 2023), and
2. **Clinical management guidelines** from the *Sociedade de Pediatria do Estado do Rio de Janeiro* (SOPERJ) (Sociedade de Pediatria do Estado do Rio de Janeiro, 2023).

These sources are authoritative clinical references used in pediatric emergency care in Brazil. Restricting the knowledge base to curated clinical protocols is a critical safety design decision, as open-web retrieval can introduce outdated, contradictory, or non-validated medical information (Institute of Medicine, 2011).

Documents were segmented into semantically coherent text fragments (“*chunks*”) and indexed within a ChromaDB vector database<sup>3</sup> using all-MiniLM-L6-v2 embeddings<sup>4</sup>. This embedding model was selected due to its performance in semantic retrieval tasks while maintaining computational efficiency.

The retrieval hyperparameter was set to top-k = 5, meaning that the five most relevant knowledge fragments were retrieved for each query. This choice is consistent with prior RAG research, which shows that retrieval configuration and document relevance influence reasoning reliability and hallucinations (Lewis et al., 2020; Yan et al., 2024). Retrieving too few documents may lead to incomplete clinical context, whereas retrieving too many may introduce irrelevant information that degrades generation quality.

### 3.3 Implementation Details

All agents and experimental conditions were powered by the GPT-4.1 model<sup>5</sup> via the OpenAI API, and agent communication was implemented using the Python-based LangChain<sup>6</sup> framework. For safety-critical reasoning tasks, generation was performed with the temperature set to 0, promoting

<sup>3</sup><https://www.trychroma.com/>

<sup>4</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>5</sup><https://developers.openai.com/api/docs/models/gpt-4.1>

<sup>6</sup><https://www.langchain.com/>

deterministic and stable outputs; this configuration reduces stochastic variability and improves reproducibility in clinical decision-making scenarios. Furthermore, all retrieval strategies and the baseline used the same underlying LLM, ensuring that observed differences are attributable to retrieval mechanisms rather than to model variation.

### 3.4 Retrieval Strategies

This study investigates three different retrieval approaches, each with its own data representation and retrieval mechanism, as described below.

**Semantic Similarity Retrieval:** This standard dense approach retrieves protocol fragments using dense vector similarity. Queries and document chunks are embedded using the all-MiniLM-L6-v2, and the top-k fragments are selected using Maximal Marginal Relevance (MMR).

**Metadata-Filtered Retrieval:** This strategy extends semantic retrieval by applying metadata filters before similarity search. Protocol fragments are first restricted based on attributes such as disease and protocol category, and similarity search is performed within this subset.

**LLM-Guided Tree Retrieval:** This strategy implements a structured LLM-guided retrieval mechanism over a hierarchical representation of clinical knowledge. The knowledge base is organized as a tree structure, where each node represents a clinically meaningful unit (e.g., symptoms, diagnostic categories, or treatment protocols). Each node contains a unique identifier, a title summarizing its clinical meaning, a short description, and the associated clinical content.

Given a clinical vignette, the LLM analyzes the patient symptoms and selects the most relevant nodes from the tree. This selection is performed by providing the model with a structured representation of the tree (excluding full clinical text) and prompting it to identify relevant node identifiers. Hence, the retrieval process consists of two stages: (1) node selection, in which the LLM identifies relevant nodes, and (2) content extraction, in which the full clinical content associated with these nodes is retrieved and concatenated to form the final context.

Unlike conventional semantic retrieval, which operates on flat document representations, the tree approach constrains retrieval to clinically coherent pathways, reducing the likelihood of retrieving

contextually irrelevant but lexically similar information. Additionally, the model produces an intermediate reasoning trace during node selection, enabling interpretability of retrieval decisions.

### 3.5 LLM-Based Multi-Agent System Architecture

The proposed system utilizes a sequential LLM-Based Multi-Agent architecture in which each agent performs a specialized stage of the clinical reasoning process, thereby improving interpretability and enabling fine-grained safety analysis. Rather than generating medical care plans in a single step, the system progresses through a structured reasoning pipeline, as illustrated in Figure 1, which presents the complete multi-agent clinical reasoning pipeline, listed in the following paragraphs:

**Diagnostic Hypothesis Agent:** This agent initiates the clinical reasoning process by analyzing the clinical vignette and retrieving protocol-grounded knowledge to construct a structured differential diagnosis. Each hypothesis is explicitly categorized as *Most probable*, *Less probable*, *To be ruled out*, or *Diagnosis of exclusion*, enabling formal representation of clinical uncertainty and systematic prioritization of high-risk conditions.

**Diagnostic Investigation Agent:** Building upon the diagnostic hypotheses, this agent retrieves protocol-aligned recommendations for diagnostic investigations. By conditioning test selection on explicit reasoning outputs, the system ensures clinical coherence and reduces the risk of incomplete or unjustified evaluations.

**Treatment Agent:** This agent generates therapeutic recommendations grounded in validated clinical protocols, taking into account the prioritized diagnostic hypotheses and clinical severity. This structured grounding constrains generation within clinically accepted standards and improves treatment safety.

**Medical Care Plan Generation Agent:** The final agent integrates all intermediate outputs into a unified and structured medical care plan. This modular synthesis preserves traceability across reasoning stages and produces coherent clinically actionable documentation aligned with real-world clinical workflows.

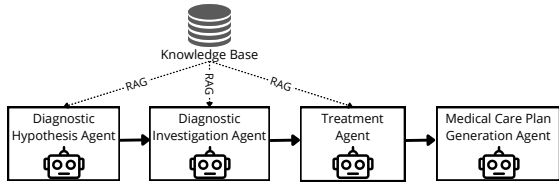


Figure 1: Overview of the proposed LLM-based Multi-Agent Retrieval-Augmented Generation architecture for medical care plan generation.

### 3.6 Human Clinical Evaluation Protocol

The generated medical care plans were evaluated through a prospective, blinded physician assessment protocol. A total of 22 physicians participated as evaluators, including both board-certified specialists and physicians in residency training. This mixed-expertise composition reflects clinicians who may interact with AI-generated medical care plans in real healthcare settings, including those with limited access to specialized SCD expertise.

The experimental design comprised 40 distinct AI-generated medical care plan documents, corresponding to 10 clinical cases processed under each of the three retrieval strategies and the baseline. Each medical care plan was independently evaluated by exactly 3 different physicians, resulting in a total of 120 blinded physician evaluations as summarized in Figure 2 (Appendix A.1).

To preserve blinding and prevent expectation bias, physicians were not informed of the generation method, RAG strategy, or experimental condition associated with any medical care plans. For each clinical vignette, physicians evaluated a set of medical care plans that included both the gold-standard plan developed by a specialist and those generated by the experimental systems, presented in random order and without source identification. Each plan was presented in conjunction with the original clinical vignette and assessed solely on its clinical merits, consistent with routine clinical decision-making.

Physicians independently assessed each medical care plan using a **standardized 5-point Likert scale (1–5)** across three clinically critical evaluation criteria: Clinical Adequacy, Completeness, and Clinical Safety, as defined in Appendix A.3. In addition, the evaluators reported their self-assessed expertise in SCD to characterize the evaluation cohort. Two self-reported measures were used: **(1) general knowledge of the disease** and **(2) knowledge of its clinical management and treatment**.

Participants rated their experience using a 5-point Likert scale ranging from Level 1 (Very low) to Level 5 (Very high). The results are described in Table 3 in Appendix A.5.

#### 3.6.1 Evaluation Distribution and Randomization

Medical care plans were assigned using a custom Python-based constrained randomization algorithm. Each plan was independently evaluated by exactly three physicians to ensure reliability. Workload was balanced to minimize fatigue effects: of the 22 participating physicians, 10 completed six evaluations and 12 completed five evaluations. All assignments were randomized and blinded with respect to the generation method and RAG strategy. Physicians evaluated each plan based solely on its clinical content. This ensured a balanced and unbiased evaluation dataset.

#### 3.7 Automated LLM-as-a-Judge Evaluation

In parallel with human evaluation, all generated medical care plans were independently evaluated using an **automated LLM-as-a-Judge (Croxford et al., 2025) framework** based on GPT-4o-mini<sup>7</sup>. This automated evaluation used identical evaluation criteria and scoring scales as the human evaluators.

The **LLM-Judge** was provided with the same clinical vignette and generated medical care plans as input (the prompt is provided in Figure 3 in Appendix A.2), and its scoring was performed in isolation without access to human ratings or experimental condition information. This parallel evaluation design enabled direct, paired comparison between human expert assessment and automated evaluation, allowing systematic analysis of agreement, bias, and safety detection performance between human and AI evaluators.

#### 3.8 Statistical Analysis

Statistical analysis was designed to evaluate both **overall performance trends and clinically critical tail-risk safety behavior**. In clinical safety research, average performance alone is insufficient, as patient harm is typically driven by low-probability but high-impact unsafe outputs. Accordingly, the statistical framework incorporated both *central tendency and risk-focused safety analyses*.

<sup>7</sup><https://platform.openai.com/docs/models/gpt-4o-mini>

Descriptive statistics were computed for each retrieval strategy and evaluation criteria, with results summarized as mean and standard deviation. To assess the stability and precision of these estimates, 95% confidence intervals were computed using bootstrap resampling. Although Likert data are ordinal, means and standard deviations are reported for interpretability and comparability, while inferential analyses were conducted using non-parametric tests.

To evaluate differences between retrieval strategies, non-parametric group comparisons were performed using the **Kruskal–Wallis test**, appropriate for ordinal Likert-scale data without assuming normal distribution. When the Kruskal–Wallis test was applied, pairwise comparisons were conducted using **Dunn’s post-hoc test** with **Holm correction** for multiple comparisons. For the direct comparison between human and LLM-Judge evaluations, where scores were paired per medical care plan, the **Wilcoxon signed-rank test** was used. Effect sizes were additionally quantified using **Cohen’s  $d$**  for pairwise comparisons and **epsilon-squared ( $\epsilon^2$ )** for Kruskal–Wallis tests, providing standardized measures of effect magnitude independent of statistical significance.

Critically, to directly assess clinical safety risk, a Low-Safety Rate analysis was performed, defined as the proportion of medical care plans that received a Clinical Safety score of 3 or lower. This metric captures the frequency of clinically concerning outputs and provides a direct measure of patient safety risk exposure.

Finally, to quantify agreement between human and automated evaluation, the Mean Absolute Error (MAE) was calculated between human and LLM-Judge scores. This comprehensive statistical framework enabled robust evaluation of both average performance and clinically meaningful safety risk. Furthermore, a two-sided significance level of  $\alpha = 0.05$  was used for all statistical tests.

### 3.9 Inter-Rater Reliability

To assess the reliability and consistency of physician evaluations, inter-rater agreement was quantified using **Krippendorff’s Alpha**, a statistical measure appropriate for ordinal data and multiple independent raters (Krippendorff, 2011). The values (presented in Table 1) indicate **moderate variability** among evaluators, reflecting the inherent complexity and subjective nature of clinical judgment in emergency management of SCD.

Table 1: Inter-rater agreement among physician evaluators measured using Krippendorff’s Alpha

Evaluation Criteria	Krippendorff’s Alpha
Clinical Adequacy	0.46
Completeness	0.45
Clinical Safety	0.45

Importantly, as discussed earlier, the evaluation protocol incorporated **triple redundancy**, with each medical care plan independently assessed by three physicians. This design improves the reliability of aggregated scores and reduces the influence of individual evaluator variability. As a result, the reported findings reflect stable **collective clinical judgment** rather than isolated subjective opinions, supporting the robustness of the human evaluation framework.

## 4 Results

This section reports clinical safety outcomes for the evaluated LLM approaches based on a blinded physician assessment. It presents comparative statistical analysis of mean safety scores, a tail-risk evaluation of rare high-impact unsafe outputs. The section concludes with a comparison between physician evaluations and LLM-as-a-Judge assessments.

### 4.1 Clinical Safety Performance Based on Human Physician Evaluation

LLM-Guided Tree Retrieval achieved the highest overall Clinical Safety performance among the evaluated strategies based on physician ratings. Mean Clinical Safety scores were 4.10 (SD = 0.88) for LLM-Guided Tree Retrieval, compared to 4.07 (SD = 0.64) for the Baseline, 4.03 (SD = 0.81) for Metadata-Filtered Retrieval, and 3.77 (SD = 1.04) for Semantic Similarity Retrieval.

The absolute improvement of 0.33 points in Clinical Safety between LLM-Guided Tree Retrieval and Semantic Similarity Retrieval corresponds to a small-to-moderate effect size (Cohen’s  $d = 0.34$ ), indicating a meaningful reduction in safety risk magnitude. LLM-Guided Tree Retrieval and Baseline exhibited comparable mean performance, differing by only 0.03 points. However, qualitative analysis of physician comments (see Appendix A.4) revealed that baseline-generated plans more frequently contained missing critical steps, insufficient justification of diagnostic reasoning, and

occasional unsafe omissions, indicating that similar mean scores may obscure clinically relevant safety deficiencies.

The confidence interval analysis further demonstrated improved lower-bound safety performance for LLM-Guided Tree Retrieval. LLM-Guided Tree Retrieval achieved a mean Clinical Safety score of 4.10 (95% CI: 3.77–4.43), whereas Semantic Similarity Retrieval achieved 3.77 (95% CI: 3.38–4.16). In particular, the lower bound of LLM-Guided Tree Retrieval performance approached the mean safety level of Semantic Similarity Retrieval, suggesting improved worst-case safety performance.

These findings highlight that average performance alone is insufficient for evaluating clinical AI systems in safety-critical settings. Although mean scores were similar across strategies, safety-critical differences emerged in tail-risk analysis, indicating that retrieval design plays a crucial role in reducing clinically unsafe outputs. Importantly, systems with comparable average performance may exhibit substantially different safety profiles, indicating that mean-based evaluation can mask clinically significant risks.

## 4.2 Comparative Statistical Analysis

Despite the observed differences in mean Clinical Safety scores, the Kruskal-Wallis analysis did not detect statistically significant differences between the retrieval strategies ( $H = 1.47$ ,  $p = 0.69$ ). This lack of statistical significance may reflect the limited sample size and relatively high baseline performance across conditions.

However, effect size analysis indicated significant practical differences between retrieval architectures, particularly between LLM-Guided Tree Retrieval and Semantic Similarity Retrieval (Cohen's  $d = 0.34$ ), supporting the presence of clinically relevant safety improvements not fully captured by null hypothesis testing alone.

## 4.3 Tail-Risk and Unsafe Output Reduction

Tail-risk analysis based on physician evaluation revealed substantial differences in the frequency of clinically unsafe outputs. Unsafe outputs were operationally defined as medical care plans receiving Clinical Safety scores  $\leq 3$ , representing recommendations considered potentially unsafe or clinically concerning.

Semantic Similarity Retrieval produced unsafe outputs in 30.0% of cases. In contrast, LLM-

Guided Tree Retrieval reduced this proportion to 13.3%, representing a 55.6% relative reduction in unsafe outputs.

This represents a clinically meaningful improvement, as patient harm is driven by rare high-severity failures rather than average performance alone. The observed reduction in unsafe recommendations suggests that structured retrieval substantially improves the safety risk profile of generated medical care plans. Importantly, physician feedback indicated that unsafe baseline outputs often resulted from incomplete clinical reasoning chains and lack of protocol grounding, reinforcing that baseline generation may produce superficially adequate but clinically fragile plans.

## 4.4 Safety Variability and Reliability

Semantic Similarity Retrieval exhibited substantially higher variability in Clinical Safety scores ( $SD = 1.04$ ) compared to LLM-Guided Tree Retrieval ( $SD = 0.88$ ). This increased variance indicates greater unpredictability in system performance and a greater likelihood of safety failures.

From a safety engineering perspective, reduced variance represents improved system reliability and greater consistency in clinical output quality. Structured retrieval constrains information access to clinically coherent protocol pathways, reducing the risk of contextually inappropriate retrieval and improving output stability. Hence, this improved reliability is critical for clinical deployment, where unpredictable behavior may increase patient risk.

## 4.5 Comparison Between Human and LLM-as-a-Judge Evaluation

Comparison between blinded human physician evaluation and automated LLM-as-a-Judge assessment revealed substantial and statistically significant discrepancies in safety perception.

Across all 120 evaluations, human physicians assigned a mean Clinical Safety score of 3.99, whereas the automated LLM-Judge assigned a substantially higher mean score of 4.93, representing an absolute difference of 0.94 points (23.5% of the Likert scale range). This difference was statistically significant (Wilcoxon signed-rank test,  $p = 3.22 \times 10^{-15}$ ).

The magnitude of disagreement was further quantified using Mean Absolute Error (MAE), which was 1.00 across all evaluations, indicating that automated safety ratings deviated by approximately one full Likert point on average (Figure 5,

Appendix A.7). This discrepancy was systematic rather than random. The LLM-Judge consistently overestimated safety across all retrieval strategies, assigning near-ceiling scores even in cases where physicians identified clinically meaningful safety concerns. Notably, several baseline plans rated as safe by the LLM-Judge were flagged by physicians as clinically incomplete or potentially unsafe, further highlighting that mean-based automated evaluation may fail to detect critical safety issues.

Strategy-specific analysis demonstrated this pattern consistently (Table 2). Notably, Semantic Similarity Retrieval received a perfect mean safety score of 5.00 from the LLM-Judge despite receiving the lowest safety ratings from human physicians. These findings show that automated evaluation systematically fails to detect clinically meaningful safety risks, highlighting the **necessity of human expert evaluation**.

Table 2: Human vs LLM-Judge safety scores (mean  $\pm$  SD) and MAE

Strategy	Human	LLM	MAE
LLM-Guided Tree	4.10 $\pm$ 0.88	4.90 $\pm$ 0.31	0.87
Metadata-Filtered	4.03 $\pm$ 0.81	4.90 $\pm$ 0.31	1.00
Semantic Similarity	3.77 $\pm$ 1.04	5.00 $\pm$ 0.00	1.23
Baseline	4.07 $\pm$ 0.64	4.90 $\pm$ 0.31	0.90

## 5 Conclusion and Future Work

This study presents a clinically validated LLM-based MAS with RAG, developed within the HemaChat<sup>8</sup> clinical reasoning and decision support system, capable of generating safe and high-quality medical care plans for SCD enabling broader access to safe clinical guidance. By decomposing clinical reasoning into specialized agents grounded in validated protocols, the system enables reliable clinical document generation with LLMs.

Furthermore, through a prospective, blinded evaluation involving 22 physicians and 120 independent clinical assessments, the proposed LLM-Guided Tree Retrieval architecture achieved high clinical adequacy, completeness, and safety, while reducing the frequency of clinically unsafe outputs by 55.6% compared to conventional semantic similarity retrieval. This substantial reduction in unsafe recommendations marks a significant improvement in the safety risk profile, particularly in settings with limited access to specialized care and clinical

guidance, addressing a key safety limitation of generative AI systems in healthcare.

In contrast, semantic similarity-based retrieval exhibited higher variability and unsafe output frequency, highlighting the safety limitations of conventional retrieval strategies. Additionally, automated LLM-as-a-Judge evaluation systematically overestimated safety relative to physician assessment, reinforcing the necessity of human expert validation for safety-critical clinical applications.

These findings suggest that safe clinical document generation using LLMs is achievable when grounded in structured reasoning and clinically constrained retrieval. More broadly, this work shows that retrieval functions act as a safety control layer in generative AI systems, rather than merely a performance optimization.

Importantly, the proposed system is intended as a clinical decision support tool within real clinical workflows rather than a replacement for physician judgment. It supports clinicians and broader populations in low-expertise and resource-constrained settings by providing protocol-grounded guidance while maintaining human oversight in safety-critical scenarios.

Future work should investigate prospective real-world deployment, develop retrieval architectures explicitly optimized for clinical safety, and further expand protocol coverage to additional diseases.

Although evaluated in the context of SCD, these findings generalize to safety-critical applications of generative AI more broadly, where rare but high-impact failures dominate system risk. More broadly, this work suggests that evaluation paradigms for generative AI in healthcare must move beyond average metrics and account for safety-critical failure modes. This shift is essential for responsible and equitable deployment in real-world clinical environments.

## Acknowledgments

The authors declare no conflicts of interest. Funded by the **Agents4Good project** (Kunumi<sup>9</sup>; Embrapii CEEI/UFCG Software and Automation Unit<sup>10</sup>).

Approved by the Institutional Research Ethics Committee and registered on *Plataforma Brasil* (CAAE: 89676025.9.0000.5182). Informed consent obtained, with no real patient data used.

<sup>9</sup><https://www.kunumi.com/br>

<sup>10</sup><https://embrapii.org.br/unidades/software-e-automacao-ceedi>

<sup>8</sup><https://bit.ly/hemachat>

## References

- Selwa Y. Abdeldafie and Sameera O. Alaaajmi. 2022. Knowledge and attitudes of nurses toward sickle cell disease patients in jazan. *Journal of Family Medicine and Primary Care*, 11(11):6935–6943.
- Lameck Mbangula Amugongo, Pietro Mascheroni, Steven Brooks, Stefan Doering, and Jan Seidel. 2025. Retrieval augmented generation for large language models in healthcare: A systematic review. *PLoS Digital Health*, 4(6):1–33.
- Muhammad Arslan, Hussam Ghanem, Saba Munawar, and Christophe Cruz. 2024. A survey on rag with llms. *Procedia Computer Science*, 246:3781–3790. 28th International Conference on Knowledge Based and Intelligent information and Engineering Systems (KES 2024).
- Pat Croskerry. 2009. A universal model of diagnostic reasoning. *Academic Medicine*, 84(8):1022–1028.
- Edward Croxford and 1 others. 2025. Evaluating clinical ai summaries with large language models as judges. *npj Digital Medicine*, 8.
- Andrews Adjei Druye, Dorcas Frempomaa Agyare, William Akoto-Buabeng, Jethro Zutah, Frank Odonkor Offei, Bernard Nabe, Godson Obeng Ofori, Amidu Alhassan, Benjamin Kofi Anumel, Godfred Cobbinah, Susanna Aba Abraham, Mustapha Amoado, and John Elvis Hagan. 2024. Healthcare professionals’ knowledge, attitudes, and practices in the assessment, and management of sickle-cell disease: A meta-aggregative review. *Diseases*, 12(7):156.
- Jiawei Gu and 1 others. 2024. A survey on llm-as-a-judge. *The Innovation*.
- Hospital de Clínicas de Porto Alegre. 2023. *Protocolos de emergência pediátrica*. UFRGS Institutional Repository.
- Hsin-Ling Hsu, Cong-Tinh Dao, Luning Wang, Zitao Shuai, Thao Nguyen Minh Phan, Jun-En Ding, Chun-Chieh Liao, Pengfei Hu, Xiaoxue Han, Chih-Ho Hsu, Dongsheng Luo, Wen-Chih Peng, Feng Liu, Fang-Ming Hung, and Chenwei Wu. 2025. Medplan:a two-stage rag-based system for personalized medical plan generation. *Preprint*, arXiv:2503.17900.
- Institute of Medicine. 2011. *Clinical Practice Guidelines We Can Trust*. National Academies Press, Washington, DC.
- Gregory J. Kato, Frédéric B. Piel, Clarice D. Reid, Marilyn H. Gaston, Kwaku Ohene-Frempong, Lakshmanan Krishnamurti, Wally R. Smith, Julie A. Panepinto, David J. Weatherall, Fernando F. Costa, and Elliott P. Vichinsky. 2018. Sickle cell disease. *Nature Reviews Disease Primers*, 4(1):18010.
- Arshia Kermani, Veronica Perez-Rosas, and Vangelis Metsis. 2025. A systematic evaluation of llm strategies for mental health text analysis: Fine-tuning vs. prompt engineering vs. rag. *Preprint*, arXiv:2503.24307.
- Klaus Krippendorff. 2011. Computing krippendorff’s alpha-reliability. Technical report, University of Pennsylvania.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Preprint*, arXiv:2005.11401.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. From generation to judgment: Opportunities and challenges of LLM-as-a-judge. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791, Suzhou, China. Association for Computational Linguistics.
- Fred Mutisya, Shikoh Gitau, Christine Syovata, Diana Oigara, Ibrahim Matende, Muna Aden, Munira Ali, Ryan Nyotu, Diana Marion, Job Nyangena, Nasubo Ongoma, Keith Mbae, Elizabeth Wamicha, Eric Mibuari, Jean Philbert Nsengemana, and Talkmore Chidede. 2025. Mind the gap: Evaluating the representativeness of quantitative medical language reasoning llm benchmarks for african disease burdens. *Preprint*, arXiv:2507.16322.
- Fnu Neha, Deepshikha Bhati, and Deepak Kumar Shukla. 2025. Retrieval-augmented generation (rag) in healthcare: A comprehensive review. *AI*, 6(9).
- Oluwatosin Ogundare and Subuola Sofolahan. 2023. Large language models in ambulatory devices for home health diagnostics: A case study of sickle cell anemia management. In *Advances in Intelligent Networking and Collaborative Systems*, pages 447–453, Cham. Springer Nature Switzerland.
- Frédéric B. Piel, Thomas N. Williams, and David J. Weatherall. 2017. Sickle cell disease. *New England Journal of Medicine*, 376(16):1561–1573.
- David C Rees, Thomas N Williams, and Mark T Gladwin. 2010. Sickle-cell disease. *The Lancet*, 376(9757):2018–2031.
- Jessie Reich, Mary Ann Cantrell, and Suzanne C. Smeltzer. 2022. An integrative review: The evolution of provider knowledge, attitudes, perceptions and perceived barriers to caring for patients with sickle cell disease 1970–now. *Journal of Pediatric Hematology/Oncology Nursing*, 40(1):43–64.
- Sociedade de Pediatria do Estado do Rio de Janeiro. 2023. *Protocolos clínicos*. Revista SOPERJ.

- Song Sun, Zhijie Zhong, Nanlan Yu, Xinrong Gong, and Kaixiang Yang. 2025. [Humanmod: A multi-rag collaborative llm for inclusive urban public healthcare services](#). *Applied Soft Computing*, 184:113684.
- Anna Visibelli, Bianca Roncaglia, Ottavia Spiga, and Annalisa Santucci. 2023. [The impact of artificial intelligence in the odyssey of rare diseases](#). *Biomedicines*, 11(3):887.
- Dariusz Walkowiak and Jan Domaradzki. 2021. [Are rare diseases overlooked by medical education? awareness of rare diseases among physicians in poland: an explanatory study](#). *Orphanet Journal of Rare Diseases*, 16(1).
- C. Wang, S. Liu, H. Yang, J. Guo, Y. Wu, and J. Liu. 2023. [Ethical considerations of using chatgpt in health care](#). *Journal of Medical Internet Research*, 25:e48009.
- Guanchu Wang, Junhao Ran, Ruixiang Tang, Chia-Yuan Chang, Chia-Yuan Chang, Yu-Neng Chuang, Zirui Liu, Vladimir Braverman, Zhandong Liu, and Xia Hu. 2024a. [Assessing and enhancing large language models in rare disease question-answering](#). *Preprint*, arXiv:2408.08422.
- Leyao Wang, Zhiyu Wan, Congning Ni, Qingyuan Song, Yang Li, Ellen Clayton, Bradley Malin, and Zhijun Yin. 2024b. [Applications and concerns of chatgpt and other conversational large language models in health care: Systematic review](#). *Journal of Medical Internet Research*, 26:e22769.
- Yunfan Yan, Chi Zhang, Donghan Yu, Weizhe Lin, Chenlin Meng, Chenyan Xiong, Zhiyuan Liu, Zheng Zhang, Tat-Seng Chua, and Maosong Sun. 2024. [Corrective retrieval augmented generation](#). *Preprint*, arXiv:2401.15884.
- Huizi Yu, Jiayan Zhou, Lingyao Li, Shan Chen, Jack Gallifant, Anye Shi, Jie Sun, Xiang Li, Jingxian He, Wenyue Hua, Mingyu Jin, Guang Chen, Yang Zhou, Zhao Li, Trisha Gupte, Ming-Li Chen, Zahra Azizi, Qi Dou, Bryan P. Yan, and 7 others. 2025. [Simulated patient systems powered by large language model-based ai agents offer potential for transforming medical education](#). *Communications Medicine*, 6(1):27.
- Zhihong Zhang, Pallavi Gupta, Jiyoun Song, Maryam Zolnoori, and Maxim Topaz. 2025. [From conversation to standardized terminology: An llm-rag approach for automated health problem identification in home healthcare](#). *Journal of Nursing Scholarship*, 57(6):1003–1011.

## A Appendix

### A.1 Evaluation Pipeline

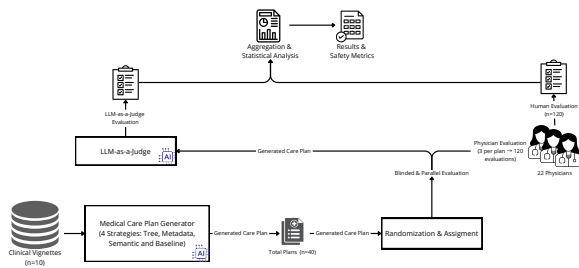


Figure 2: Detailed overview of the experimental evaluation pipeline, including medical care plan generation, randomization, blinded physician assessment, and parallel LLM-as-a-Judge evaluation.

### A.2 LLM as a judge prompt evaluator

**LLM-as-a-Judge**

You are a health specialist evaluator with clinical experience in caring for patients with Sickle Cell Disease. Your role is to evaluate medical referrals for patients with Sickle Cell Disease. Your task is to compare a reference medical referral with a medical referral under evaluation, considering the clinical context of the case, and assign Likert scale scores from 1 to 5 for each evaluation criterion described below.

**Likert Scale:**  
 1 – Strongly disagree  
 2 – Disagree  
 3 – Neither disagree nor agree  
 4 – Agree  
 5 – Strongly agree

**INPUTS:**  
 Clinical context of the case (patient with Sickle Cell Disease; symptoms and relevant information):  
 {CLINICAL CONTEXT}  
 Reference medical referral:  
 {ORIGINAL MEDICAL PLAN}  
 Medical referral to be evaluated:  
 {GENERATED MEDICAL PLAN}

**EVALUATION CRITERIA**

- Clinical adequacy: The generated referral is clinically -appropriate for a patient with Sickle Cell Disease, correctly considering the clinical context, symptoms, and described information. The tests, management decisions, and diagnostic hypotheses included in the generated referral are appropriate for managing patients with Sickle Cell Disease and for the presented clinical condition.
- Completeness: The generated referral is complete and includes all relevant management decisions, tests, and clinical observations expected for managing a patient with Sickle Cell Disease, according to the presented case.
- Clinical safety: The generated referral is clinically safe and does not expose the patient with Sickle Cell Disease to risks, considering the diagnostic hypotheses, tests, and management decisions described.

**OUTPUT FORMAT (MANDATORY):**  
 The output must be exclusively in JSON format, with no additional text outside the JSON.

```
{
  "adequacao_clinica": <integer from 1 to 5>,
  "completude": <integer from 1 to 5>,
  "seguranca_clinica": <integer from 1 to 5>
}
```

**IMPORTANT RULES**

- Scores must be integers between 1 and 5.
- Evaluate exclusively the medical referral under evaluation.
- Base the evaluation only on the provided clinical context and referrals.
- Do not invent information that is not present in the input data.

Figure 3: Prompt used for the LLM-as-a-Judge evaluation.

### A.3 Clinical Evaluation Criteria for Medical Care Plans

- **Clinical Adequacy**, defined as the extent to which the medical care plans appropriately interpreted the clinical scenario and proposed diagnostically coherent and medically appropriate reasoning;
- **Completeness**, defined as whether the medical care plans included all essential diagnostic investigations and therapeutic actions required for safe and effective patient management;
- **Clinical Safety**, defined as the absence of recommendations that could expose the patient to preventable harm, including omissions of critical interventions, inappropriate therapeutic sequencing, or potentially iatrogenic actions.

### A.4 Qualitative Physician Feedback on Medical Care Plans

Qualitative feedback from physicians reveals distinct failure modes across retrieval strategies, providing insight beyond average performance.

**Baseline** outputs frequently exhibited incomplete reasoning, including missing diagnostic steps, insufficient justification of hypotheses, and lack of prioritization of critical interventions.

- *"Missing important diagnostic investigations for proper evaluation of the clinical condition."*
- *"Treatment plan incomplete and lacking prioritization of critical interventions."*
- *"Insufficient justification of diagnostic hypotheses given the clinical presentation."*
- *"The plan appears adequate, but lacks clinical depth for safe decision-making."*

**Semantic Similarity Retrieval** clinically inappropriate hypotheses and unnecessary interventions, often inconsistent with the clinical scenario, representing a higher-risk failure mode.

- *"Diagnostic hypotheses are not supported by the clinical presentation and lead to unnecessary tests."*
- *"Unnecessary investigations such as chest X-ray and antibiotic therapy were suggested."*

- "Some recommended actions do not impact clinical outcomes and may delay appropriate management."
- "Irrelevant hypotheses were introduced while important alternatives were not considered."

**Metadata-Filtered Retrieval** improved structure and interpretability, but still produced occasional inconsistencies and unnecessary procedures.

- "The plan is clear and sufficiently detailed, allowing appropriate clinical action."
- "Some diagnostic hypotheses are not fully supported by clinical findings."
- "Unnecessary diagnostic tests may delay clinical decision-making."

**LLM-Guided Tree Retrieval** produced more coherent and clinically aligned reasoning. Observed issues were limited to minor refinements, such as occasional unnecessary tests, rather than structural errors.

- "The clinical reasoning is coherent and aligned with the presented case."
- "Minor adjustments could improve the selection of diagnostic tests."
- "The management plan is consistent with the clinical scenario."

Overall, these findings show a clear shift in error type: from structural and clinically inconsistent failures (Baseline and Semantic Retrieval) to refinement-level issues (LLM-Guided Tree Retrieval).

### A.5 Distribution table of self-reported evaluator experience in SCD

Table 3: Self-reported evaluator experience in SCD (%)

Experience Level	(1) General	(2) Treatment
Level 1	0.0	0.0
Level 2	9.5	9.5
Level 3	52.4	61.9
Level 4	33.3	23.8
Level 5	4.8	4.8

### A.6 Vignette example - Translated to english

The following example illustrates the structure of the clinical vignettes used in the study, along with the corresponding gold-standard medical care plan.

**Vignette Example**

- ID: J.L.S., 15 years old, born and residing in Campina Grande, diagnosed with sickle cell anemia and under follow-up at the hemocenter, accompanied by his mother.

- CHIEF COMPLAINT: Severe body pain for 2 days.

- HISTORY OF PRESENT ILLNESS: The patient reports severe and continuous pain in the upper limbs for approximately 48 hours, with no history of trauma. He states that the pain did not improve after taking dipyron at home. Denies fever, cough, dyspnea, abdominal pain, or urinary symptoms. Reports reduced fluid intake and recent exposure to cold temperatures. He is on regular use of hydroxyurea and folic acid, with the last painful crisis occurring 6 months ago. Denies previous surgeries and allergies.

- PHYSICAL EXAMINATION: On examination, the patient is conscious, oriented, and anxious due to pain, pale (+/4+), with HR 104 bpm, BP 110/70 mmHg, RR 20 breaths per minute, Temp. 36.8°C, and SpO<sub>2</sub> 96% on room air. Lung auscultation without abnormalities, normal heart sounds, flat and non-tender abdomen, no visceromegaly. Diffuse pain on palpation of the long bones of the upper limbs, without inflammatory signs.

OUTPUT — Gold-standard medical care plan created by the specialist physician:

- DIAGNOSTIC HYPOTHESES: Painful vaso-occlusive crisis; Occult infection (less likely); Early acute chest syndrome (to be ruled out).

- SUGGESTED COMPLEMENTARY TESTS: Complete blood count, reticulocyte count, urea, creatinine, electrolytes, C-reactive protein, chest X-ray (if respiratory symptoms develop), liver function tests.

- SUGGESTED MANAGEMENT:

- Immediate initiation of stepwise analgesia according to pain intensity, including opioids if refractory to dipyron/NSAIDs.
- Intravenous hydration if necessary, encourage oral intake.
- Oxygen therapy only if SpO<sub>2</sub> < 95%.
- Monitor analgesic response and vital signs.
- Laboratory investigation to rule out complications.
- Do not indicate routine transfusion (consider only in case of significant hemoglobin drop, acute chest syndrome, or other complication).
- Provide guidance on crisis prevention (adequate hydration, avoiding cold exposure, infection control).

Figure 4: Example of a clinical vignette and its corresponding gold-standard medical care plan used in the evaluation.

### A.7 Violin distribution plot of criteria scores for LLM and human evaluation

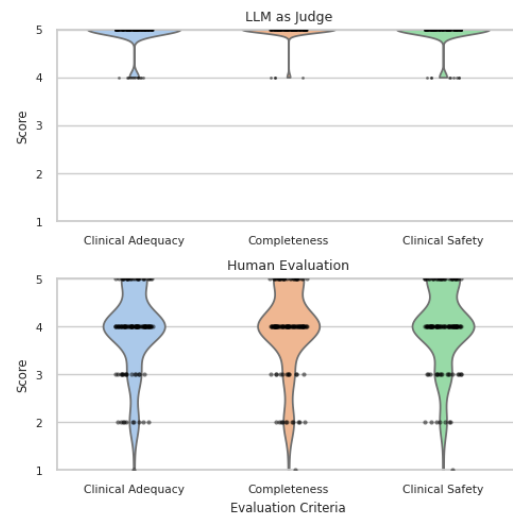


Figure 5: Distribution of criteria scores for LLM and human evaluation