Investigating RAG-based Approaches in Clinical Trial and Patient Matching

Daniel Leon Tramontini Shrestha Ghosh Carsten Eickhoff

University of Tuebingen, Germany

DANIEL.LEON-TRAMONTINI@STUDENT.UNI-TUEBINGEN.DE SHRESTHA.GHOSH@UNI-TUEBINGEN.DE CARSTEN.EICKHOFF@UNI-TUEBINGEN.DE

Abstract

The task of matching clinical trials and patients involves predicting whether a patient meets the eligibility criteria of a clinical trial, via evidences from patient records, such as clinical notes. Given that both the trial eligibility criteria and the clinical notes of patients are unstructured texts, Large Language Models (LLMs) hold the potential for improving performance on this task. Current methods use Retrieval-Augmented Generation (RAG) in order to predict patient eligibility for an eligibility criterion.

In this work, we systematically investigate three aspects of RAG-based approaches: (i) the complexity of the task, (ii) data retrieval for longitudinal records, and (iii) the effect of abstention on prediction quality. We show that criteria complexity and abstention have noticeable effects on model performance, while the choice of embedding models and ranking methods has no effect on the retrieved evidences from patient history. We hope that findings from our study encourage research in understanding the impact of RAG components in other clinical decision-making tasks.

Keywords: clinical trial and patient matching, retrieval-augmented generation, clinical decision-making

Data and Code Availability This paper uses the N2C2 cohort selection dataset by Stubbs et al. (2019). It comprises 288 de-identified patient records and their eligibility labels for 13 eligibility criteria. More details are provided in Section 4. The data is currently unavailable for download. Our code is available online

at https://github.com/leontramontini97/clinical_trial-patient_matching.

Institutional Review Board (IRB) This research does not require an IRB approval.

1. Introduction

The task of matching clinical trials and patients is challenging. Key clinical information determining patient eligibility is often buried in the clinical notes of longitudinal patient records. Recent works by Jin et al. (2024); Wornow et al. (2025); Li et al. (2025) highlight the potential of LLMs in predicting patient eligibility per criterion and generating explanations for the same.

The state-of-the-art methods tackle the trial and patient matching problem via trial-centric or patient-centric approaches. Trial-centric approaches identify patients relevant to a particular clinical trial, and, patient-centric approaches identify clinical trials relevant to a patient. Trial-centric solutions, such as, Beattie et al. (2024); Wornow et al. (2025); Li et al. (2025), use the N2C2 Cohort Selection benchmark (Stubbs et al., 2019). This benchmark focuses on 13 inclusion criteria for diabetic patients with heart conditions, and 288 longitudinal patient records. Patient-centric methods, such as, Jin et al. (2024); Rybinski et al. (2024), use the SIGIR (Koopman and Zuccon, 2016) and the TREC (TREC Biomedical Tracks) benchmarks. These benchmarks have short and few patient descriptions (< 100) and cohort keywords used for retrieving relevant trials from a large trial registry, such as the ClinicalTrials.gov, with hundreds and thousands of clinical trials.

Evaluations of these methods under the current benchmarks (Stubbs et al., 2019; Koopman

and Zuccon, 2016; TREC Biomedical Tracks) focus on accuracy-based metrics. Bedi et al. (2025); Omar et al. (2024); Nemati et al. (2025) have stressed that focusing only on accuracy-based measures does not provide any information about LLM uncertainty, fairness and bias.

From identifying clinical conditions and nonclinical data (demographic details) to reasoning over clinical data and handling ambiguous conditions, eligibility criteria have a range of complexity. In their N2C2 cohort selection task report, Stubbs et al. (2019) note that the lowest performing criteria are those that require complicated reasoning with temporal modifiers, and those which required inference. LLM-based systems deploy strategies from RAG to prompt engineering. Yet, it is unclear how these strategies tackle criteria complexity.

In this work, we set up a modular RAG pipeline for the task of trial and patient matching and investigate three aspects that affect LLM-based systems tackling this task.

Firstly, we characterize the complexity of the task by annotating the number of entities and relations in the eligibility criteria. We propose generalized strategies to infer implicit entities and report their effect on the final performance. Secondly, we evaluate different embedding models and ranking strategies of the data retrieval for longitudinal patient records. Thirdly, we examine the effect of abstention on prediction quality via accuracy and verbalized confidence.

2. Related Work

Matching clinical trials and patients is resourceintensive and was done manually by experts with limited automation until only a decade ago (Penberthy et al., 2012). Starting from early attempts with supervised machine learning on manually selected features and learned features (Zhang and Demner-Fushman, 2017; Vazquez et al., 2021), the current transformer-based LLMs have created unprecedented leaps in trial recruitment tasks (Jin et al., 2024; Wornow et al., 2025). Due to the sensitive nature of the task and the lack of public datasets of retrospective matches, a majority of the LLM-based methods (Roberts et al., 2022; Jin et al., 2024; Rybinski et al., 2024) are evaluated on the patient-centric trial recommendation task from the Text REtrieval Conference (TREC) tracks on clinical trials¹. The notable trial-centric public dataset with longitudinal patient data is the 2018 N2C2 cohort selection task (Stubbs et al., 2019). Another line of work transforms clinical eligibility criteria into logical formats to be applied on structured patient databases (Yuan et al., 2019).

Previous benchmarks rely on accuracy-based measures of precision, recall and F1 for evaluating model performance in matching clinical trials and patients (Stubbs et al., 2019; Soboroff, 2021; Roberts et al., 2022). However, LLMs present new challenges in transparency and accountability. This has led to studies focused on multidimensional evaluation of LLM applications (Bedi et al., 2025) and new checklists for studies on the development and evaluation of LLMs (Tripathi et al., 2025).

3. Trial-Patient Matching

We investigate recent methods, such as, Beattie et al. (2024); Wornow et al. (2025); Li et al. (2025), that use LLMs for the trial and patient matching task. These methods follow the RAG approach, where the generation model is prompted with relevant parts of the patient record as evidences to guide the prediction of the LLMs. Table 1 breaks down the methods by the embedding models, the generative models, the prompting techniques and output complexity.

Model Size We see that all models use smaller embedding models (order of parameters less than a billion), with MiniLM (Reimers and Gurevych, 2019) and BioBERT (Lee et al., 2019) having 22.7M and 65M parameters, respectively. The generation models are much larger, from 8B parameters for Llama-3-8B-Instruct (Dubey et al., 2024) to estimated hundreds of billions of parameters in GPT-4.

Prompting Technique Two of three methods supplement criteria definition with modifications to the original criteria (Wornow et al., 2025) or with tips on how to resolve each criterion (Beattie et al., 2024). They do not provide when and why these criteria modifications work.

Output Complexity The nature of output from the generation model ranges from simple

^{1.} https://www.trec-cds.org/

Method	Embedding Model	Generation Model	Prompting Technique	Output Complexity
Beattie et al. (2024)	ada-002 (Greene et al., 2022)	GPT-4 (Achiam et al., 2023)	Personified role Manually crafted tips for all criteria Response format: JSON object	Return criterion name, supporting evidence, and binary prediction.
Wornow et al. (2025)	MiniLM-L6-v2 (Reimers and Gurevych, 2019)	GPT-4	Personified role Manually modified criteria definitions (6 of 13) Response format: JSON object	Return criterion name, list of medications, rationale, binary prediction, and confidence rating
Li et al. (2025)	BioBERT (Lee et al., 2019)	Llama-3-8B-Instruct (Dubey et al., 2024)	AI assistant role Response format: Text	Return a binary prediction

Table 1: Break down of the models used for embedding and generation, prompting technique and output complexity.

text-based binary values (Li et al., 2025) to JSON objects comprising prediction, criteria name and text evidence (Beattie et al., 2024). Wornow et al. (2025) expect additional items, such as, medication list, rationale and confidence rating (high, medium, low).

3.1. RAG Pipeline

Data Processing We divide each patient record into smaller chunks and create a vector-store of the chunk embeddings. We additionally maintain the timestamp of the patient visit corresponding to each chunk. The eligibility criterion are annotated (more details in Section 5) and prompting techniques are adopted to make the criteria more explicit. We track the effects of these prompting techniques on the final performance.

Data Retrieval Given a criterion, it is embedded using the same embedding model we use for the patients. We follow standard information retrieval and fetch the top-most relevant chunks computed using the cosine similarity of the criterion and chunk embeddings. In Section 6, we investigate if retrievers indeed return chunks with relevant information for eligibility prediction.

Answer Generation We formulate prompts with the criterion and the retrieved chunks for the LLM to predict eligibility. The prompt template comprises general instructions, the task and the output format.

Prompt template Based on the following clinical record excerpt, determine if the

patient meets this criterion:

Criterion: [criteria]

Relevant Clinical Context: [Clinical note excerpt]

Provide your analysis in JSON format with the following structure:

"status": "met" or "not met",

"justification": "Brief explanation with specific evidence from the clinical context"

If the information is insufficient, you must still make a determination of either "met" or "not met" based on the available evidence.

In Section 7, we examine the uncertainty in the model via verbalized certainty values and how this changes when the model is given the option to abstain.

4. Experiment Setup

Dataset We work with the 2018 N2C2 cohort selection Stubbs et al. (2019). This public dataset comprises eligibility labels for 288 patients on 13 eligibility criteria. The patient records are de-identified longitudinal records, with an average of 2711 tokens per patient. The criteria labels and their definitions are in the Appendix A. We report our results on the test set of 86 patients, comprising (86×13) 1118 patient-criteria labels. We chose the N2C2 dataset over patient-centric datasets, such as, the SI-GIR 2016 (Koopman and Zuccon, 2016) and

Methods	Strategy	Macro-F1	Micro-F1
Oleynik et al. (2019) (N2C2 2018 best)	Rule-based classifier	0.75	0.91
Beattie et al. $(2024)^a$ Wornow et al. (2025)	RAG RAG	0.75* 0.81	$0.86 \\ 0.93$
Li et al. (2025) (LLM-Match)	RAG with a trained classification head	0.86	-
Our method	RAG with original crieria RAG with improved crieria	0.70 0.81	0.76 0.86

^{*}Derived from per criteria scores reported by Beattie et al. (2024).

Table 2: SOTA performance on N2C2 cohort selection dataset of 288 patients and 13 criteria (train = 202; test = 86). We derived the missing metrics when data was available.

the clinical trial tasks from the TREC Biomedical tracks (TREC Biomedical Tracks), since the N2C2 dataset has criterion-level labels on longitudinal patient datasets. This helps us to explore the relationship between complexity of the eligibility criteria and the system performance.

Baselines We compare our method with the following state-of-the-art (SOTA) LLM-based methods, Wornow et al. (2025), Li et al. (2025) and Beattie et al. (2024), in Table 2. We also report the best method from the original N2C2 task for completeness Oleynik et al. (2019). All methods, except for Oleynik et al. (2019), use Retrieval-Augmented Generation (RAG). Li et al. (2025) also train a classifier on top of the LLM, but this provides only marginal improvement (macro-F1 increases by 0.01 to 0.86).

Implementation Details Each patient record is split text into 500-character chunks with a 50-character overlap. We use OpenAI's GPT-40 class of models for answer generation and for our LLM-as-a-judge evaluations. We use Open-AI's model text-embedding-ada-002 to generate patient embeddings and store it in a FAISS index. Our code is available online at https://github.com/leontramontini97/clinical_trial-patient_matching.

5. Criteria Complexity

The eligibility criteria, expressed in natural language, is often modified to improve matching with patient data Beattie et al. (2024); Wornow et al. (2025). For instance, in the original N2C2 dataset by Stubbs et al. (2019), two criteria "Advanced cardiovascular disease (CAD)" and

Criteria Label	#Entities	#Implicit	#Relations	
DRUG-ABUSE	2	1	1	
ALCOHOL-ABUSE	3	2	2	
ENGLISH	1	1	-	
MAKES-DECISIONS	1	1	-	
ABDOMINAL	5	2	3	
MAJOR-DIABETES	3	2	2	
ADVANCED-CAD	2	1	1	
MI-6MOS	2	1	1	
KETO-1YR	3	1	2	
DIETSUPP-2MOS	4	2	3	
ASP-FOR-MI	4	1	2	
HBA1C	3	1	2	
CREATININE	2	1	1	
Average	2.6	1.2	1.4	

Table 3: Annotating defining characteristics of the N2C2 eligibility criteria.

"Major diabetes-related complication" were further clarified, specifically the conditions satisfying the terms "advanced" and "major complications" were laid out for gold standard annotation by experts.

Criteria Characteristics Trial eligibility criteria vary in semantic complexity. Each criterion describes a central entity via relationships to attributes and logical operations. A criterion can range from being fully objective (HbA1c value between 6.5% and 9.5%) to subjective and ambiguous (major complications). Criteria can also be wither disease-specific or disease-agnostic, such as patient demographics and decision-making capability.

As such, this variability in criteria makes comparison between trials quite challenging. Following the Chia annotation model Kury et al. (2020), we annotate the entities and relations in N2C2 selection criteria (see B). Often criteria assume im-

^aScores reported on a subset of the test set with 40 patients.

plicit medical knowledge, such as measurements of a lab value above normal levels, adding another layer of complexity. Hence, we also annotate each entity as being implicit or explicit. In Table 3, we provide the number of entities, relations, and implicit entities in every criterion of the N2C2 dataset. The average criteria in the dataset has 2.6 entities, about half of which are implicit (1.2) with 1.4 relations between the entities. We find two criteria, ABDOMINAL and DIETSUPP-2MOS with the highest relations (3) and entities (5 and 4, respectively).

We identified the following classes of criteriaspecific amendments. We test these amendments on the same set of retrieved patient chunks, to not introduce any confounding effect from the retriever. The effects of criteria-specific amendments are reported in Table 4.

Extended Description In the original task, the definitions of two criteria were extended with explicit conditions to reduce ambiguity and subjectivity for the expert human annotators. For instance, the term "advanced" in ADVANCED-CAD was defined to be constrained to two or more of four specific observations. Furthermore, the term "major complication" for MAJOR-DIABETES was confined to any of six conditions that are strongly correlated with uncontrolled diabetes. We include these extended definitions in the criteria description, during the data processing. We also extend the definition of ABDOM-INAL to include examples of intra-abdominal surgeries and rephrase the original criteria to improve clarity.

We see a huge jump in recall for ABDOMI-NAL, from 0.167 to 0.667 with the extended definition. Clarifying the original definition into two separate conditions - history of an intra-abdominal surgery or small bowel obstruction - and explicitly specifying examples of the former condition spanning types of intra-abdominal procedures, we see that the model becomes better at picking up the more instances of criterion eligibility.

We also see an improvement in precision for MAJOR-DIABETES from 0.7 to 0.875. This makes sense, since the original definition requires the model to infer whether a condition is diabetes related or not and whether it is a major complication. The new definition simplifies this

to specifically define the conditions that qualify these conditions.

While we measure positive changes in two criteria, interestingly, ADVANCED-CAD does not improve much. Upon reexamining the modified version of the criterion, we find the new definition in fact introduces more complexity. At least two of four conditions must be satisfied for eligibility. The high false positive rate, 0.585, and a low false negative rate, 0.062 also support this. Upon examining the reasoning for some false positives, we find that even though LLM provides two conditions that are met, these are in fact incorrect.

Explicit Default Decision The two criteria ENGLISH and MAKES-DECISION have an implicit condition that it is assumed that the patient is eligible for both unless evidence to the contrary is present in the dataset. In this case, the LLM is explicitly instructed to return that the criteria is met unless there is contradictory evidence.

The recall for both these criteria increase, from 0.603 to 1 for ENGLISH, and from 0.662 to 0.903 for MAKES-DECISION. While both criteria are defined as presence of an observation (speaks English, is capable of making decisions), they come up in the clinical notes only when they are not met. With non-English speakers, it is explicitly noted which language they speak and whether they have an interpreter. The evidence that a patient cannot make a decision is implicitly conveyed through the observations and diagnosis of mental health of the patient.

Explicit Temporal Tagging There are three criteria that require temporal decision making. This involves first determining the most recent record of the patient, and then whether the criterion is fulfilled within the time frame mentioned in the criterion.

In order to tackle this, we add the date of the most recent patient visit for the criteria, KETO-1YR, MI-6MOS and DIETSUPP-2MOS. Further, we include specific instructions on how to handle temporal context.

Temporal-specific instructions

IMPORTANT TEMPORAL CONTEXT: This criterion has a time constraint. Pay special at-

C++	Chitania tan		Before			After		
Strategy	Criteria tag	P	\mathbf{R}	$\mathbf{F1}$	P	\mathbf{R}	$\mathbf{F}1$	
Extended Description	ABDOMINAL	1	0.167	0.286	0.870	0.667	0.755	
Extended Description	MAJOR-DIABETES	0.7	0.814	0.753	0.875	0.814	0.843	
	ADVANCED-CAD	0.636	0.933	0.757	0.624	1	0.769	
Default Decision	ENGLISH	0.957	0.603	0.740	0.913	1	0.954	
Default Decision	MAKES-DECISIONS	0.982	0.662	0.791	0.974	0.903	0.937	
	KETO-1YR*	-	-	-		-	-	
Temporal Tagging	MI-6MOS	0.417	0.625	0.5	1	0.875	0.933	
	DIETSUPP-2MOS	0.813	0.886	0.847	0.875	0.814	0.843	
N	HBA1C	1	0.657	0.793	1	0.657	0.793	
Numerical limits	CREATININE	1	0.627	0.769	1	0.627	0.769	
	ALCOHOL-ABUSE	0.75	1	0.857	0.75	1	0.857	
N	DRUG-ABUSE	0.3	1	0.462	0.3	1	0.462	
None	ASP-FOR-MI	0.9	0.926	0.913	0.9	0.926	0.913	
	Overall average	0.788	0.742	0.706	0.840	0.857	0.819	

Table 4: Performance per criterion before and after modification. KETO-1YR has no positive labels.

tention to dates and timing. The clinical record uses synthetic dates where 2-digit years should be interpreted as 2 XXX (e.g., "2/16/51" = "2051-02-16" or "2051-02-16" depending on the case, not "1951-02-16"). The reference date (most recent clinical note) is: \langle FROM-PATIENT-RECORD \rangle . This is our present moment. When evaluating time-based criteria, calculate time intervals from events to this reference date. Example: If the reference date is 2151-04-11 and an event occurred on "2/16/51" (2151-02-16), that's about 2 months prior, which IS within The past \langle CRITERION-SPECIFIC-TIME \rangle .

While we cannot comment on the effect of the temporal tagging and instructions on KETO-1YR since it is not present in the test data, we see a sharp increase in precision and recall for MI-6MOS, from 0.417 to 1 precision and 0.625 to 0.875 recall. There is no overall effect in DIETSUPP-2MOS. On further inspection, we find that the errors stem from the LLM's inability to identify dietary supplements.

Explicit Numerical Limits Two criteria require interpretation of lab results: one mentioned explicitly, *i.e.*, HBA1C, whose values must be between 6.5% and 9.5%, and creatinine, and one implicitly, *i.e.*, CREATININE, whose values should be above normal limits. A third criterion, ALCOHOL-ABUSE, also requires an inference of

whether current alcohol consumption is over recommended limits.

The original task does not specify the actuals limits for creatinine and alcohol consumption used by the annotators. Given the variance of these limits and a lack of common standards, it is difficult to apply explicit limits, but we can inspect the reasonings LLMs provide to infer the limits the LLM might have used to predict the eligibility. In the case of creatinine, we find that the LLMs assume the upper limit in the range of 1.3 to 1.5 mg/dL and make the decision, while in other cases, it refuses to make a decision to the lack of a specified upper limit. In the case of alcohol abuse, the model makes predictions only when the evidence very clear, such as the patient consumes alcohol multiple times daily. In many other cases, the model specifies the recommended limit as 7 or 14 drinks per week for men, but does not make a decision, citing lack of specific limits.

It is interesting to note that the model defaults to known limits for men. This can lead to biased decision when the model is not provided with clear specifications, especially in cases where men and women have different standard limits.

6. Data Retrieval

Due to the longitudinal nature of the patient data, it is stored in chunks. Then, for matching

Information per prompt		
Default (1 criterion, relevant chunks)	0.76	
1 criterion, full patient record	0.63	
all criteria, full patient record	0.67	

Table 5: Effect of using the entire patient record vs relevant chunks on the F1 score.

a criterion only the relevant chunks are retrieved and sent to the LLM to reduce noise.

Effect of Chunking In Table 5, we report the F1 scores of the model when it is prompted with the entire patient record under two conditions. Once with a single criterion per patient, and then with all criteria per patient. The LLM's performance decreases when provided with the entire patient record, to 0.67 F1 when all criteria are prompted at once and to 0.63 when prompted with a single criterion. In a separate experiment Wornow et al. (2025) showed that as the number of relevant chunks provided keeps increasing, the performance plateaus, but never drops. This provides evidence of position bias in LLMs. The LLMs focus on the top evidences, hence showing a plateaued gain when presented with all chunks in a record, but ranked versus a drop in performance when the full record is passed as is.

Effect of Embedding Models We compare three embedding models used by our baselines: Open AI's text-embedding-ada-002, Sentence Transformer's all-MiniLM-L6-v2 and BioBERT model dmis-lab/biobert-v1.1 and found no difference in the chunks and the order in which they were returned. We took the top 10 chunks returned by these models for every patientcriterion pair and measured the Jaccard index between all pairs of models, computed as the intersection over the union of two sets. The average Jaccard index of the three pairs was 0.298 ± 0.05 . An LLM-as-a-judge evaluates the top-10 chunks for sufficiency averaged over a random selection of 10 patients. It returns insufficient information 63% of the time for BioBERT, followed by MiniLM (61%) and ada-002 (60%). In cases, where there is enough information to make a decision, the accuracy of the ada-002 embedding model is the highest at 80%, followed by MiniLM at 78% and BioBERT at 77%.

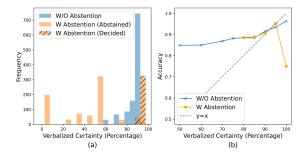


Figure 1: (a) Confidence distribution with and without abstention. (b) Accuracy at confidence thresholds, s.t. accuracy at x is computed for predictions with a conf. $\geq x$.

Effect of Ranking We compare ranking vector embeddings using the FAISS similarity search with BM25 and diversity rankings using maximal marginal rankings (MMR). These methods do return quite different chunks, the average Jaccard similarity of the retrieved chunks with the original vector embeddings in 0.18. When we ask an LLM-as-a-judge whether the top-10 returned chunks are sufficient to make a prediction, it returns yes 25% of the times for the FAISS similarity search, only 18.9% of the time for MMR, and 16.5% of the time for BM25.

7. Answer Consistency

Certainty vs. Abstension We conduct two experiment where the LLM is asked to additionally to output a verbalized confidence of its prediction, between 0-100%. In the second run, we additionally provide the option to the model to abstain. Figure 1 shows the main results. We also experiment with confidence expressed in levels between 1-5. We found a high Pearson correlation of 0.8 (p-value; 0.05) between the different modes of verbalizations.

From Figure 1 (a) we can see that when forced to predict either "met" or "not met", the model mostly also labels its predicts as being highly confident, with a mean confidence level of 88% \pm 8.2%. The distribution is more spread out when the model is allowed to abstain. We now see that there are cases where the model is much less confident, with the average confidence equal to 54%

 \pm 33.5%. Further, the model abstains from prediction if certainty drops below 75%.

Abstaining helps the model improve overall accuracy from 0.84 to 0.88, with a trade-off of an abstention rate of 66.2%. From Figure 1 (b), we see that, the model is overconfident only at 100% confidence level where the accuracy of the model drops to 0.75. The change in certainty in cases where the model does make a decision compared to non-abstention is $0.98\% \pm 2.4\%$. This means that the verbalized confidence is stable across prompts for high certainty values (>80%).

Abstention Justifications The justification for abstention in 59% of the cases provided is "insufficient evidence", followed by justifications that mention terms denoting lack of information or no evidence (38.1%). There were 7 instances where the model reasons that there is conflicting information. At a criterion level, the inability to determine if at least two conditions are met for ADVANCED-CAD (2.5%) and a lack of reference range for lab values for creatinine (1.7%) are the common justifications for abstention.

8. Discussion and Limitations

Automated Criteria Complexity The results in Section 7, highlight the importance of criteria complexity and strategies that can be targeted to criteria classes, once we know the entities and relations in the criteria. We highlight the importance of implicit criterion and how this affects the overall LLM performance. Especially sensitive are cases where the model might assume standards applicable only to certain groups, such as the normal upper creatine values or weekly alcohol limits for men, putting other groups at a disadvantage. While our study covers only one clinical trial, extending this to other trials on a scale requires an automated method of annotating eligibility criteria and finding potential implicit information needs crucial for decisionmaking.

Evaluating LLM Rationale While we provide qualitative anecdotes of LLM rationale, there is a need for a more systematic study of evaluating LLM rationale. These rationales indicate absence of sufficient information, indicating information gaps or criteria which depend heavily on lack of evidence as evidence itself, such as

no information in English-speaking abilities imply patient can speak English. They also expose model biases, as we saw in the cases of ALCOHOL-ABUSE and CREATININE, where the model may default to using known information applicable to a particular group, such as recommended limits for men, putting other groups at a disadvantage.

Lack of Rich Data Our study on the N2C2 data comes with caveats: it is trial-specific and has high skew in the label distribution of some criteria. Our next step would be to check the generalizability of our findings regarding criteria importance and model stability on other datasets, and to achieve that it is important to develop automated criteria annotators, as discussed above. We also lack data on ground-truth explanations and standard quality tests for evaluating LLM rationale.

In our experiments on abstention in Section 7, the abstention rate was 66.2% with more than 97% of the justifications related to insufficient or no evidence, that the LLM was predicting under ordinary circumstances with no option of abstaining. This highlights the importance of tracking implicit information needs which medical experts do not necessarily need, but is required for LLM-based predictions for more transparent decision-making.

9. Conclusion

In this paper, we investigate RAG-based approaches for the task of matching clinical trial and patients via three aspects: criteria complexity, retrieval and answer generation. We characterize the complexity of eligibility criteria by the number of entities and relationships they contain and the number of implicit entities that need resolving. We show generalizable techniques to effectively tackle groups of implicit entities that can lead to a performance metrics (F1 $0.706 \rightarrow$ 0.819). We find that while different embedding models and ranking methods retrieve quite different chunks, LLM-as-a-judge evaluates the majority of the cases as insufficient evidence for predicting eligibility with little variation in overall accuracy. Finally, we show that LLMs selfreporting confidence can be unreliable and abstention reveals decisions LLMs make in the face of insufficient information. We hope that our findings encourage future research in improving LLM-based clinical decision-making.

Acknowledgments

The first author was financially supported as a student research assistant by Boehringer Ingelheim .

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Jacob Beattie, Sarah Neufeld, Daniel Yang, Christian Chukwuma, Ahmed Gul, Neil Desai, Steve Jiang, and Michael Dohopolski. Utilizing large language models for enhanced clinical trial matching: A study on automation in patient screening. Cureus, 16(5), 2024.
- Suhana Bedi, Yutong Liu, Lucy Orr-Ewing, Dev Dash, Sanmi Koyejo, Alison Callahan, Jason A. Fries, Michael Wornow, Akshay Swaminathan, Lisa Soleymani Lehmann, Hyo Jung Hong, Mehr Kashyap, Akash R. Chaurasia, Nirav R. Shah, Karandeep Singh, Troy Tazbaz, Arnold Milstein, Michael A. Pfeffer, and Nigam H. Shah. Testing and Evaluation of Health Care Applications of Large Language Models: A Systematic Review. JAMA, January 2025. doi: 10.1001/jama.2024.21700. URL https://doi.org/10.1001/jama.2024.21700.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv e-prints, pages arXiv-2407, 2024.
- Ryan Greene, Ted Sanders, Lilian Weng, and Arvind Neelakantan. New and improved embedding model. https://openai.com/index/new-and-improved-embedding-model/, 2022. Last accessed: 2025 Sep 03.

- Qiao Jin, Zifeng Wang, Charalampos S Floudas, Fangyuan Chen, Changlin Gong, Dara Bracken-Clarke, Elisabetta Xue, Yifan Yang, Jimeng Sun, and Zhiyong Lu. Matching patients to clinical trials with large language models. *Nature communications*, 15(1):9074, 2024.
- Bevan Koopman and Guido Zuccon. A test collection for matching patients to clinical trials. In Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval, pages 669–672, 2016.
- Fabrício Kury, Alex Butler, Chi Yuan, Li-heng Fu, Yingcheng Sun, Hao Liu, Ida Sim, Simona Carini, and Chunhua Weng. Chia, a large annotated corpus of clinical trial eligibility criteria. *Scientific Data*, 7(1), 2020.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4): 1234–1240, 2019.
- Xiaodi Li, Shaika Chowdhury, Chung Il Wi, Maria Vassilaki, Xiaoke Liu, Terence T Sio, Owen Garrick, Young J Juhn, James R Cerhan, Cui Tao, et al. Llm-match: An open-sourced patient matching model based on large language models and retrieval-augmented generation. arXiv preprint arXiv:2503.13281, 2025.
- Ali Nemati, Mohammad Assadi Shalmani, Qiang Lu, and Jake Luo. Benchmarking large language models from open and closed source models to apply data annotation for free-text criteria in healthcare. Future Internet, 17(4): 138, 2025.
- Michel Oleynik, Amila Kugic, Zdenko Kasáč, and Markus Kreuzthaler. Evaluating shallow and deep learning strategies for the 2018 n2c2 shared task on clinical text classification. *Journal of the American Medical Informatics Association*, 26(11):1247–1254, 2019.
- Mahmud Omar, Girish N Nadkarni, Eyal Klang, and Benjamin S Glicksberg. Large language

models in medicine: a review of current clinical trials across healthcare applications. *PLOS Digital Health*, 3(11):e0000662, 2024.

Lynne T Penberthy, Bassam A Dahman, Valentina I Petkov, and Jonathan P DeShazo. Effort required in eligibility screening for clinical trials. *Journal of Oncology Practice*, 8(6): 365–370, 2012.

Nils Reimers and Iryna Gurevych. Sentencebert: Sentence embeddings using siamese bertnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, 2019.

Kirk Roberts, Dina Demner-Fushman, Ellen M Voorhees, Steven Bedrick, and William R Hersh. Overview of the trec 2022 clinical trials track. In *TREC*, 2022.

Maciej Rybinski, Wojciech Kusa, Sarvnaz Karimi, and Allan Hanbury. Learning to match patients to clinical trials using large language models. *Journal of Biomedical Informatics*, 159:104734, 2024.

Ian Soboroff. Overview of trec 2021. In TREC, 2021.

Amber Stubbs, Michele Filannino, Ergin Soysal, Samuel Henry, and Özlem Uzuner. Cohort selection for clinical trials: n2c2 2018 shared task track 1. Journal of the American Medical Informatics Association, 26(11):1163–1171, 2019.

TREC Biomedical Tracks. https://www.trec-cds.org/, 2024. Last accessed: 2024 Oct 31.

Satvik Tripathi, Dana Alkhulaifat, Florence X Doo, Pranav Rajpurkar, Rafe McBeth, Dania Daye, and Tessa S Cook. Development, evaluation, and assessment of large language models (deal) checklist: A technical report, 2025.

Janette Vazquez, Samir Abdelrahman, Loretta M Byrne, Michael Russell, Paul Harris, and Julio C Facelli. Using supervised machine learning classifiers to estimate likelihood of participating in clinical trials of a de-identified version of researchmatch. *Journal of Clinical and Translational Science*, 5(1):e42, 2021.

Michael Wornow, Alejandro Lozano, Dev Dash, Jenelle Jindal, Kenneth W Mahaffey, and Nigam H Shah. Zero-shot clinical trial patient matching with llms. *NEJM AI*, 2(1): AIcs2400360, 2025.

Chi Yuan, Patrick B Ryan, Casey Ta, Yixuan Guo, Ziran Li, Jill Hardin, Rupa Makadia, Peng Jin, Ning Shang, Tian Kang, and Chunhua Weng. Criteria2Query: a natural language interface to clinical databases for cohort definition. Journal of the American Medical Informatics Association, 26(4):294–305, February 2019. doi: 10.1093/jamia/ocy178. URL https://doi.org/10.1093/jamia/ocy178.

Kevin Zhang and Dina Demner-Fushman. Automated classification of eligibility criteria in clinical trials to facilitate patient-trial matching for specific patient populations. *Journal of the American Medical Informatics Association*, 24(4):781–787, 2017.

Appendix A. Eligibility Criteria Definition

Table 6, lists the original criteria definitions for each criteria label from the 2018 N2C2 cohort selection dataset Stubbs et al. (2019). Further annotation guidelines were provided for two criteria.

The term "major complication" for MAJOR-DIABETES was defined as any of the following that are a result of (or strongly correlated with) uncontrolled diabetes: amputation, kidney damage, skin conditions, retinopathy, nephropathy, neuropathy.

The term "advanced" in ADVANCED-CAD was defined as having 2 or more of the following: Taking 2 or more medications to treat CAD; History of myocardial infarction (MI); Currently experiencing angina; Ischemia, past or present.

Criteria Label	Definition		
DRUG-ABUSE	Drug abuse, current or past		
ALCOHOL-ABUSE	Current alcohol use over weekly recommended limits		
ENGLISH	Patient must speak English		
MAKES-DECISIONS	Patient must make their own medical decisions		
ABDOMINAL	History of intra-abdominal surgery, small or large intestine resection, or small bowel obstruction		
MAJOR-DIABETES	Major diabetes-related complication.		
ADVANCED-CAD	Advanced cardiovascular disease (CAD).		
MI-6MOS	MI in the past 6 months		
KETO-1YR	Diagnosis of ketoacidosis in the past year		
DIETSUPP-2MOS	Taken a dietary supplement (excluding vitamin D) in the past 2 months		
ASP-FOR-MI	Use of aspirin to prevent MI		
HBA1C	Any hemoglobin A1c (HbA1c) value be-		
CDE ATINIDA	tween 6.5% and 9.5%		
CREATININE	Serum creatinine > upper limit of normal		

Table 6: Criteria labels and their definitions for he N2C2 cohort selection dataset.

Appendix B. Eligibility Criteria Annotation

In Table 7, we annotate the 2018 N2C2 cohort selection eligibility criteria using the Chia Annotation Model (Kury et al., 2020). Additionally, we label an entity as expressed directly (D) or implicitly (I).

LEON TRAMONTINI GHOSH EICKHOFF

Cuitania I ab al		Entities			Relations
Criteria Label	Item	Entity (Text)	Expression	Item	Relation (arg1, arg2)
DRUG-ABUSE	T1 T2	Condition (drug abuse) Temporal (current or past)	D I/D	R1	has_temporal (T1, T2)
ALCOHOL-ABUSE	T1 T2 T3	Condition (alcohol use) Temporal (Current) Qualifier (over weekly recommended limits)	D I/D I	R1 R2	has_temporal (T1, T2) has_qualifier (T1, T3)
ENGLISH	T1	Observation (speak English)	I/D		
MAKES-DECISIONS	T1	Observation (make their own medical decision)	I		
ABDOMINAL	T1 T2 T3 T4 T5	Observation (History of) Procedure (intra-abdominal surgery) Procedure (small or large intestine resection) Condition (small bowel obstruction) Scope (T3, T4)	I/D I D D	R1 * R2	has.temporal (T2, T1) or (T3, T4) subsumes (T2, T5)
MAJOR-DIABETES	T1 T2 T3	Observation (complication) Qualifier (major) Qualifier (diabetes-related)	D I I	R1 R2	has_qualifier (T1, T2) has_qualifier (T1, T3)
ADVANCED-CAD	T1 T2	Condition (cardiovascular disease (CAD)) Qualifier (advanced)	D I	R1	has_qualifier(T1, T2)
MI-6MOS	T1 T2	Condition (MI) Temporal (past 6 months)	D I	R1	has_temporal (T1, T2)
KETO-1YR	T1 T2 T3	Condition (ketoacidosis) Temporal (past year) Context (diagnosis)	D I D	R1 R2	has_temporal (T1, T2) has_context (T1, T3)
DIETSUPP-2MOS	T1 T2 T3 T4	Observation (dietary supplement) Temporal (past 2 months) Context (vitamin D) Negation (excluding)	I I D	R1 R2 R3	has_negation (T3, T4) has_temporal (T1, T2) has_context (T1, T3)
ASP-FOR-MI	T1 T2 T3 T4	Drug (Aspirin) Condition (MI) Context (prevent) Scope (to prevent MI)	D D I/D	R1 R2	has_context (T2, T3) has_scope (T1, T4)
HBA1C	T1 T2 T3	Measurement (hemoglobin A1c (HbA1c)) Value (value between 6.5% to 9.5%) Qualifier (Any)	D D I	R1 R2	has_value (T1, T2) has_qualifier (T1, T3)
CREATININE	T1 T2	Measurement (Serum creatinine) Value (> upper limit of normal)	D I	R1	has_value (T1, T2)

Table 7: Entity and relation annotations for N2C2 eligibility criteria according to the Chia Annotation Model. We additionally label an entity as expressed directly (D) or implicitly (I).