

Are large-language models enough to generate questionnaires? A use-case with clinical trials

Anonymous ACL submission

Abstract

We propose DivAC-LLM, a hybrid pipeline that automatically generates physician-facing questionnaires from clinical trial eligibility criteria. The pipeline integrates rule-based systems, shallow machine learning models, fine-tuned biomedical transformers, and prompt-engineered Large Language Models (LLMs). Unlike purely generative approaches, our system explicitly predicts the number of questions derived from each criterion and injects structured cues into LLM prompts to improve control and accuracy. We validate the approach on two datasets: CTQG, a corpus we created from cancer trials, and QG-SQuAD, a modified version of SQuAD for question generation. Experimental results show that combining symbolic and neural components reduces hallucinations and improves semantic alignment. The results highlight that LLMs benefit from lightweight structure and explicit intermediate reasoning when generating interpretable clinical questionnaires.

1 Introduction

Since the mid-20th century, clinical trials have been the gold standard for evaluating the safety and efficacy of medical interventions (Weed, 1968). Randomized Controlled Trials (RCTs) have guided drug approvals, treatment guidelines, and healthcare policies (Meinert, 2012). However, a bottleneck in this process is eligibility pre-screening (Brøgger-Mikkelsen et al., 2020). Although inclusion and exclusion criteria are designed to identify suitable patient subgroups, they are typically expressed in unstructured free text. As a result, these criteria are often difficult to interpret, even for clinicians, and pose significant challenges to automatic parsing and operationalization (Raghavan et al., 2014).

With the widespread adoption of Electronic Health Records (EHRs) in the 2000s (Blumenthal and Tavenner, 2010), and the public availability

of over 450,000 trials on ClinicalTrials.gov¹, the vision of automated patient-trial matching became increasingly plausible. However, the gap between free-text criteria and structured patient data remains unresolved (Kaskovich et al., 2023).

One intermediate solution is the use of structured questionnaires automatically derived from eligibility criteria. These questionnaires, intended for physicians, can be completed semi-automatically using EHR data (Figure 1). They offer interpretability, partial standardization, and actionable guidance—enabling healthcare providers to assess trial eligibility in a more streamlined and reproducible manner.

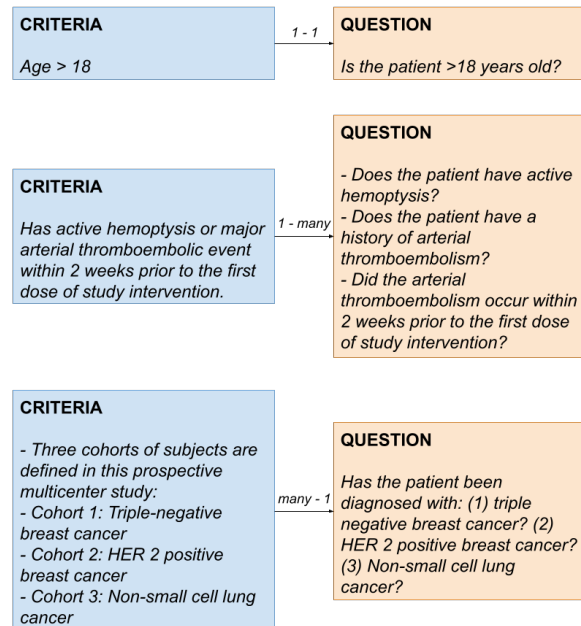


Figure 1: Examples of question generation.

Prior literature has focused on query generation rather than on physician-oriented questionnaire generation (Park et al., 2024a). However, recent advances in LLMs enable prompt-based rephrasing

¹<https://ClinicalTrials.gov>

of eligibility criteria into questions, although LLMs alone remain insufficient (Smith et al., 2023). Our hypothesis is that LLMs need to be explicitly informed about the criteria values and the number of questions that a criterion generates, to improve their ability to create a questionnaire from a clinical trial.

To explore this idea, we propose a hybrid pipeline that identifies criteria in trial descriptions, extracts the clinical attributes, and predicts whether a criterion should yield one or multiple questions before prompting an LLM. This produces questionnaires where clinicians fill in patient information rather than interpreting raw text.

We also create a dataset, called CTQG (Clinical Trial Question Generation), extracted from lung cancer clinical trials from ClinicalTrials.gov written exclusively in English, built to train models to translate criteria into questions.

We evaluate our approach on an annotated corpus of real-world clinical trial criteria. In order to compare it with existing methods, we also apply our approach to the SQuAD dataset (Rajpurkar et al., 2018).

The contributions² of this paper are:

- An overview of related work on patient pre-screening and automated processing of clinical trial criteria.
- A reproducible and operational pipeline for transforming clinical trial criteria into physician-facing questionnaires.
- A set of guidelines for the creation of the CTQG dataset and the public release of the QG-SQuAD dataset.

2 Related Work

Regarding clinical trial eligibility criteria processing, early efforts focused on rule-based extraction pipelines such as EliIE (Kang et al., 2017) and Criteria2Query (Yuan et al., 2019), as well as annotated corpora like Chia (Kury et al., 2020), which transform unstructured criteria into structured representations aligned with database queries.

Text mining from clinical trial documents has also enabled alignment between eligibility criteria and structured patient data (Weng et al., 2011). This includes entity normalization, logical form

parsing and OMOP-compatible cohort definitions (Ziletti and D’Ambrosi, 2025). The Leaf Clinical Trials corpus (Dobbins et al., 2022) has provided a large-scale semantic benchmark for such tasks. More recently, LeafAI (Dobbins et al., 2023) introduced neural query generators that achieve retrieval performance comparable to experienced programmers.

Some approaches convert criteria into interactive questions to assist in patient recruitment. DQueST (Liu et al., 2019) generates dynamic "Yes/No" questions to match patients with trials. However, it lacks domain-specific clinical phrasing and does not integrate EHR signals.

LLMs have opened new possibilities for clinical trial matching. (Wornow et al., 2025) apply zero-shot prompting to evaluate trial-patient compatibility and generate natural language justifications. TrialGPT (Jin et al., 2024) fine-tunes generative models for trial matching across thousands of real-world protocols, while FollowupBench (Gatto et al., 2025) addresses clinical question generation from patient-provider dialogues. However, these systems aim for classification or summarization, not the generation of interpretable, structured questions. Besides, they often suffer from factual hallucinations, constraint omission, and semantic overgeneralization (García-Barragán et al., 2025; Pal et al., 2022b).

Hybrid architectures have emerged to address these shortcomings by combining symbolic, statistical, and neural methods. (Li et al., 2024) propose a post-processing framework to reduce hallucinations in medical QA. (Yi et al., 2024) incorporate entity-based pre-filtering before LLM reasoning in trial-patient matching. Finally, (Shi and Bucher, 2025) combine prompt engineering with rule-based parsing for better generalization across eligibility criteria. It is the closest work to ours.

Recent works unleashed the power of LLMs to enhance patient-trial matching. Criteria2Query3 (Park et al., 2024b) used GPT-4 to transform the text from ClinicalTrials.gov into SQL queries. Finally, the PRISM framework (Gupta et al., 2024) fine-tuned an LLM to both simplify eligibility criteria and realize patient prescreening, but the question generation module does not consider that a criteria can become several questions. However, it is the closest work to ours.

Despite recent progress, no work has tackled structured question generation from clinical trial criteria for physician-facing tools. There is no stan-

²The clinical trial codes and the final prompts will be available with the final version of this paper

159 dard protocol or template to transform criteria into
 160 questions: identical meanings can yield distinct
 161 formulations across trials.

162 3 Methodology

163 3.1 Data

164 3.1.1 CTQG dataset

165 The main dataset used in this work has been de-
 166 signed based on a set of clinical trials obtained
 167 from the ClinicalTrials.gov platform. This is an
 168 online database, administered by the US National
 169 Institutes of Health (NIH), which contains infor-
 170 mation on more than 440,000 clinical studies from
 171 220 countries covering a wide range of diseases
 172 and conditions.

173 These criteria were then refined with medical
 174 expert input to establish an explicit correspondence
 175 between each criterion and the questions required
 176 to operationalize it. An oncologist generated and
 177 curated 729 questions using ChatGPT-4o, ensuring
 178 normalization so that identical criteria across trials
 179 map to the same question. The main prompting
 180 setup is provided in Appendix A (Section A.1.1).
 181 Criteria were presented in blocks to obtain initial
 182 question drafts, which were edited and standard-
 183 ized. CTQG is designed as a high-precision expert-
 184 normalized benchmark, prioritizing semantic con-
 185 sistency over scale. The details of the questions
 186 creation is available at appendix B.

187 Some prompts we crafted for ChatGPT are avail-
 188 able at Appendix A. No explicit rules were defined
 189 to sort criteria according to how many questions
 190 they generate.

191 However, the correspondence between criterion
 192 and question is not strictly one-to-one, but can be
 193 also one-to-many or many-to-one (see Figure 1).
 194 Data have been processed representing this infor-
 195 mation.

196 3.1.2 SQuAD dataset

197 The second dataset employed in this study is the
 198 Stanford Question Answering Dataset (SQuAD)
 199 (Rajpurkar et al., 2018). Released in 2016, SQuAD
 200 1.1 was designed as a realistic benchmark for
 201 machine-reading comprehension. It consists of
 202 Wikipedia passages paired with human-written
 203 questions whose answers correspond to contigu-
 204 ous spans in the text, totalling over 100,000 ques-
 205 tion-answer pairs.

206 For this work, SQuAD was transformed to mir-
 207 ror the structure of the CTQG dataset, yielding

101,888 entries. Each passage was split into sen-
 208 tences, their character offsets were recorded, and
 209 each question was linked to the sentence contain-
 210 ing its answer span (Algorithm 1). We used the
 211 NLTK toolkit (Bird and Loper, 2004) for sentence
 212 segmentation. 213

Algorithm 1 Map questions to sentence intervals

Require: Paragraphs \mathcal{P} ; contexts \mathcal{C}_p per paragraph; for each context c : QA
 pairs $\mathcal{A}_c = \{(a_{\text{start}}^{(k)}, q^{(k)})\}_k$

```

1: for all  $p \in \mathcal{P}$  do
2:   for all  $c \in \mathcal{C}_p$  do
3:      $[s_1, \dots, s_L] \leftarrow \text{SPLITINTOSENTENCES}(c)$ 
4:      $o \leftarrow 0$ ;
5:     for  $i = 1$  to  $L$  do
6:        $\ell_i \leftarrow \text{INDEXOF}(c, s_i, o)$ ;  $r_i \leftarrow \ell_i + |s_i|$ ;  $o \leftarrow r_i$ 
7:     end for
8:     for all  $(a_{\text{start}}, q) \in \mathcal{A}_c$  do
9:       find  $j$  s.t.  $\ell_j \leq a_{\text{start}} < r_j$  ▷ binary search or scan
10:      map  $q \rightarrow (s_j, [\ell_j, r_j])$ 
11:    end for
12:  end for
13: end for
```

214 Furthermore, correspondences were identified
 215 between each original text segment (or “Criterion”)
 216 and the number of questions generated. The distri-
 217 bution of the number of questions associated with
 218 each criterion is as follows: 40,961 (40.2%) with
 219 zero questions, 40,467 (39.72%) with one ques-
 220 tion, 14,536 (14.27%) with two questions, 3,985
 221 (3.91%) with three questions, 1,251 (1.23%) with
 222 four questions, and 688 (0.67%) with five questions.
 223 This distribution highlights that, unlike CTQG, this
 224 dataset includes many cases in which no questions
 225 are generated. QG-SQuAD reflects a positional
 226 task: most sentences contain zero or one answer-
 227 able segment, and the number of questions depends
 228 on span alignment. In other words, QG-SQuAD
 229 reflects positional span alignment, not conceptual
 230 decomposition. As a result, QG-SQuAD offers
 231 scale and robustness tests but evaluate another skill:
 232 local span-driven generation rather than domain-
 233 level reasoning. 233

234 Both processed datasets were partitioned into
 235 training, validation, and test subsets, whose sizes
 236 are given in Table 1. 236

237 3.1.3 Comparative of datasets

238 Throughout this section, both datasets have been
 239 compared from different perspectives, both glob-
 240 ally and with respect to the train, test, and valida-
 241 tion partitions. These measures are presented in
 242 Table 1. 242

243 As can be seen, there is homogeneity within the
 244 partitions of each dataset, although certain differ-
 245 ences can indeed be detected between them. For ex-
 246 ample, SQuAD texts are longer in terms of the num- 246

| Dataset | Partition | Size | Vocab. size | Entropy | Characters | | | | Words | | | | Sentences | | | |
|---------|-------------------|---------|-------------|---------|------------|-------|--------|--------|-------|-----|-------|-------|-----------|-----|------|------|
| | | | | | Min | Max | Avg | StD | Min | Max | Avg | StD | Min | Max | Avg | StD |
| CTQG | <i>Train</i> | 473 | 1,879 | 9.19 | 5 | 602 | 126.62 | 99.94 | 1 | 76 | 18.77 | 14.91 | 1 | 4 | 1.15 | 0.43 |
| | <i>Test</i> | 148 | 925 | 8.68 | 9 | 742 | 125.40 | 111.04 | 1 | 111 | 18.90 | 16.86 | 1 | 4 | 1.15 | 0.46 |
| | <i>Validation</i> | 108 | 818 | 8.63 | 13 | 573 | 130.55 | 125.63 | 3 | 93 | 19.43 | 19.02 | 1 | 4 | 1.15 | 0.49 |
| | <i>Total</i> | 729 | 2,244 | 9.25 | 5 | 742 | 126.95 | 106.24 | 1 | 111 | 18.90 | 15.96 | 1 | 4 | 1.15 | 0.44 |
| SQuAD | <i>Train</i> | 81,252 | 91,287 | 11.27 | 1 | 2,458 | 150.93 | 78.67 | 0 | 341 | 24.37 | 12.54 | 0 | 12 | 1.04 | 0.25 |
| | <i>Test</i> | 10,407 | 26,644 | 10.95 | 2 | 781 | 150.39 | 79.04 | 1 | 127 | 24.29 | 12.55 | 0 | 4 | 1.04 | 0.23 |
| | <i>Validation</i> | 10,229 | 26,877 | 10.93 | 1 | 1,412 | 155.68 | 85.23 | 0 | 230 | 25.13 | 13.58 | 0 | 11 | 1.05 | 0.31 |
| | <i>Total</i> | 101,888 | 103,675 | 11.30 | 1 | 2,458 | 151.35 | 79.40 | 0 | 341 | 24.44 | 12.65 | 0 | 12 | 1.04 | 0.26 |

Table 1: Statistical summary of the datasets used (CTQG, QG-SQuAD) with respect to the different partitions (train, test, validation, total), including values for size, vocabulary size, Shannon entropy, and the minimum, maximum, mean, and standard deviation of the number of characters, words and sentences. Note that for the SQuAD dataset, a sample of the total was taken for the development of this research.

ber of characters and words, although they contain fewer sentences than CTQG. Similarly, SQuAD shows a broader vocabulary as well as higher entropy, which makes sense given the size differences between the two datasets.

3.2 DivAC-LLM pipeline

To tackle the challenges mentioned in Section 2, our pipeline is divided into two tasks: predicting the number of questions to generate per criterion and the generation itself with criteria extraction. Figure 2 illustrates our pipeline and the next sections explain the different modules.

3.2.1 Question division prediction

This task can either be seen as a regression (approximate the real number of questions to be generated from a criterion) or a classification (categorize the number of questions to be generated). In both cases, it is a supervised learning approach, as we know how many questions were generated for each criterion. To ease the task, we decided to treat it as a ternary classification problem: 0 (no question to be generated), 1 (one question to be generated) and 2 (several questions to be generated).

As it is a predictive task, LLMs are not as strong as most classic fine-tuning approaches (Trust and Minghim, 2024). In addition, since criteria are a short sentence (18.9 words on average), simple solutions can be effective, thus reducing computing time and energy consumptions. We created a voting algorithm in three steps, giving more weight to precision than recall, arguing that it is best that the system ignores how many questions have to be generated rather than guessing a number:

1. We extracted dozens of features from the criterion with SpaCy (Honnibal et al., 2020) (whose list is available in Table 15 at Appendix E) and established correlation pat-

terns between the feature values and the three classes in order to implement a rule-based system. To complete it, we implemented a Random Forest algorithm and we predicted the remaining classes that had not been labeled by the rule-based system, that has great precision but low recall.

2. We fine-tuned two pre-trained transformers specialized in the bio-medical domain (BioBERT (Lee et al., 2019) and BioMedLM (Bolton et al., 2024)) to predict the right class for each criterion.
3. We crafted a voting algorithm with several conditions, where in case of disagreement between the rule-based + Random Forest and the different fine-tuned models, the priorities could be given to a different predictor, or to the "I don't know" class. The description of the different strategies is given in Table 16 in Appendix E.

We chose the voting configuration that maximizes precision with acceptable recall.

The training process is illustrated in Algorithm 2 and the voting system in Algorithm 3, both in Appendix C.

3.2.2 Question generation and criteria management

Two approaches are state-of-the-art, as seen in Section 2: supervised approaches with transformers fine-tuning and zero-shot learning with LLMs. They are end-to-end methods that take as input a criterion in natural language and transform it into questions. The fine-tuning approach can be considered as a machine translation task (Raffel et al., 2020), where the criterion is the source language and the questions the target. The advantage of this

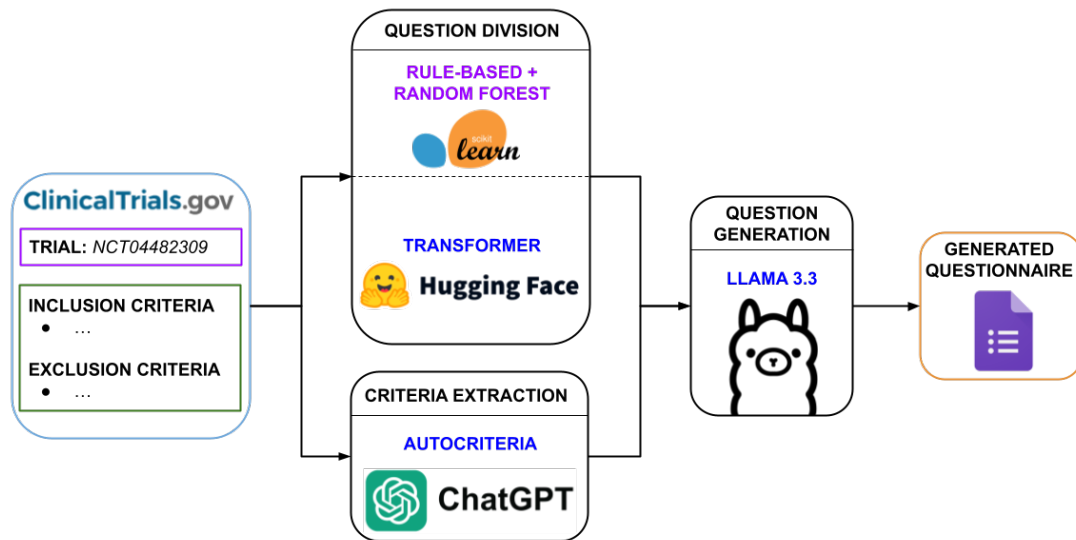


Figure 2: Automatic questionnaire generation pipeline.

technique is that the model learns from real question examples and thus is accurate when "translating". However, in order to work, it requires a parallel corpus of a minimum size (Koehn and Knowles, 2017) in order to be effective. On the other side, the LLM approach requires a few or even no question examples in order to work, due to its huge language model and generative capacities, but the issues are in the risks of hallucinations and variability in its generation.

Nevertheless, as the LLMs approach is based on prompt engineering and instructions crafting (Liu et al., 2023), it is possible to control the output by adding to the prompt information about the criteria. In other words, the LLMs approach can be stronger if used with a previous step where the criteria values are extracted in a NER fashion and passed explicitly to the prompt. We decided to combine two LLMs:

1. One is the state-of-the-art in criteria extraction called AutoCriteria (Datta et al., 2024), a framework based on the GPT-4 model to extract granular eligibility criteria information from clinical trial documents.
2. Another is an instance of Llama 3.3 (Grattafiori et al., 2024), an open-source LLM, easier to fine-tune than commercial ones, that we used in order to generate more accurate questions for each criterion with less variability among clinical trials with a few-shot learning approach where we feed the generic LLM with examples of criteria-questions pairs. The

prompts that have been used are presented in Appendix A.1.

The final result is a three-module system (question division, AutoCriteria and LLM) that we called DivAC-LLM, optimized for questionnaire generation from clinical trials.

4 Experimental results and evaluation

4.1 Experimental setup

Regarding the baselines, we used, depending on the dataset:

1. **CTQG**: three state-of-the-art pre-trained transformers, two specialized in the medical domain (BioBART (Yuan et al., 2022), BERT2BERT (Moll et al., 2025)), one generic (Flan-t5 (Smilga and Alabiad, 2024)), ChatGPT-4o and Llama 3.3 against the full DivAC-LLM.
2. **QG-SQuAD**: two state-of-the-art generic pre-trained transformers (Prophetnet (Qi et al., 2020), Pegasus-Xsum (Zhang et al., 2019)), ChatGPT-4o and Llama3.3 against DivAC-LLM.

For the experiments, the hardware, training time and costs are available in Tables 18 and 19 in Appendix E. When fine-tuning transformers models for question generation, we performed a grid search, where the ranges are available at Table 21 in Appendix E. The best hyperparameters for each dataset are available at Table 17 in Appendix E.

| Dataset & class | Precision | Recall | F-score |
|-----------------|-------------|-------------|-------------|
| CTQG 0 | 0.74 | 0.33 | 0.45 |
| CTQG 1 | 0.84 | 0.71 | 0.77 |
| CTQG 2 | 0.33 | 0.43 | 0.38 |
| QG-SQuAD 0 | 0.69 | 0.41 | 0.52 |
| QG-SQuAD 1 | 0.38 | 0.41 | 0.40 |
| QG-SQuAD 2 | 0.45 | 0.68 | 0.54 |

Table 2: Results of question division prediction with the best voting strategy for both CTQG and QG-SQuAD.

4.2 Results evaluation

For question division, we realized a classic automatic evaluation. For question generation, we conducted two distinct evaluations, a human evaluation based on (Nema and Khapra, 2018) and an automatic evaluation based on the QGEval metrics benchmark (Fu et al., 2024).

4.2.1 Question division evaluation

As it is a module that can be activated or deactivated, we evaluate separately the question division prediction on both CTQG and QG-SQuAD. We show in Table 2 the results with the best voting strategy. While the question division module underperforms on QG-SQuAD, due to this module reliance on linguistic heuristics rather than domain-informed labeling, the module achieves better results on CTQG, in the precision of classes 0 and 1. In the Section 4.2.3 we observe the impact of the module for questionnaire generation.

4.2.2 Human evaluation

We compared the results of the benchmarked models with ours (DivAC-LLM) using a human evaluation. We used three metrics with three annotators (3 Spanish males between 30 and 40, 2 NLP engineers, 1 domain expert) to evaluate generated questions quality: *fluency* (how fluent is the generated question), *adequacy* (how close are the reference and the generated question) and *understandability* (how understandable is the generated question) (Callison-Burch et al., 2006). Annotators followed a shared guideline defining adequacy, fluency and understandability, with ties allowed. For the three of them, a scale from 1 to 5 is asked for, where the higher the value, the better the text is with respect to such metric. These evaluations have been made individually, and no information about the model has been provided (blinded revision).

We merged the mean of each annotator for each model and each metric and summarized it in the

Table 3. The results show that for fluency and understandability, the metrics that depend on the quality of the generation, the gold standard and the DivAC-LLM have similar scores, whereas DivAC-LLM matches or slightly exceeds reference adequacy, reflecting normalization consistency rather than superior clinical judgment.

Additionally, two bio-medical experts (1 male, 1 female both Spanish and between 30 and 40) were in charge of detecting which questions were generated and which were the original reference questions. In Table 4 we can observe that they are unable to detect which of the questions belong to the reference standard and which to Llama 3.3, proving how hard it is to detect generated content for a non-computing linguist. Finally, these two annotators were asked to choose between DivAC-LLM and their own ground truth. In 69.8% of the cases, they chose DivAC-LLM.

Three annotators (two computational linguists and one biomedical expert) rated all generated questions for adequacy, understanding, and fluency. Table 5 reports standard inter-annotator agreement (Fleiss’ κ , Krippendorff’s α , Scott’s π) (Artstein and Poesio, 2008) together with a Wilcoxon dominance score (W) (Dror et al., 2018). The W score summarises all pairwise one-sided Wilcoxon tests as the number of significant wins minus losses for each model, providing a single comparative indicator that is not affected by chance-correction artifacts.

Agreement values are heterogeneous across systems: encoder-decoder baselines show moderate agreement, whereas LLMs display low or negative $\kappa/\alpha/\pi$, reflecting the instability of chance-corrected measures under skewed or near-ceiling rating distributions. The Wilcoxon scores offer a clearer picture: DivAC-LLM and ChatGPT obtain the highest W values across all three dimensions, while Pegasus underperforms with negative scores. Mid-range models (BioBART, BERT2BERT, FLAN-T5, ChatGPT-lite, Llama3.3) show small positive or negative values, indicating mixed performance.

4.2.3 QGEval evaluation

To validate DivAC-LLM, instead of relying on traditional machine translation metrics such as BLEU (Papineni et al., 2001), ROUGE (Lin, 2004), or BERT-score (Zhang et al., 2020), we used the QGEval benchmark. This open source framework brings together more than a dozen evaluation met-

| Model | Annotator 1 | | | Annotator 2 | | | Annotator 3 | | | Mean | | |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Adeq. | Underst. | Flu. | Adeq. | Underst. | Flu. | Adeq. | Underst. | Flu. | Adeq. | Underst. | Flu. |
| biobart | 4.01 | 4.48 | 4.47 | 3.87 | 4.47 | 4.47 | 3.81 | 4.25 | 4.14 | 3.90 | 4.40 | 4.36 |
| BERT2BERT | 3.44 | 4.18 | 4.15 | 3.32 | 4.17 | 4.21 | 3.40 | 3.74 | 3.72 | 3.38 | 4.03 | 4.03 |
| FLAN-T5 | 3.83 | 4.77 | 4.68 | 3.60 | 4.70 | 4.70 | 3.93 | 4.50 | 3.88 | 3.79 | 4.66 | 4.42 |
| pegasus-xsum | 2.86 | 3.46 | 3.41 | 2.76 | 3.45 | 3.46 | 2.50 | 2.99 | 2.40 | 2.71 | 3.30 | 3.09 |
| ChatGPT | 4.28 | 4.99 | 4.99 | 3.83 | 5.00 | 5.00 | 4.67 | 4.95 | 4.88 | 4.26 | 4.98 | 4.96 |
| ChatGPT-lite | 4.59 | 4.99 | 4.97 | 4.18 | 5.00 | 5.00 | 4.91 | 5.00 | 4.95 | 4.56 | 5.00 | 4.97 |
| Llama3.3 | 4.28 | 5.00 | 5.00 | 3.99 | 5.00 | 5.00 | 4.58 | 4.68 | 4.60 | 4.28 | 4.89 | 4.87 |
| DivAC-LLM | 4.71 | 4.98 | 4.92 | 4.87 | 4.94 | 5.00 | 4.93 | 4.98 | 4.94 | 4.84 | 4.97 | 4.95 |
| human | 4.54 | 4.98 | 4.94 | 4.83 | 4.97 | 4.98 | 4.74 | 4.99 | 5.00 | 4.70 | 4.98 | 4.97 |

Table 3: Human evaluation scores for adequacy, understanding, and fluency from three annotators (Annotator 1-3) and their mean.

| Model | Percentage (%) |
|-----------------|----------------|
| FLAN-T5 | 6.4 |
| ChatGPT | 11.1 |
| BERT2BERT | 13.4 |
| bioBART | 13.4 |
| ChatGPT lite | 14.6 |
| Llama3.3 | 20.5 |
| Manual labeling | 20.5 |

Table 4: Percentage of times each model was selected as the best system.

rics, ranging from the aforementioned ones to others such as QRelScore (Wang et al., 2022) (which combines word-level matching with sentence-level generation) and UniEval (Li et al., 2025) (a set of fine-grained evaluation tags designed to better correlate with human judgments). A detailed description of these metrics is provided in Tables 13 and 14 in Appendix D. Table 6 shows the results of the metrics for the CTQG dataset. Our method outperforms all the other baselines in both lexical and semantic metrics. We also made an ablation study by removing the question division and then AutoCriteria modules. It is interesting to notice that a fine-tuned Llama 3.3 works best with both modules and then with none, so we hypothesize there could be some dependency between them. To confirm the importance of question division for this task, we also compared our method against an oracle, where the LLM knows the real number of questions to be generated. The oracle outperforms all the other configurations but the difference is not as big with DivAC-LLM, showing that our module is well trained.

Regarding QG-SQuAD, as there are no criteria, we combined the question division module with different LLMs and compared the results with fine-tuned transformers. Table 7 shows the results for the QG-SQuAD dataset. In this case, the question division module improves LLMs models, but

the fine-tuned models outperform the LLMs for ROUGE and, more surprisingly, for BERTscore, proving that the question division is critical for clinical trials, but not generalizable for all question generation tasks.

5 Discussion

Our results support three observations. First, conditioning prompts with predicted question counts and extracted criterion values improves CTQG generation across lexical and semantic metrics, and better matches expert phrasing than bare LLM prompting. Second, separating division, extraction, and generation yields more stable outputs than end-to-end prompting, especially for multi-clause criteria. Third, precision-oriented division reduces multi-question splits and hallucinations.

The comparison with QG-SQuAD clarifies scope. Because its “criteria” are generic sentences and alignment is position-based, division accuracy has limited semantic impact, and fine-tuned seq2seq baselines remain competitive on ROUGE and BERTScore. In contrast, CTQG benefits from explicit domain cues (entities, thresholds, exception lists), which allow rule- and feature-based models to anchor generation.

Human evaluation mirrors these trends: fluency and understandability vary little across strong systems, while adequacy is more sensitive to prompt conditioning and explicit metadata. Inter-annotator agreement is moderate, consistent with prior subjective NLG evaluation. While the questionnaire format improves clarity, errors may still affect eligibility decisions. DivAC-LLM mitigates this risk via rule-guided prompting and explicit metadata, but remains a decision-support tool requiring clinical supervision.

Finally, the pipeline ablation suggests that LLMs

| Model | Adequacy | | | | Understanding | | | | Fluency | | | |
|--------------|-----------------|-----------------|---------------|-----------|-----------------|-----------------|---------------|-----------|-----------------|-----------------|---------------|-----------|
| | Fleiss κ | Kripp. α | Scott's π | W | Fleiss κ | Kripp. α | Scott's π | W | Fleiss κ | Kripp. α | Scott's π | W |
| biobart | 0.326 | 0.318 | 0.316 | -1 | 0.503 | 0.498 | 0.496 | +1 | 0.471 | 0.466 | 0.463 | +1 |
| BERT2BERT | 0.333 | 0.321 | 0.318 | -1 | 0.457 | 0.447 | 0.445 | +1 | 0.465 | 0.455 | 0.453 | +1 |
| FLAN-T5 | 0.255 | 0.245 | 0.242 | -1 | 0.384 | 0.381 | 0.378 | +1 | 0.156 | 0.134 | 0.130 | +1 |
| pegasus-xsum | 0.308 | 0.281 | 0.278 | -7 | 0.514 | 0.488 | 0.486 | -7 | 0.355 | 0.323 | 0.321 | -7 |
| ChatGPT | 0.199 | 0.163 | 0.160 | +4 | 0.163 | 0.162 | 0.159 | +4 | 0.055 | 0.045 | 0.041 | +5 |
| ChatGPT-lite | 0.092 | 0.043 | 0.039 | -1 | 1.000 | 1.000 | 1.000 | +1 | 0.192 | 0.191 | 0.188 | +1 |
| Llama3.3 | 0.225 | 0.202 | 0.199 | -1 | 1.000 | -0.017 | -0.021 | +1 | 1.000 | -0.022 | -0.026 | +1 |
| human | 0.220 | 0.214 | 0.211 | NA | 0.149 | 0.150 | 0.147 | NA | 0.071 | 0.068 | 0.065 | NA |
| DivAC-LLM | 0.015 | -0.299 | -0.304 | +6 | -0.001 | -0.434 | -0.439 | +4 | 1.000 | -0.434 | -0.439 | +5 |

Table 5: Inter-annotator agreement (Fleiss’ κ , Krippendorff’s α , Scott’s π) and Wilcoxon signed-rank dominance score (W). W is the number of significant wins minus losses for each model against all others for that dimension.

| metrics | QRelScore | BERTScore | ROUGE-L | Q-BLEU4 |
|---|--------------|--------------|--------------|--------------|
| bioBART | 0.363 | 0.772 | 0.721 | 0.462 |
| BERT2BERT | 0.393 | 0.696 | 0.438 | 0.523 |
| Flan-t5 | 0.658 | 0.710 | 0.471 | 0.443 |
| ChatGPT | 0.658 | 0.710 | 0.471 | 0.443 |
| Llama3.3 | 0.670 | 0.709 | 0.475 | 0.456 |
| Autocriteria + Llama3.3 | 0.745 | 0.891 | 0.399 | 0.395 |
| Question division + Llama3.3 | 0.742 | 0.893 | 0.408 | 0.410 |
| Fine tuned Llama3.3 | 0.710 | 0.944 | 0.694 | 0.618 |
| DivAC-LLM | 0.746 | 0.945 | 0.758 | 0.682 |
| DivAC-LLM oracle mode (no autocriteria) | 0.747 | 0.956 | 0.750 | 0.696 |
| divAC-LLM oracle mode | 0.752 | 0.953 | 0.722 | 0.653 |

Table 6: Automatic evaluation results on CTQG with respect to different models and combinations compared with DivAC-LLM.

| Model | BERTScore | BLEU | ROUGE |
|------------------|--------------|--------------|-------------|
| Gemini | 0.868 | 0.149 | 2.27 |
| Gemini + div | 0.855 | 0.209 | 2.98 |
| ChatGPT | 0.868 | 0.153 | 2.19 |
| ChatGPT + div | 0.874 | 0.210 | 3.07 |
| Llama3.3 | 0.869 | 0.177 | 2.37 |
| Llama3.3 + div | 0.884 | 0.170 | 3.17 |
| ProphetNet-large | 0.888 | 0.143 | 3.41 |
| Pegasus-XSum | 0.898 | 0.173 | 3.21 |

Table 7: Automatic evaluation results on QG-SQuAD comparing fine-tuned models and LLMs with question division module.

are strong surface realizers but weak controllers without explicit inputs. Lightweight predictors provide control signals that keep generation within clinically useful bounds.

6 Future work

In future work, we will replace the static voting mechanism with a learned calibration layer to improve per-criterion decision reliability, following recent work on expectation consistency calibration for neural networks (Clarté et al., 2023). We also plan to enrich criterion extraction with structured

representations and apply constrained decoding to enforce template-level consistency, extending recent findings in controlled text generation (Hokamp and Liu, 2017). Finally, we will integrate clinician feedback in a human-in-the-loop workflow so that their edits continuously fine-tune both division and generation modules, aligning with current advances in active learning for foundation models (Wan et al., 2023).

7 Conclusion

We introduced DivAC-LLM, a hybrid pipeline for generating physician-facing questionnaires from clinical trial criteria. By combining rule-based signals and classical models with LLM prompting, and by explicitly injecting extracted values, the system improves automatic metrics on CTQG and yields human judgments comparable to strong baselines for fluency and understandability while improving adequacy. The findings indicate that lightweight structure and prompt conditioning compensate for limitations of end-to-end prompting in clinical settings where multi-clause logic and exceptions are common.

571 Limitations

572 Certain limitations can be identified in this ap-
573 proach:

- 574 • **Dataset scope.** The dataset is limited to
575 English-language lung cancer trials, which
576 narrows both linguistic and clinical variabil-
577 ity. Broader multilingual evaluation is needed
578 to ensure robustness across registries and
579 pathologies (Neumann et al., 2024).
- 580 • **Human evaluation scale.** Expert assessment
581 was conducted by only three annotators within
582 one specialty, potentially biasing adequacy
583 and fluency judgments. Larger collaborative
584 annotation efforts such as MedMCQA illus-
585 trate the benefits of scaling expert input (Pal
586 et al., 2022a).
- 587 • **Model transparency.** Although DivAC-LLM
588 improves interpretability through explicit in-
589 termediate steps, it still depends on opaque
590 LLM components. Reliability and calibration
591 remain open challenges before any clinical
592 deployment.
- 593 • **External validation.** No downstream testing
594 has yet been performed on real patient screen-
595 ing tasks, leaving practical effectiveness and
596 safety unverified.
- 597 • **Energy consumption** Although the efforts
598 were optimized in order to avoid a huge en-
599 ergy consumption, the training and fine-tuning
600 periods were long enough to be costly and
601 energy-greedy, which has to be improved in
602 order to be more sustainable.
- 603 • **Why not fine-tuning an LLM to do every-**
604 **thing?** Given the size of the dataset and its
605 high-quality expert-based annotation, a few-
606 shot end-to-end supervised fine-tuning of a
607 small Llama could have been used to enrich
608 the benchmark. However, as we have very
609 short inputs (and outputs), we prioritize pre-
610 cision over recall (the "I always reply some-
611 thing" from the LLM can decrease precision).
612 Furthermore, the energy consumption would
613 be even bigger, so we decided to skip this
614 experiment.

References

- Ron Artstein and Massimo Poesio. 2008. [Survey article: Inter-coder agreement for computational linguistics](#). *Computational Linguistics*, 34(4):555–596. 616
617
618
- Steven Bird and Edward Loper. 2004. [NLTK: The natural language toolkit](#). In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics. 619
620
621
622
623
- David Blumenthal and Marilyn Tavenner. 2010. The "meaningful use" regulation for electronic health records. *N. Engl. J. Med.*, 363(6):501–504. 624
625
626
- Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, and Christopher D. Manning. 2024. [Biomedlm: A 2.7b parameter language model trained on biomedical text](#). *Preprint*, arXiv:2403.18421. 627
628
629
630
631
632
- Mette Brøgger-Mikkelsen, Zarqa Ali, John R Zibert, Anders Daniel Andersen, and Simon Francis Thomsen. 2020. Online patient recruitment in clinical trials: Systematic review and meta-analysis. *J. Med. Internet Res.*, 22(11):e22179. 633
634
635
636
637
- Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. [Re-evaluating the role of Bleu in machine translation research](#). In *11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 249–256, Trento, Italy. Association for Computational Linguistics. 638
639
640
641
642
643
- François Clarté, Yann Chevalayre, Guillaume Contardo, and Massimiliano Tommasi. 2023. [Expectation consistency for calibration of neural networks](#). In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, volume 216 of *Proceedings of Machine Learning Research*, pages 5431–5450. 644
645
646
647
648
649
- Surabhi Datta, Kyeryoung Lee, Hunki Paek, Frank J Manion, Nneka Ofoegbu, Jingcheng Du, Ying Li, Liang-Chin Huang, Jingqi Wang, Bin Lin, Hua Xu, and Xiaoyan Wang. 2024. [AutoCriteria: a generalizable clinical trial eligibility criteria extraction system powered by large language models](#). *J. Am. Med. Inform. Assoc.*, 31(2):375–385. 650
651
652
653
654
655
656
- Nicholas J Dobbins, Bin Han, Weipeng Zhou, Kristine F Lan, H Nina Kim, Robert Harrington, Özlem Uzuner, and Meliha Yetisgen. 2023. [LeafAI: query generator for clinical cohort discovery rivaling a human programmer](#). *J. Am. Med. Inform. Assoc.*, 30(12):1954–1964. 657
658
659
660
661
662
- Nicholas J. Dobbins, Tony Mullen, Özlem Uzuner, and Meliha Yetisgen. 2022. [The leaf clinical trials corpus: a new resource for query generation from clinical trial eligibility criteria](#). *Scientific Data*, 9(1):490. 663
664
665
666
- Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. [The hitchhiker’s guide to testing statistical significance in natural language processing](#). 667
668
669

| | | | |
|-----|---|--|-------------------|
| 670 | In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics. | | |
| 671 | | | |
| 672 | | | |
| 673 | | | |
| 674 | Weiping Fu, Bifan Wei, Jianxiang Hu, Zhongmin Cai, and Jun Liu. 2024. QGEval: Benchmarking multi-dimensional evaluation for question generation . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 11783–11803, Miami, Florida, USA. Association for Computational Linguistics. | | |
| 675 | | | |
| 676 | | | |
| 677 | | | |
| 678 | | | |
| 679 | | | |
| 680 | | | |
| 681 | Álvaro García-Barragán, Ahmad Sakor, Maria-Esther Vidal, Ernestina Menasalvas, Juan Cristobal Sanchez Gonzalez, Mariano Provencio, and Víctor Robles. 2025. Nssc: a neuro-symbolic ai system for enhancing accuracy of named entity recognition and linking from oncologic clinical notes . <i>Medical & Biological Engineering & Computing</i> , 63(3):749–772. | | |
| 682 | | | |
| 683 | | | |
| 684 | | | |
| 685 | | | |
| 686 | | | |
| 687 | | | |
| 688 | Joseph Gatto, Parker Seegmiller, Timothy Burdick, Inas S. Khayal, Sarah DeLozier, and Sarah M. Preum. 2025. Follow-up question generation for enhanced patient-provider conversations . <i>Preprint</i> , arXiv:2503.17509. | | |
| 689 | | | |
| 690 | | | |
| 691 | | | |
| 692 | | | |
| 693 | Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783. | | |
| 694 | | | |
| 695 | | | |
| 696 | | | |
| 697 | | | |
| 698 | | | |
| 699 | | | |
| 700 | | | |
| 701 | Shashi Gupta, Aditya Basu, Mauro Nievas, Jerrin Thomas, Nathan Wolfrath, Adhitya Ramamurthi, Bradley Taylor, Anai N Kothari, Regina Schwind, Therica M Miller, Sorena Nadaf-Rahrov, Yanshan Wang, and Hrituraj Singh. 2024. PRISM: Patient records interpretation for semantic clinical trial matching system using large language models . <i>NPJ Digit. Med.</i> , 7(1):305. | | |
| 702 | | | |
| 703 | | | |
| 704 | | | |
| 705 | | | |
| 706 | | | |
| 707 | | | |
| 708 | | | |
| 709 | Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics. | | |
| 710 | | | |
| 711 | | | |
| 712 | | | |
| 713 | | | |
| 714 | | | |
| 715 | | | |
| 716 | Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python . | | |
| 717 | | | |
| 718 | | | |
| 719 | Qiao Jin, Zifeng Wang, Charalampos S. Floudas, Fangyuan Chen, Changlin Gong, Dara Bracken-Clarke, Elisabetta Xue, Yifan Yang, Jimeng Sun, and Zhiyong Lu. 2024. Matching patients to clinical trials with large language models . <i>Nature Communications</i> , 15(1):9074. | | |
| 720 | | | |
| 721 | | | |
| 722 | | | |
| 723 | | | |
| 724 | | | |
| 725 | Tian Kang, Shaodian Zhang, Youlan Tang, Gregory W. Hruby, Alexander Rusanov, and Noémie Elhadad. | | |
| 726 | | | |
| | | 2017. EliIE: An open-source information extraction system for clinical trial eligibility criteria . <i>J. Am. Med. Inform. Assoc.</i> , 24(6):1062–1071. | 727 728 729 |
| | Samuel Kaskovich, Kirk D Wyatt, Tomasz Oliwa, Luca Graglia, Brian Furner, Jooho Lee, Anoop Mayampurath, and Samuel L Volchenbom. 2023. Automated matching of patients to clinical trials: A patient-centric natural language processing approach for pediatric leukemia . <i>JCO Clin. Cancer Inform.</i> , 7:e2300009. | 730 731 732 733 734 735 736 | |
| | Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation . In <i>Proceedings of the First Workshop on Neural Machine Translation</i> , pages 28–39, Vancouver. Association for Computational Linguistics. | 737 738 739 740 741 | |
| | Fabrcio Kury, Alex Butler, Chi Yuan, Li-heng Fu, Yingcheng Sun, Hao Liu, Ida Sim, Simona Carini, and Chunhua Weng. 2020. Chia, a large annotated corpus of clinical trial eligibility criteria . <i>Scientific Data</i> , 7(1):281. | 742 743 744 745 746 | |
| | Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining . <i>Bioinformatics</i> , 36(4):1234–1240. | 747 748 749 750 751 | |
| | Songda Li, Yunqi Zhang, Chunyuan Deng, Yake Niu, and Hui Zhao. 2024. Better late than never: Model-agnostic hallucination post-processing framework towards clinical text summarization . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 995–1011, Bangkok, Thailand. Association for Computational Linguistics. | 752 753 754 755 756 757 758 | |
| | Yi Li, Haonan Wang, Qixiang Zhang, Boyu Xiao, Chengchang Hu, Hualiang Wang, and Xiaomeng Li. 2025. Unieval: Unified holistic evaluation for unified multimodal understanding and generation . <i>Preprint</i> , arXiv:2505.10483. | 759 760 761 762 763 | |
| | Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . <i>Annu Meet Assoc Comput Linguistics</i> , pages 74–81. | 764 765 766 | |
| | Cong Liu, Chi Yuan, Alex M Butler, Richard D Carvajal, Ziran Ryan Li, Casey N Ta, and Chunhua Weng. 2019. DQueST: dynamic questionnaire for search of clinical trials . <i>J. Am. Med. Inform. Assoc.</i> , 26(11):1333–1343. | 767 768 769 770 771 | |
| | Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing . <i>ACM Comput. Surv.</i> , 55(9). | 772 773 774 775 776 | |
| | Curtis L Meinert. 2012. <i>ClinicalTrials: design, conduct and analysis</i> , volume 39. OUP USA. | 777 778 | |
| | Johannes Moll, Louisa Fay, Asfandyar Azhar, Sophie Ostmeier, Tim Lueth, Sergios Gatidis, Curtis Langlotz, and Jean-Benoit Delbrouck. 2025. Structuring | 779 780 781 | |

| | | |
|-----|--|---|
| 782 | radiology reports: Challenging llms with lightweight models. <i>arXiv preprint arXiv:2506.00200</i> . | |
| 783 | | |
| 784 | Preksha Nema and Mitesh M. Khapra. 2018. Towards a better metric for evaluating question generation systems. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3950–3959, Brussels, Belgium. Association for Computational Linguistics. | |
| 785 | | |
| 786 | | |
| 787 | | |
| 788 | | |
| 789 | | |
| 790 | Mark Neumann, Emily Alsentzer, Alistair Johnson, and Abhinav Singh. 2024. Medalign: Benchmarking foundation models for multilingual clinical text understanding. <i>Journal of Biomedical Informatics</i> , 152:104512. | |
| 791 | | |
| 792 | | |
| 793 | | |
| 794 | | |
| 795 | Ankit Pal, Logesh Umaphathi, and Malaikannan Sankarasubbu. 2022a. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. <i>NeurIPS Datasets and Benchmarks Track</i> . | |
| 796 | | |
| 797 | | |
| 798 | | |
| 799 | Ankit Pal, Logesh Kumar Umaphathi, and Malaikannan Sankarasubbu. 2022b. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In <i>Proceedings of the Conference on Health, Inference, and Learning</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260. PMLR. | |
| 800 | | |
| 801 | | |
| 802 | | |
| 803 | | |
| 804 | | |
| 805 | | |
| 806 | Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2001. BLEU: A method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02</i> , pages 311–318, Morristown, NJ, USA. Association for Computational Linguistics. | |
| 807 | | |
| 808 | | |
| 809 | | |
| 810 | | |
| 811 | | |
| 812 | | |
| 813 | Jimyung Park, Yilu Fang, Casey Ta, Gongbo Zhang, Betina Idnay, Fangyi Chen, David Feng, Rebecca Shyu, Emily R. Gordon, Matthew Spotnitz, and Chunhua Weng. 2024a. Criteria2query 3.0: Leveraging generative large language models for clinical trial eligibility query generation. <i>Journal of Biomedical Informatics</i> , 154:104649. | |
| 814 | | |
| 815 | | |
| 816 | | |
| 817 | | |
| 818 | | |
| 819 | | |
| 820 | Jimyung Park, Yilu Fang, Casey Ta, Gongbo Zhang, Betina Idnay, Fangyi Chen, David Feng, Rebecca Shyu, Emily R Gordon, Matthew Spotnitz, and Chunhua Weng. 2024b. Criteria2Query 3.0: Leveraging generative large language models for clinical trial eligibility query generation. <i>J. Biomed. Inform.</i> , 154(104649):104649. | |
| 821 | | |
| 822 | | |
| 823 | | |
| 824 | | |
| 825 | | |
| 826 | | |
| 827 | Weizhen Qi, Yu Yan, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and Ming Zhou. 2020. ProphetNet: Predicting future n-gram for sequence-to-SequencePre-training. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2401–2410, Online. Association for Computational Linguistics. | |
| 828 | | |
| 829 | | |
| 830 | | |
| 831 | | |
| 832 | | |
| 833 | | |
| 834 | Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1). | |
| 835 | | |
| 836 | | |
| 837 | | |
| 838 | | |
| | Preethi Raghavan, James L Chen, Eric Fosler-Lussier, and Albert M Lai. 2014. How essential are unstructured clinical narratives and information fusion to clinical trial recruitment? <i>AMIA Summits Transl. Sci. Proc.</i> , 2014:218–223. | 839 840 841 842 843 |
| | Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , Stroudsburg, PA, USA. Association for Computational Linguistics. | 844 845 846 847 848 849 |
| | Jianlin Shi and Brian T. Bucher. 2025. Initial investigation of llm-assisted development of rule-based clinical nlp system. <i>Preprint</i> , arXiv:2506.16628. | 850 851 852 |
| | Veronika Smilga and Hazem Alabiad. 2024. TüDuo at SemEval-2024 task 2: Flan-t5 and data augmentation for biomedical NLI. In <i>Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)</i> , pages 737–744, Mexico City, Mexico. Association for Computational Linguistics. | 853 854 855 856 857 858 |
| | Andrew L Smith, Felix Greaves, and Trishan Panch. 2023. Hallucination or confabulation? neuroanatomy as metaphor in large language models. <i>PLOS Digit. Health</i> , 2(11):e0000388. | 859 860 861 862 |
| | Paul Trust and Rosane Minghim. 2024. A study on text classification in the age of large language models. <i>Machine Learning and Knowledge Extraction</i> , 6(4):2688–2721. | 863 864 865 866 |
| | Shuo Wan, Tianyang Gao, Yichi Zhang, and Zhen Wang. 2023. A survey of deep active learning for foundation models. <i>Intelligent Computing</i> , 2(2):0058. | 867 868 869 |
| | Xiaoqiang Wang, Bang Liu, Siliang Tang, and Lingfei Wu. 2022. QRelScore: Better evaluating generated questions with deeper understanding of context-aware relevance. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 562–581, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | 870 871 872 873 874 875 876 |
| | L L Weed. 1968. Medical records that guide and teach. <i>N. Engl. J. Med.</i> , 278(11):593–600. | 877 878 |
| | Chunhua Weng, Xiaoying Wu, Zihui Luo, Mary Regina Boland, Dimitri Theodoratos, and Stephen B Johnson. 2011. EliXR: an approach to eligibility criteria extraction and representation. <i>J. Am. Med. Inform. Assoc.</i> , 18 Suppl 1(Supplement 1):i116–24. | 879 880 881 882 883 884 |
| | Michael Wornow, Alejandro Lozano, Dev Dash, Jenelle Jindal, Kenneth W Mahaffey, and Nigam H Shah. 2025. Zero-shot clinical trial patient matching with llms. <i>NEJM AI</i> , 2(1):AICs2400360. | 885 886 887 888 |
| | Xiaoyang Yi, Yuru Bao, Jian Zhang, Yifang Qin, and Faxin Lin. 2024. Integrating structural semantic knowledge for enhanced information extraction pre-training. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , | 889 890 891 892 893 |

| | | | |
|-----|--|---|------|
| 894 | pages 2156–2171, Miami, Florida, USA. Association | | |
| 895 | for Computational Linguistics. | | |
| 896 | Chi Yuan, Patrick B. Ryan, Casey Ta, Yixuan Guo, Zi- | | |
| 897 | ran Li, Jill Hardin, Rupa Makadia, Peng Jin, Ning | | |
| 898 | Shang, Tian Kang, and Chunhua Weng. 2019. Crite- | | |
| 899 | ria2Query: a natural language interface to clinical | | |
| 900 | databases for cohort definition. <i>J. Am. Med. Inform.</i> | | |
| 901 | <i>Assoc.</i> , 26(4):294–305. | | |
| 902 | Hongyi Yuan, Zheng Yuan, Ruyi Gan, Jiaying Zhang, | | |
| 903 | Yutao Xie, and Sheng Yu. 2022. BioBART: Pretrain- | | |
| 904 | ing and evaluation of a biomedical generative lan- | | |
| 905 | guage model . In <i>Proceedings of the 21st Workshop</i> | | |
| 906 | <i>on Biomedical Language Processing</i> , pages 97–109, | | |
| 907 | Dublin, Ireland. Association for Computational Lin- | | |
| 908 | guistics. | | |
| 909 | Jingqing Zhang, Yao Zhao, Mohammad Saleh, and | | |
| 910 | Peter J. Liu. 2019. Pegasus: Pre-training with ex- | | |
| 911 | tracted gap-sentences for abstractive summarization . | | |
| 912 | <i>Preprint</i> , arXiv:1912.08777. | | |
| 913 | Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. | | |
| 914 | Weinberger, and Yoav Artzi. 2020. Bertscore: | | |
| 915 | Evaluating text generation with bert . <i>Preprint</i> , | | |
| 916 | arXiv:1904.09675. | | |
| 917 | Angelo Ziletti and Leonardo D’Ambrosi. 2025. Gener- | | |
| 918 | ating patient cohorts from electronic health records | | |
| 919 | using two-step retrieval-augmented text-to-sql gener- | | |
| 920 | ation . <i>Preprint</i> , arXiv:2502.21107. | | |
| 921 | A Appendix | | |
| 922 | A.1 Prompts | | |
| 923 | A.1.1 CTQG initial dataset generation | | |
| 924 | <i>I’ll give you (one by one) a set of clinical trial inclu-</i> | | |
| 925 | <i>sion/exclusion criteria and I need you to turn them into a</i> | | |
| 926 | <i>question, as if in a form. e.g. for criteria "patient is > 18" the</i> | | |
| 927 | <i>question would be "is the patient > 18?"</i> | | |
| 928 | A.1.2 Question generation for Clinical Trials | | |
| 929 | <i>You are a clinical assistant specialized in communicating el-</i> | | |
| 930 | <i>igibility criteria for clinical trials. Your task is to transform</i> | | |
| 931 | <i>inclusion or exclusion criteria into understandable questions</i> | | |
| 932 | <i>directed at patients in English, while maintaining both clinical</i> | | |
| 933 | <i>accuracy and clarity of language.</i> | | |
| 934 | <i>Additionally, you must briefly reason out loud before generat-</i> | | |
| 935 | <i>ing each question, explaining step by step how you break down</i> | | |
| 936 | <i>the criterion, what terms you identify, and how you decide to</i> | | |
| 937 | <i>formulate the questions. This reasoning must be clear, logical,</i> | | |
| 938 | <i>and concise, and should be included just before each group of</i> | | |
| 939 | <i>questions.</i> | | |
| 940 | <i>Objective: Convert each criterion into one or more clear and</i> | | |
| 941 | <i>understandable questions in English, without losing medical</i> | | |
| 942 | <i>rigor.</i> | | |
| 943 | <i>Important instructions (read carefully and follow all of them):</i> | | |
| 944 | 1. Analyze each criterion carefully: if the same criterion | | |
| 945 | contains multiple clinical conditions, diseases, values, | | |
| 946 | histories, or options, you must generate a separate ques- | | |
| 947 | tion for each of those elements. Do not combine multiple | | |
| 948 | conditions into a single question. | | |
| | | 2. Before generating the questions, write a brief explana- | 949 |
| | | tion of your reasoning in English (chain of thought), | 950 |
| | | identifying the relevant parts of the clinical criterion | 951 |
| | | and justifying how you are going to split them. | 952 |
| | | 3. Always use the following format for your answers: | 953 |
| | | • “Reasoning: [your step-by-step explanation in En- | 954 |
| | | glish]” | 955 |
| | | • “New question N: [generated question]” | 956 |
| | | • If there is more than one question for that crite- | 957 |
| | | rion: “New question N: [question 1] [question | 958 |
| | | 2] [question 3] ...” | 959 |
| | | 4. Keep the original abbreviations (e.g., VIH, BMI, COPD). | 960 |
| | | Do not explain or expand them. | 961 |
| | | 5. Do not rewrite or interpret the criteria. Only transform | 962 |
| | | them into questions. | 963 |
| | | 6. Always generate the questions in English. | 964 |
| | | <i>Reasoning example:</i> | 965 |
| | | <i>Criterion: Age ≥ 18 years old Reasoning: The criterion refers</i> | 966 |
| | | <i>to age, with a threshold of 18 years. It is a single condition</i> | 967 |
| | | <i>related to minimum age. New question 1: Is the patient ≥18</i> | 968 |
| | | <i>years old?</i> | 969 |
| | | <i>Other examples without reasoning:</i> | 970 |
| | | <i>Example 1:</i> | 971 |
| | | <i>Criterion: Age ≥ 18 years old</i> | 972 |
| | | <i>Question: Is the patient ≥18 years old?</i> | 973 |
| | | <i>Example 2:</i> | 974 |
| | | ... | 975 |
| | | A.1.3 Question generation for SQuAD | 976 |
| | | <i>You are a documentalist, expert in converting sentences into</i> | 977 |
| | | <i>questions. Your task is to extract from the original text under-</i> | 978 |
| | | <i>standable questions in English which answer is included in</i> | 979 |
| | | <i>the same text, while maintaining both accuracy and clarity of</i> | 980 |
| | | <i>language.</i> | 981 |
| | | <i>Additionally, you must briefly reason out loud before generat-</i> | 982 |
| | | <i>ing each question, explaining step by step how you break</i> | 983 |
| | | <i>down the text, what terms you identify, and how you decide to</i> | 984 |
| | | <i>formulate the questions. This reasoning must be clear, logical,</i> | 985 |
| | | <i>and concise, and should be included just before each group of</i> | 986 |
| | | <i>questions.</i> | 987 |
| | | <i>Objective: Convert each sentence into one or more clear and</i> | 988 |
| | | <i>understandable questions in English, without losing rigor.</i> | 989 |
| | | <i>Important instructions (read carefully and follow all of them):</i> | 990 |
| | | 1. Analyze each text carefully: if the same text contains | 991 |
| | | multiple relevant questions, you must generate them. If | 992 |
| | | you consider that there is no relevant question to be | 993 |
| | | generated, simply write “Unknown” | 994 |
| | | 2. Before generating the questions, write a brief explana- | 995 |
| | | tion of your reasoning in English (chain of thought), | 996 |
| | | identifying the relevant parts of the text and justifying | 997 |
| | | how you are going to split them. | 998 |
| | | 3. Always use the following format for your answers: | 999 |
| | | • “Reasoning: [your step-by-step explanation in En- | 1000 |
| | | glish]” | 1001 |
| | | • “New question N: [generated question]” | 1002 |
| | | • If there is more than one question for that crite- | 1003 |
| | | rion: “New question N: [question 1] [question | 1004 |
| | | 2] [question 3] ...” | 1005 |
| | | 4. Keep the original abbreviations (e.g., VIH, BMI, COPD). | 1006 |
| | | Do not explain or expand them. | 1007 |

| | | | | |
|------|---|---|--|------|
| 1008 | 5. Do not rewrite or interpret the text. Only transform | 2. | Binary questions by default: questions are | 1053 |
| 1009 | them into questions. | | yes/no. The logical polarity (inclusion vs. ex- | 1054 |
| 1010 | 6. Always generate the questions in English. | | clusion) is not encoded in the wording. For | 1055 |
| 1011 | <i>Remember: if necessary, you should generate various ques-</i> | | example, both “patients with X are excluded” | 1056 |
| 1012 | <i>tions separated by “?”, and avoid including new lines in the</i> | | and “patients must have X” yield questions of | 1057 |
| 1013 | <i>questions and between them if many are generated</i> | | the form “Does the patient have X?” | 1058 |
| 1014 | <i>Example 1:</i> | | | |
| 1015 | <i>Criterion: Cultural reactions have varied from ridicule to ad-</i> | 3. | Preserve clinical content and terminology: | 1059 |
| 1016 | <i>miration; many common stereotypes exist regarding redheads</i> | | clinical constructs (e.g., ECOG, TNM, PD- | 1060 |
| 1017 | <i>and they are often portrayed as fiery-tempered.</i> | | L1, drug names, numeric thresholds) are pre- | 1061 |
| 1018 | <i>Question: What type of temper are people with red hair con-</i> | | served verbatim. Acronyms and units are not | 1062 |
| 1019 | <i>sidered to have?</i> | | simplified or explained. | 1063 |
| 1020 | <i>Example 2:</i> | | | |
| 1021 | ... | | | |
| 1022 | B Annotation guidelines for | 4. | Abstract away protocol management: | 1064 |
| 1023 | criterion-to-question mapping | | protocol-management phrasing (“must be”, | 1065 |
| 1024 | | | “is required”, “prior to enrolment”) is removed | 1066 |
| 1025 | Unit of annotation The primary unit is an <i>atomic</i> | | when it does not alter the clinical condition. | 1067 |
| 1026 | <i>clinical condition</i> expressed in an eligibility crite- | | Time constraints that affect eligibility are kept | 1068 |
| 1027 | rion (e.g., age constraint, ECOG status, prior treat- | | (e.g., “within N days before treatment”). | 1069 |
| 1028 | ment, specific comorbidity, organ function, preg- | B.2.2 Numbers, thresholds and comparators | | 1070 |
| 1029 | nancy/contraception). A single textual criterion | | 5. Preserve numeric operators and units: in- | 1071 |
| 1030 | may yield one or multiple questions if it contains | | equalities (\geq , \leq , $>$, $<$) are kept as symbols | 1072 |
| 1031 | several independent conditions or clinically rele- | | in the question text. Units (g/dL, \times ULN, | 1073 |
| 1032 | vant exceptions. | | mL/min, mmHg, months, weeks, years, etc.) | 1074 |
| 1033 | B.1 Inputs and annotator setup | | are preserved. Example pattern: | 1075 |
| 1034 | | | Criterion: “Age $\geq X$ years” \rightarrow | 1076 |
| 1035 | • Source documents: full inclusion/exclusion | | Question: “Is the patient $\geq X$ years | 1077 |
| 1036 | sections of oncology trial protocols. | | old?” | 1078 |
| 1037 | • Annotator: one medically literate annota- | | | |
| 1038 | tor, working alone (no adjudication, no inter- | 6. | Multi-constraint tests: when a criterion ag- | 1079 |
| 1039 | annotator agreement). | | gregates several lab thresholds or organ func- | 1080 |
| 1040 | • LLM assistance: ChatGPT is used to draft | | tion conditions that jointly define a single con- | 1081 |
| 1041 | questions from criteria; the annotator then | | cept (e.g., “adequate organ function”), a single | 1082 |
| 1042 | manually edits the output to follow the rules | | question is created: | 1083 |
| 1043 | in this section. | | “Does the patient meet the criteria | 1084 |
| 1044 | • Effort: typical time per question is a few sec- | | for adequate [organ] function (e.g., | 1085 |
| 1045 | onds, with a single prompt per batch of criteria | | [list of thresholds])?” | 1086 |
| 1046 | (no iterative prompt engineering per item). | | If thresholds appear elsewhere as independent | 1087 |
| 1047 | B.2 Question generation rules | | criteria, they may be annotated separately. | 1088 |
| 1048 | B.2.1 General phrasing | B.2.3 Splitting criteria into multiple questions | | 1089 |
| 1049 | 1. Patient-centred formulation: questions are | | 7. One atomic condition \rightarrow one question: if a | 1090 |
| 1050 | formulated about the patient, not about ab- | | criterion includes several independent clauses | 1091 |
| 1051 | stract “participants”. Typical patterns: | | (e.g., age + performance status + disease | 1092 |
| 1052 | • “Is the patient ... ?” | | stage), these are split into separate questions, | 1093 |
| | • “Does the patient have ... ?” | | one per atomic condition. | 1094 |
| | • “Has the patient ... ?” | | 8. Exceptions and complex logic: for patterns | 1095 |
| | | | of the form “X except Y and Z”: | 1096 |

- At least one question captures the broad condition (e.g., “Does the patient have any previous or current malignancy?”).
- Additional questions may be created for clinically important exceptions when they are meaningful variables (e.g., major cardiovascular disease).
- Benign exceptions (e.g., vitiligo, alopecia) are not systematically turned into separate questions; they remain implicit in the logic.

B.2.4 Multi-choice vs. yes/no questions

9. **Yes/No for presence/absence:** conditions that are naturally binary (pregnancy, presence of a comorbidity, history of a given treatment, presence of brain metastases, etc.) are encoded as yes/no questions.
10. **List when categories are enumerated:** if a criterion explicitly enumerates mutually exclusive categories (e.g., tumour types, histologic subtypes, study cohorts), it is encoded as a single multi-choice question:

“Which of the following applies to the patient? (1) ..., (2) ..., (3) ...”

The corresponding answer type is a list. Follow-up criteria that merely continue the list are not considered separate items.

B.2.5 Target population

11. **Question population:** for each question, a separate field encodes the subpopulation for which the question is meaningful, e.g.:
 - “All patients”
 - “Women of childbearing potential”
 - “Male participants”
 - “Only if [brain metastases = yes]”

Optionally, population restrictions may be repeated in the question text as a leading clause (e.g., “For women of childbearing potential: ...”).

B.3 Summary of the annotation process

All questions are created by a single medically literate annotator using ChatGPT as a drafting tool, followed by manual post-editing to enforce the guidelines above. Each criterion is processed once (no double annotation, no adjudication), and the order of questions within each trial follows the original

ordering of criteria. When a criterion yields multiple questions, these appear immediately after the parent in the internal sheet.

C Question division experimental details

C.1 Data preparation details

- CTQG: there is a large imbalance in classes, where class 1 is over-represented compared to the other two (Table 8).

| Class | Description | Proportion |
|-------|---------------------------------------|------------|
| 0 | criterion generates no question | 0.196 |
| 1 | criterion generates one question | 0.647 |
| 2 | criterion generates several questions | 0.157 |

Table 8: Data class support for CTQG question division prediction.

- QG-SQuAD: there is more balance than in CTQG, but this time there is one class that is under-represented, class 2 (Table 9).

| Class | Description | Proportion |
|-------|---------------------------------------|------------|
| 0 | criterion generates no question | 0.405 |
| 1 | criterion generates one question | 0.405 |
| 2 | criterion generates several questions | 0.190 |

Table 9: Data class support for QG-SQuAD question division prediction.

C.2 Rules sample

The rules we extracted from the CTQG features are available at Table 10 and the ones for QG-SQuAD at Table 11.

Algorithm 2 Question division algorithm

Require: Labeled criterion $\{(x_i, y_i)\}$

- 1: Extract features $\Phi(x_i)$
- 2: Fit/keep high-precision rules \mathcal{R}
- 3: Train RandomForest on $(\Phi(x_i), y_i)$
- 4: Compute embeddings $E(x_i)$
- 5: Fine-tune transformer models $\{f_m\}$ on $(E(x_i), y_i)$

D Complementary results

Tables 13 and 14 present the complementary results of Tables 6 and 7, respectively, with additional metrics introduced section 4.2.3.

E Complementary information

In this last Appendix, we include all details on the methods we used (Tables 15, 16) or on the experimental set-up (Tables 17, 19, 18, 20).

| Rule condition | Class | Precision | Support |
|---|-------|-----------|---------|
| raw_word_len \geq 36 AND shannon_char_entropy \leq 4.1787 | 2 | 1.00 | 5 |
| raw_word_len \geq 39 AND num_clauses \leq 4 | 2 | 1.00 | 5 |
| raw_word_len \leq 2 AND num_words \leq 2 | 0 | 1.00 | 5 |
| raw_char_len \leq 179 AND raw_char_2grams_voc \geq 113 | 2 | 1.00 | 5 |
| raw_char_len \geq 263 AND raw_char_2grams_voc \leq 142 | 2 | 1.00 | 5 |
| raw_char_len \geq 216 AND raw_char_3grams_voc \leq 174 | 2 | 1.00 | 5 |
| raw_char_len \geq 216 AND raw_char_4grams_voc \leq 194 | 2 | 1.00 | 6 |
| raw_char_len \geq 273 AND raw_word_2grams_voc_lemma \leq 38 | 2 | 1.00 | 5 |
| raw_char_len \leq 188 AND raw_word_2grams_voc_pos \geq 15 | 2 | 1.00 | 5 |
| raw_char_len \geq 271 AND num_clauses \leq 4 | 2 | 1.00 | 5 |
| raw_char_len \leq 12 AND num_words \leq 2 | 0 | 1.00 | 6 |
| raw_voc_len \geq 28 AND shannon_char_entropy \leq 4.1705 | 2 | 1.00 | 5 |
| raw_voc_len \leq 2 AND num_words \leq 2 | 0 | 1.00 | 5 |
| raw_voc_len \leq 33 AND num_words \geq 47 | 2 | 1.00 | 5 |

Table 10: Validated rule set for the main classifier. Each rule is defined by a condition, target class, precision, and support.

Algorithm 3 Inference and voting (precision-prioritized)

Require: Criterion x , rules \mathcal{R} , RF, models $\{f_m\}$, thresholds

- 1: Compute $\Phi(x)$ and $E(x)$
 - 2: If a rule in \mathcal{R} fires with confidence $\geq \tau_r$, **return** class
 - 3: Get RF prediction (y_{rf}, c_{rf})
 - 4: Get FT predictions $\{(y_m, c_m)\}$ and majority y_{maj} with support p_{maj} and min confidence C_{maj}
 - 5: **if** $p_{maj} \geq \tau_{maj}$ **and** $C_{maj} \geq \tau_{ft}$ **then return** y_{maj}
 - 6: **end if**
 - 7: **if** $c_{rf} \geq \tau_{ft}$ **then return** y_{rf}
 - 8: **end if**
 - 9: **return** IDK
-

| Rule condition | Class | Precision | Support |
|--|-------|-----------|---------|
| raw_char_2grams_voc \leq 0 AND shan- non_word_entropy \geq 4 | 0 | 0.90 | 10 |
| raw_voc_len \geq 5 AND raw_word_2grams_voc_pos \leq 0 | 0 | 0.83 | 6 |
| raw_char_len \leq 2 AND raw_word_voc_pos \geq 5 | 1 | 0.80 | 5 |
| raw_char_4grams_voc \geq 5 AND raw_word_2grams_voc_pos \leq 0 | 0 | 0.80 | 5 |
| raw_char_4grams_voc \leq 0 AND shan- non_word_entropy \geq 3 | 0 | 0.78 | 36 |
| raw_char_len \geq 5 AND raw_word_2grams_voc_pos \leq 0 | 0 | 0.75 | 8 |
| raw_char_2grams_voc \leq 0 AND num_words \geq 4 | 0 | 0.75 | 8 |
| raw_char_len \leq 0 AND shannon_word_entropy \geq 3 | 0 | 0.74 | 42 |
| raw_voc_len \leq 2 AND raw_word_2grams_voc_pos \geq 5 | 1 | 0.71 | 7 |

Table 11: Validated rule set derived from the SQuAD-style subset.

| Hyperparameter | value |
|---------------------|---------------|
| Activation function | Softmax |
| loss function | cross-entropy |
| dropout | 0.1 |
| class weighting | No |
| weight decay | 0.000001 |
| epochs | 3 |
| batch size | 32 |
| learning rate | 0.000005 |
| warm-up | 1000 steps |

Table 12: Best hyper parameters for transformers fine-tuning for question division prediction with QG-SQuAD.

| Metrics (I) | QRelScore | METEOR | ROUGE-L | QSTS | QBLEU | BARTScore |
|------------------|---------------|--------------|--------------------|-------------|-------------|--------------|
| BERT2BERT | 0.34 | 0.32 | 0.47 | 0.69 | 0.38 | -4.65 |
| bioBART | 0.43 | 0.39 | 0.52 | 0.79 | 0.41 | -4.37 |
| Llama3.3 | 0.68 | 0.44 | 0.47 | 0.79 | 0.46 | -3.55 |
| FLAN-T5 | 0.66 | 0.36 | 0.41 | 0.72 | 0.35 | -3.83 |
| ChatGPT | 0.66 | 0.44 | 0.47 | 0.74 | 0.44 | -3.49 |
| DivAC-LLM | 0.74 | 0.53 | 0.56 | 0.77 | 0.52 | -3.16 |
| Metrics (II) | BARTScore | fluency | clarity | conciseness | relevance | consistency |
| BERT2BERT | 0.77 | -5.22 | 0.88 | 0.76 | 0.77 | 0.78 |
| bioBART | 0.81 | -4.86 | 0.88 | 0.80 | 0.81 | 0.82 |
| Llama3.3 | 0.92 | -3.60 | 0.89 | 0.92 | 0.92 | 0.92 |
| FLAN-T5 | 0.86 | -3.85 | 0.86 | 0.85 | 0.86 | 0.86 |
| ChatGPT | 0.91 | -3.69 | 0.88 | 0.90 | 0.91 | 0.91 |
| DivAC-LLM | 0.93 | -3.38 | 0.91 | 0.92 | 0.93 | 0.92 |
| Metrics (III) | answerability | acceptance | answer consistency | overall | mean | scaled |
| BERT2BERT | 0.77 | 0.78 | 0.75 | 0.78 | 0.11 | 0.32 |
| bioBART | 0.81 | 0.81 | 0.78 | 0.82 | 0.19 | 0.46 |
| Llama3.3 | 0.93 | 0.91 | 0.84 | 0.91 | 0.36 | 0.71 |
| FLAN-T5 | 0.87 | 0.85 | 0.80 | 0.85 | 0.29 | 0.58 |
| ChatGPT | 0.92 | 0.90 | 0.85 | 0.90 | 0.35 | 0.69 |
| DivAC-LLM | 0.93 | 0.92 | 0.83 | 0.91 | 0.42 | 0.83 |

Table 13: Extended automatic evaluation results on CTQG with respect to five different models compared with DivAC-LLM.

| Metrics (I) | QRelScore | METEOR | ROUGE-L | QSTS | QBLEU | BARTScore |
|------------------|---------------|-------------|--------------------|-------------|-------------|--------------|
| Gemini | 0.46 | 0.23 | 0.23 | 0.62 | 0.15 | -4.28 |
| Gemini + div | 0.35 | 0.24 | 0.30 | 0.58 | 0.21 | -4.30 |
| ChatGPT | 0.47 | 0.22 | 0.22 | 0.62 | 0.15 | -4.33 |
| ChatGPT + div | 0.33 | 0.24 | 0.31 | 0.59 | 0.21 | -4.37 |
| Llama3.3 | 0.32 | 0.17 | 0.24 | 0.52 | 0.18 | -4.63 |
| Llama3.3 + div | 0.28 | 0.21 | 0.32 | 0.54 | 0.17 | -4.47 |
| ProphetNet-large | 0.26 | 0.28 | 0.34 | 0.30 | 0.14 | -3.75 |
| Pegasus-XSum | 0.37 | 0.24 | 0.32 | 0.36 | 0.17 | -3.79 |
| Metrics (II) | BARTScore | fluency | clarity | conciseness | relevance | consistency |
| Gemini | -3.74 | 0.81 | 0.89 | 0.90 | 0.90 | 0.90 |
| Gemini + div | -4.28 | 0.78 | 0.82 | 0.82 | 0.84 | 0.82 |
| ChatGPT | -3.77 | 0.85 | 0.91 | 0.91 | 0.92 | 0.91 |
| ChatGPT + div | -4.35 | 0.77 | 0.79 | 0.80 | 0.82 | 0.80 |
| Llama3.3 | -4.56 | 0.87 | 0.87 | 0.87 | 0.88 | 0.87 |
| Llama3.3 + div | -4.64 | 0.86 | 0.86 | 0.86 | 0.88 | 0.86 |
| ProphetNet-large | -4.19 | 0.66 | 0.76 | 0.77 | 0.77 | 0.77 |
| Pegasus-XSum | -4.21 | 0.87 | 0.90 | 0.92 | 0.92 | 0.92 |
| Metrics (III) | answerability | acceptance | answer consistency | overall | mean | scaled |
| Gemini | 0.90 | 0.89 | 0.73 | 0.86 | 0.19 | 0.62 |
| Gemini + div | 0.83 | 0.82 | 0.72 | 0.81 | 0.15 | 0.52 |
| ChatGPT | 0.91 | 0.91 | 0.75 | 0.88 | 0.19 | 0.65 |
| ChatGPT + div | 0.80 | 0.80 | 0.73 | 0.79 | 0.14 | 0.50 |
| Llama3.3 | 0.88 | 0.88 | 0.69 | 0.85 | 0.11 | 0.47 |
| Llama3.3 + div | 0.86 | 0.87 | 0.73 | 0.85 | 0.14 | 0.56 |
| ProphetNet-large | 0.77 | 0.77 | 0.67 | 0.74 | 0.15 | 0.46 |
| Pegasus-XSum | 0.92 | 0.91 | 0.76 | 0.89 | 0.22 | 0.81 |

Table 14: Extended automatic evaluation results on QG-SQuAD comparing fine-tuned models and LLMs with and without question division module.

| Feature | Mean value |
|-------------------------------|-------------------|
| AND or OR in text | 0.52 |
| has semicolon | 0.03 |
| has conjunction | 0.57 |
| number of clauses | 2.49 |
| number of words | 18.32 |
| word length | 15.73 |
| character length | 108.58 |
| vocab. length | 13.82 |
| word bigrams | 14.44 |
| character bigrams | 68.92 |
| word trigrams | 13.66 |
| character trigrams | 89.26 |
| word quadrigrams | 12.72 |
| character quadrigrams | 94.41 |
| word lemmatized | 13.73 |
| word bigrams_lemma | 14.42 |
| word trigrams_lemma | 13.65 |
| word quadrigrams_lemma | 12.72 |
| word POS | 7.39 |
| word bigrams POS | 7.00 |
| word trigrams POS | 6.13 |
| word quadrigrams POS | 5.28 |
| word NER | 1.24 |
| word bigrams NER | 0.66 |
| word trigrams NER | 0.31 |
| word quadrigrams NER | 0.14 |
| shannon character entropy | 4.01 |
| shannon word entropy | 3.45 |
| conditional word entropy | 0.12 |
| conditional character entropy | 0.08 |
| renyi word entropy | 0.12 |
| renyi character entropy | 1.69 |

Table 15: Average feature values per criterion in CTQG.

| Voting Strategy | Description / Decision Rule |
|----------------------|--|
| pred | If both the rule-based model (prediction) and BioBERT agree on a class different from 2, use that class; otherwise, output 2. |
| default1 | If both the rule-based model and BioBERT agree on a class different from 1, use that class; otherwise, output 1. |
| majority_transformer | If the two transformers (BioBERT and BioMedLM) agree, take their class; otherwise, fall back to the rule-based model. |
| majority_rules | If the rule-based model agrees with either transformer, use the rule-based prediction; otherwise, fall back to BioBERT. |
| hybrid_priority | Prioritize class 1 if any transformer predicts 1; prioritize class 2 if rules predict 2; otherwise, if rules and BioBERT agree, keep that; else 0. |
| consensus_only | Output the class only if all three models (rules, BioBERT, BioMedLM) fully agree; otherwise, label as -1 (uncertain). |
| strict_negative | Reject (-1) if all predict -1 ; if all agree otherwise, output that class; else fall back to BioBERT. |
| 2of3 | If any two of the three models agree, take that class; otherwise, label as -1 (no consensus). |

Table 16: Overview of voting strategies combining rule-based predictions, BioBERT, and BioMedLM outputs.

| Hyperparameter | Value |
|---------------------|---------------|
| Activation function | Softmax |
| Loss function | Cross-entropy |
| Dropout | 0.5 |
| Class weighting | No |
| Weight decay | 5e-6 |
| Epochs | 8 |
| Batch size | 4 |
| Learning rate | 1e-5 |
| Warm-up | 100 steps |

Table 17: Best hyperparameters for transformer fine-tuning (question division, CTQG).

| Metric | Value | Assumptions / Notes |
|---------------------------|---------------------|--|
| Total utterances | 92 | |
| Total tokens | 24,949 | |
| Avg. tokens / utterance | 271.200 | Computed as 24,949/92 |
| GPU time | 0.153 h (9.000 min) | Inference only |
| Energy consumption | 0.046 kWh | 0.153 h \times 0.300 kW |
| CO ₂ emissions | 9.200 g | EU grid avg. 0.200 kgCO ₂ /kWh |
| Electricity cost | €0.009 | €0.20/kWh (EU, 2024) |
| Cloud rental | €0.380 | €2.50/GPUh (on-demand) |

Table 18: Approximate compute and environmental footprint for 92 utterances (24,949 tokens) on one NVIDIA A100 40GB at 300.000 W. Power and emissions factors from NVIDIA (2022) and IEA (2023).

| Experiment | Duration | GPU h | Energy (kWh) | Elec. (€) | Cloud (€) | CO ₂ (kg) |
|------------|----------|-------|--------------|-----------|-----------|----------------------|
| CTQG | ~2 weeks | 350 | 105 | 21 | 875 | 21 |
| QG-SQuAD | ~1 month | 700 | 210 | 42 | 1750 | 42 |

Table 19: Runtime and cost estimates for grid search (question generation + division) on a single A100 40GB at 300.000 W. Prices: €0.20/kWh (EU 2024), €2.50/GPUh (on-demand). Emissions: 0.200 kgCO₂/kWh.

| Parameter | CTQG | QG-SQuAD |
|------------------|-------------|-----------------|
| Batch size | 16 | 32 |
| Learning rate | 0.0001 | 0.00001 |
| Epochs | 8 | 12 |
| Warm-up steps | 50 | 1500 |

Table 20: Hyperparameter configurations used in CTQG and QG-SQuAD experiments.

| Hyperparameter | CTQG | QG-SQuAD |
|-----------------------|---|----------------------|
| Epochs | [4, 6, 8, 10, 12] | [6, 8, 10, 12, 14] |
| Batch size | [16, 32, 64, 128] | [32, 64] |
| Learning rate | $[5e^{-6}, 1e^{-6}, 1e^{-5}, 1e^{-4}, 5e^{-5}]$ | $[5e^{-5}, 1e^{-5}]$ |
| Warm-up steps | [10, 50, 100] | [1k, 5k, 10k] |

Table 21: Grid search hyperparameter ranges for CTQG and QG-SQuAD experiments.