# MIRROR: A Novel Approach for the Automated Evaluation of Open-Ended Question Generation

**Aniket Deroy**
Computer Science and Engineering
IIT Kharagpur
roydanik18@kgpian.iitkgp.ac.in

**Subhankar Maity**
Department of Artificial Intelligence
IIT Kharagpur
subhankar.ai@kgpian.iitkgp.ac.in

**Sudeshna Sarkar**
Computer Science and Engineering
IIT Kharagpur
sudeshna@cse.iitkgp.ac.in

## Abstract

Automatic question generation is a critical task that involves evaluating question quality by considering factors such as engagement, pedagogical value, and the ability to stimulate critical thinking. These aspects require human-like understanding and judgment, which automated systems currently lack. However, human evaluations are costly and impractical for large-scale samples of generated questions. Therefore, we propose a novel system, MIRROR (Multi-LLM Iterative Review and Response for Optimized Rating), which leverages large language models (LLMs) to automate the evaluation process for questions generated by automated question generation systems. We experimented with several state-of-the-art LLMs, such as GPT-4, Gemini, and Llama2-70b. We observed that the scores of human evaluation metrics, namely relevance, appropriateness, novelty, complexity, and grammaticality, improved when using the feedback-based approach called MIRROR, tending to be closer to the human baseline scores. Furthermore, we observed that Pearson's correlation coefficient between GPT-4 and human experts improved when using our proposed feedback-based approach, MIRROR, compared to direct prompting for evaluation. Error analysis shows that our proposed approach, MIRROR, significantly helps to improve relevance and appropriateness.

## 1 Introduction

Automated question generation (AQG) is crucial in education because it enhances critical thinking, promotes active learning, and provides personalized learning experiences. Currently, metrics such as BLEU, METEOR, and ROUGE are used for the purpose of evaluation of open ended questions [Zhang et al., 2021, Al Faraby et al., 2023]. However, traditional automated evaluation metrics such as BLEU, ROUGE, or METEOR, which are often used for tasks like machine translation or summarization, may not be well-suited for evaluating the quality of generated questions [Alva-Manchego et al., 2021, Bhandari et al., 2020, Reiter, 2018]. These metrics primarily measure surface-level similarity to reference questions, rather than deeper aspects of question quality such as relevance, clarity, and engagement [Nema and Khapra, 2018]. Some aspects of question quality, such as engagement, pedagogical value, and the potential to stimulate critical thinking, require a human-like understanding and judgment that automated systems cannot yet replicate. However human judgements are costly and cannot be replicated over large samples of generated questions. So we propose an LLM feedback-based system called **MIRROR** (*Multi-LLM Iterative Review and Response for Optimized Rating*) for

automating the process of human evaluation for questions generated by AQG systems. We understand that LLMs have understanding and reasoning capabilities.

So, the question we ask is *Can LLMs replace human experts in evaluating the quality of generated questions by AQG systems?* The results show that our proposed feedback-based approach for evaluating questions is worthwhile in generating quality evaluation for questions generated by AQG systems. We observe that the scores of the human evaluation metrics namely *relevance* [Maity et al., 2024b], *appropriateness* [Maity et al., 2024b], *novelty* [Maity et al., 2024b], *complexity* [Gong et al., 2022, Maity et al., 2024b], and *grammaticality* [Ushio et al., 2022, Maity et al., 2024b,a] improve on using the feedback-based approach called MIRROR and tend to be closer to the human baseline scores. Also, we observe that the Pearson's correlation coefficient between GPT-4 and human experts improves on using our proposed feedback-based approach called MIRROR as compared to using direct prompting for evaluation.

Our contributions in this work are as follows:

1. We propose a novel method, **MIRROR** (**M**ulti-LLM **I**terative **R**eview and **R**esponse for **O**ptimized **R**ating). This feedback-based procedure involves prompting various LLMs to generate scores based on human evaluation metrics, namely *grammaticality*, *relevance*, *appropriateness*, *novelty*, and *complexity*, for the task of automated open-ended question generation. We then ask the same model to list the strengths and flaws of the provided questions. These strengths and flaws are subsequently given to another LLM along with the questions, prompting it to repeat the evaluation procedure.

2. We show that our proposed LLM feedback-based approach, which provides scores based on various human evaluation metrics, produces results closer to the human baseline compared to the direct prompting approach.

3. We also measure the Pearson's correlation coefficient between the GPT-4 generated scores and those of human evaluators, demonstrating that the feedback-based approach improves correlation compared to the direct approach for evaluating generated questions. Error analysis shows that our proposed approach, MIRROR, significantly improves the *relevance* and *appropriateness* of the evaluations.

## 2  Related Work

Evaluating open-ended question generation using automated metrics (e.g., BLEU, ROUGE, METEOR, etc.) presents significant challenges [Nema and Khapra, 2018, Gong et al., 2022, Ushio et al., 2022, Maity et al., 2024a]. These metrics often fail to capture higher-order cognitive skills and overlook deeper educational values such as stimulating critical thinking [Al Faraby et al., 2023, Mulla and Gharpure, 2023]. Automated evaluations may not fully account for nuances such as context relevance and cognitive complexity, making it difficult to assess whether questions promote skills such as analysis, synthesis, and evaluation [Mulla and Gharpure, 2023]. Furthermore, scalability issues arise because human evaluation for large datasets is impractical [Al Faraby et al., 2023]. Using LLMs to evaluate human-like criteria such as *grammaticality*, *relevance*, *appropriateness*, *novelty*, and *complexity* has been a prominent research area [Chiang and Lee, 2023]. Aligning machine-generated content with human judgment is crucial for practical applicability. Incorporating feedback loops within LLMs to refine their outputs is an emerging field [Madaan et al., 2023, Chaudhari et al., 2024, Ouyang et al., 2022, Stiennon et al., 2020]. Techniques like Reinforcement Learning from Human Feedback (RLHF) fine-tune LLMs based on human preferences, enhancing the quality and contextual relevance of responses [Yuan et al., 2023, Chaudhari et al., 2024, Wu et al., 2023]. This work proposes a novel method using LLM feedback to evaluate open-ended questions based on human-like metrics: *grammaticality* [Ushio et al., 2022, Maity et al., 2024b,a], *relevance* [Maity et al., 2024b], *appropriateness* [Maity et al., 2024b], *complexity* [Gong et al., 2022, Maity et al., 2024b], and *novelty* [Maity et al., 2024b]. This approach aims to bridge the gap between automated metrics and human evaluation, enhancing the quality and applicability of generated questions in educational settings.

# 3 Dataset

We use 1000 samples from the EduProbe dataset [Maity et al., 2024b] for our experiments. The dataset consists of *<Context, Question>* pairs from subjects such as History, Geography, Economics, Environmental Studies, and Science. For EduProbe, we already have the gold-standard questions corresponding to the context. We also use 500 samples from the SciQ dataset [Welbl et al., 2017] for our experiments. The dataset consists of *<Context, Question>* pairs from subjects such as Physics, Chemistry, Biology, and Earth Science. For SciQ dataset we create a question corresponding to the context by using an educator. The contexts from both these datasets are considered for our experiments to generate the questions using the LLMs and then to evaluate them automatically via LLMs. Together EduProbe and SciQ datasets covers a wide domain of subjects necessary for generating open-ended questions and showing the wider applicability of our work.

# 4 Methodology

In this section, we discuss the generation of questions via prompting GPT-3.5 Turbo (Section 4.1), the direct approach to evaluating the quality of the generated questions (Section 4.2), the feedback-based approach for evaluating generated question quality (Section 4.3), and the correlation between the best-performing LLM and human experts (Section 4.4). In our experiments, we used state-of-the-art LLMs, such as GPT-4 [Achiam et al., 2023], Gemini [Team et al., 2023], and Llama2-70b [Touvron et al., 2023], to assess both direct and feedback-based approaches.

## 4.1 Generating Questions via Prompting

Figure 1 provides *<Context, Generated Question>* pair from the EduProbe and SciQ dataset respectively. Figure 2 provides the prompt for generating questions from a context. We prompt GPT-3.5 Turbo to generate question from a context corresponding to EduProbe and SciQ datasets. We evaluate the generated questions based on five metrics: *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity*.

| EduProbe | SciQ |
|---|---|
| **Context:** Purchasing power parity (PPP) is an economic indicator that signifies the purchasing power of the currencies of various nations of the world against each other. It helps in comparing living standards between different countries and estimating economic productivity. | **Context:** Without Coriolis Effect the global winds would blow north to south or south to north. But Coriolis makes them blow northeast to southwest or the reverse in the Northern Hemisphere. The winds blow northwest to southeast or the reverse in the southern hemisphere. |
| **Generated Question:** What does purchasing power parity do? | **Generated Question:** How does the Coriolis Effect influence the direction of global winds in both the Northern and Southern Hemispheres? |

Figure 1: A sample of *<Context, Generated Question>* pairs from the EduProbe and SciQ datasets.

Generate a Question based on the following context {Context}

Figure 2: Prompt used on GPT-3.5 Turbo to generate a question from a context.

## 4.2 Direct Approach for Evaluating Question Quality

The algorithm 1 provides the set of steps required to produce the direct approach. The overview diagram for our direct approach is shown in Figure 3. The prompt provided to the LLMs for the direct prompting approach is shown in Figure 4.

## 4.3 Feedback-based Approach for Evaluating Question Quality

The input to our proposed MIRROR approach includes human evaluation metric definitions, the initial strengths set ($S_0$), the initial flaws set ($F_0$), a generated question to be evaluated, and its corresponding context. The output from our approach is the final human evaluation metric scores. The goal of our algorithm is to perform repeated feedback between two LLMs (i.e., $LLM_1$, $LLM_2$)

---
**Algorithm 1** Direct Approach
---
**Require:** Human Evaluation Metric definitions
**Require:** Question and its context
**Ensure:** Scores
   Compute scores for human evaluation metrics via LLM prompting provided the given context, question, and metrics.
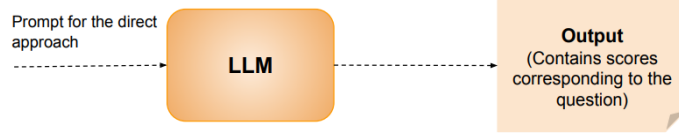
---



Figure 3: An overview of the direct prompting approach.



*Given the definition of human evaluation metrics:*

*Grammaticality: Measures the grammatical correctness of the generated question, regardless of the context;*
*Appropriateness: Examines the semantic correctness of the question irrespective of the context;*
*Relevance: Measures the degree to which the generated question is pertinent and aligned with the given context.*
*Complexity: Estimates the level of reasoning or cognitive effort required to answer the generated question.*
*Novelty: Measures the originality and distinctiveness of the generated question in comparison to the gold standard question for the given context.*

*The Context corresponding to the question is: {Context}*

*Provide a score between 1 and 5 for the question {Question} based on the human evaluation metrics namely Grammaticality, Appropriateness, Relevance, Complexity, and Novelty as well as the context corresponding to the question stated above.*

Figure 4: Prompt used in direct approach for evaluating human evaluation metrics.

so that the metric scores converge. The strengths and flaws generated by the LLMs are extracted from the entire output using a rule-based pattern-matching algorithm. The algorithm 2 provides the set of steps required to produce the MIRROR approach.

Initially, the process starts by defining a set of human evaluation metrics which are the criteria used to assess the quality of the question. These metrics include *grammaticality*, *appropriateness*, *relevance*, *complexity*, and *novelty* which contribute to the effectiveness of the question. At the outset, two sets are initialized: one for strengths ($S_0$) and another for flaws ($F_0$). Both sets start empty. The next step involves computing initial scores for the question based on the predefined evaluation metrics. From this assessment by $LLM_1$, the first sets of strengths and flaws are identified, resulting in $S_1$ and $F_1$. The process then enters an iterative loop designed to refine these initial assessments. In each iteration, the identified strengths and flaws ($S_1$, $F_1$) are provided as feedback to the second LLM, $LLM_2$. Along with this feedback, the human evaluation metrics, the question itself, and its context are also provided. $LLM_2$ then generates new scores for the evaluation metrics and updates the sets of strengths and flaws to $S_2$ and $F_2$. The updated strengths and flaws from $LLM_2$ (i.e., $S_2$, $F_2$) are then fed back to $LLM_1$. $LLM_1$ uses this information to re-evaluate the question, updating its scores and further refining the sets of strengths and flaws to $S_3$ and $F_3$. This process continues iteratively, with each model using the feedback from the other to refine its evaluation, until a convergence criterion is met. The convergence is typically defined as the point at which the scores from $LLM_1$ and $LLM_2$ become identical for two consecutive iterations. Once convergence is achieved, the loop terminates, and the final, converged scores and associated strengths and flaws are considered the accurate evaluation of the question. This iterative process ensures that the evaluation is thorough and benefits from the complementary perspectives of two different LLMs, ultimately leading to a more reliable and nuanced assessment of the question.

The overview diagram for our MIRROR approach is shown in Figure 5. The prompt provided to the LLMs for producing the MIRROR approach is shown in Figure 6. The method described is a process for refining the evaluation of a question using a feedback loop between two LLMs. The goal is to achieve a consensus on the evaluation scores and identify key strengths and flaws in the question, ensuring that the evaluation is as accurate and comprehensive as possible.

---

**Algorithm 2** MIRROR (Multi-LLM Iterative Review and Response for Optimized Rating)

---

**Require:** Human Evaluation Metric definitions
**Require:** Initial strengths set $S_0$ and flaws set $F_0$, both initialized as empty sets
**Require:** Question and its context
**Ensure:** Converged scores for the question
    $S_0 \leftarrow \{\}$
    $F_0 \leftarrow \{\}$
    Compute initial scores for human evaluation metrics with the given question and context
    **Identify initial strengths and flaws**
    $S_1, F_1 \leftarrow$ Identify strengths and flaws
    **while** convergence criteria are not met **do**
        Provide $S_1$, $F_1$ as feedback, along with metrics, question, and context, to $LLM_2$
        Ask $LLM_2$ to generate scores for human evaluation metrics
        Obtain strengths $S_2$ and flaws $F_2$ from $LLM_2$
        Provide $S_2$, $F_2$ as feedback, along with metrics, question, and context, to $LLM_1$
        Ask $LLM_1$ to provide updated scores for the question
        Obtain updated strengths $S_3$ and flaws $F_3$ from $LLM_1$
        Update: $S_1 \leftarrow S_3$, $F_1 \leftarrow F_3$
        **Check for convergence**
        **if** scores from $LLM_1$ and $LLM_2$ are identical for two consecutive iterations **then**
            **Terminate the loop**
        **end if**
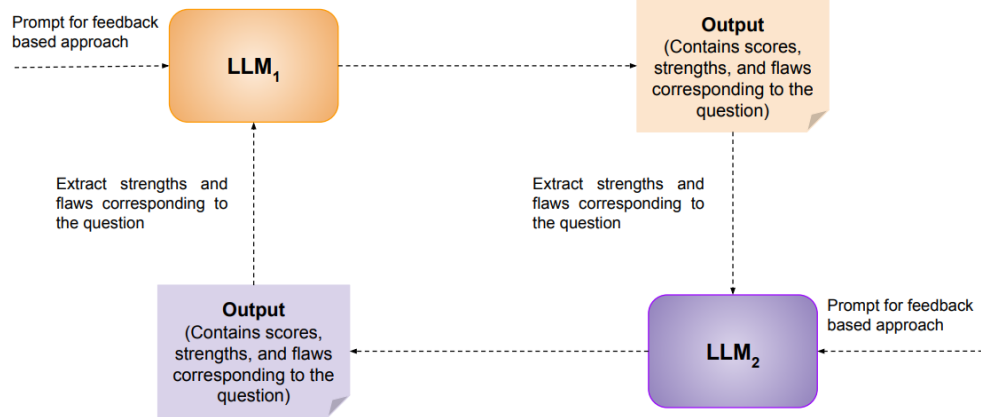    **end while**

---



Figure 5: An overview of the proposed approach called MIRROR.

## 4.4 Human Baseline Scores and Pearson's Correlation Coefficient

We recruit three educators and ask them to evaluate the questions generated by prompting GPT-3.5 Turbo in terms of *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity* to generate the human baseline scores corresponding to EduProbe and SciQ datasets. The three educators were asked to provide scores for each metric on a scale of 1 to 5. The scores given by the educators for each metric were then averaged. The inter-annotator agreement is calculated in terms of Fleiss's Kappa. The inter-annotator agreement between the three experts is 0.67, 0.64, 0.66, 0.45, and 0.59 for *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity*, respectively. Later we also calculate the Pearson's correlation coefficient [Freedman et al., 2020] between the best-performing model (i.e., GPT-4) and human experts.

## 5 Results

**Human Evaluation Results and Correlation Analysis.** Table 1 shows the human evaluation metric scores for the EduProbe dataset. Table 2 shows the human evaluation metric scores for the SciQ dataset. Table 3 shows the Pearson's correlation coefficient scores for the EduProbe and

Figure 6: Feedback-based prompt for evaluating human evaluation metrics.

SciQ datasets between GPT-4 and human experts. We observe that the human baseline scores the highest in *grammaticality*, *relevance*, and *appropriateness*, but lower in *novelty* and *complexity* for both the EduProbe and SciQ datasets. GPT-4 performs closest to the human baseline, especially in the feedback-based approach, with high scores in *grammaticality* and *appropriateness*. Gemini consistently scores slightly lower than GPT-4 across all metrics and evaluation methods. Llama2-70b has the lowest scores among the LLMs explored but still performs relatively well, with its highest score in *grammaticality*. Overall, LLMs tend to score higher in the feedback-based approach compared to the direct approach, suggesting that feedback helps improve perceived performance. In summary, GPT-4 outperforms Gemini and Llama2-70b in both evaluation methods, with scores closest to human performance, especially in *grammaticality* and *appropriateness*. Gemini performs moderately well, while Llama2-70b shows relatively lower scores across the metrics. For the EduProbe and SciQ datasets, GPT-4 shows higher correlation scores with human experts when using the feedback-based approach compared to the direct approach across all metrics. The highest correlation score is for *grammaticality* with a feedback-based approach. For the feedback-based approach, *appropriateness* and *relevance* show moderate correlations, suggesting a reasonable alignment with human expert evaluations in these areas. We observe lower correlations in *novelty* and *complexity* because both approaches (i.e., feedback-based and direct) show lower correlations in these areas, with the lowest being *complexity* in the direct approach. The direct approach generally yields lower correlation scores across all metrics, highlighting the importance of feedback in improving GPT-4's alignment with human expert evaluations. In summary, the correlation scores suggest that GPT-4's evaluations are more aligned with human experts when using the feedback-based approach, particularly in *grammaticality*. The alignment is moderate in *appropriateness* and *relevance*, while it is weaker in *novelty* and *complexity*. Direct approach results in overall lower correlation scores across all metrics.

For the SciQ and EduProbe datasets, the feedback-based approach achieves a higher correlation with human evaluations in terms of *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity*, suggesting that iteratively refining the questions using LLM feedback better aligns with human judgments on grammatical correctness (See Table 3). For *appropriateness*, the feedback-based approach shows a modest improvement over the direct approach. This indicates that feedback helps in evaluating whether the questions are suitable and contextually appropriate more closely to how humans would rate them. The feedback-based approach significantly outperforms direct prompting in *relevance*, suggesting that iterative LLM feedback effectively enhances the model's ability to judge how relevant the questions are to the given context. The feedback-based method also achieves a higher correlation in assessing *novelty*. This indicates that feedback helps the model better understand and evaluate the uniqueness and originality of the questions. While the improvement in *complexity* evaluation is smaller, the feedback-based approach still surpasses direct prompting. This suggests that feedback helps the model slightly better assess the intricacy and difficulty of the questions. Overall, the higher correlation scores for the feedback-based approach across all metrics suggest that iteratively using LLM feedback improves the alignment of GPT-4's evaluations with human expert judgments, making it a more effective method for assessing the quality of generated questions.

Table 1: Human evaluation results on the EduProbe dataset for grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue indicates the highest metric values for the corresponding methods.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| Human Baseline | 4.95 | 4.97 | 4.48 | 3.98 | 3.10 |
| *Eduprobe (Direct Approach)* | | | | | |
| GPT-4 | 4.81 | 4.73 | 4.20 | 4.12 | 3.92 |
| Gemini | 4.61 | 4.51 | 4.02 | 4.03 | 3.88 |
| Llama2-70b | 4.38 | 4.20 | 3.84 | 4.01 | 3.88 |
| *EduProbe (Feedback-based Approach)* | | | | | |
| GPT-4 | 4.87 | 4.82 | 4.30 | 4.25 | 4.05 |
| Gemini | 4.72 | 4.64 | 4.14 | 4.10 | 4.00 |
| Llama2-70b | 4.60 | 4.62 | 4.08 | 4.06 | 3.83 |

Table 2: Human evaluation results on the SciQ dataset for grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue indicates the highest metric values for the corresponding methods.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| Human Baseline | 4.90 | 4.93 | 4.38 | 3.99 | 3.20 |
| *SciQ (Direct Approach)* | | | | | |
| GPT-4 | 4.70 | 4.44 | 4.03 | 4.01 | 3.74 |
| Gemini | 4.42 | 4.34 | 3.92 | 3.84 | 3.65 |
| Llama2-70b | 4.23 | 4.10 | 3.73 | 3.67 | 3.28 |
| *SciQ (Feedback-based Approach)* | | | | | |
| GPT-4 | 4.77 | 4.74 | 4.24 | 4.20 | 4.01 |
| Gemini | 4.64 | 4.58 | 4.08 | 4.04 | 3.93 |
| Llama2-70b | 4.58 | 4.55 | 3.94 | 3.91 | 3.80 |

**Error Analysis.** We conducted a human study using questions from the EduProbe and SciQ datasets, which were evaluated by each LLM using both feedback-based and direct approaches. We observed that the scores provided by human experts and the direct approach matched in 54%, 39%, 36%, 42%, and 46% of the cases for *grammaticality*, *relevance*, *appropriateness*, *complexity*, and *novelty*, respectively. We observed that the scores provided by human experts and the feedback-based approach matched in 67%, 64%, 62%, 55%, and 61% for *grammaticality*, *relevance*, *appropriateness*, *complexity*, and *novelty*, respectively. This suggests that the direct approach has significant differences compared to human experts, and the use of the feedback-based approach significantly improves the results, especially in terms of *relevance* and *appropriateness*. So, *relevance* and *appropriateness* are two areas where our feedback-based approach (i.e., MIRROR) significantly improves, in addition to moderate improvements in other metrics.

# 6 Examples of Outputs Generated by LLMs Using the Feedback-based Approach

The feedback approach consists of feedback occurring between two LLMs, namely $LLM_1$ and $LLM_2$. Figure 7 shows the output provided by $LLM_1$ on a generated question based on Economics framed from the EduProbe dataset. Figure 8 shows the output provided by $LLM_2$ on a generated question based on Economics framed from the EduProbe dataset. We observe that the scores provided by $LLM_1$ and $LLM_2$ for different human evaluation metrics have become the same after completion of the Feedback Approach.

Table 3: Pearson's correlation coefficient scores for the EduProbe and SciQ datasets between GPT-4 and human experts on grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue denotes the highest correlation values for a particular dataset and approach.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| *EduProbe* | | | | | |
| GPT-4 (Direct Approach) | 0.41 | 0.32 | 0.26 | 0.25 | 0.28 |
| GPT-4 (Feedback-based Approach) | 0.62 | 0.48 | 0.51 | 0.42 | 0.38 |
| *SciQ* | | | | | |
| GPT-4 (Direct Approach) | 0.40 | 0.36 | 0.30 | 0.32 | 0.30 |
| GPT-4 (Feedback-based Approach) | 0.65 | 0.44 | 0.55 | 0.44 | 0.34 |

```
Scores and Rationale for the Question:  "What does purchasing power parity
do?"

Grammaticality:  5/5
Strengths:  The question is well-formed and free from grammatical errors.
Flaws:  None.

Appropriateness:  4.5/5
Strengths:  The question is semantically correct and understandable.
Flaws:  The question is very general and could be more specific about PPP.

Relevance:  5/5
Strengths:  The question directly addresses the topic of the context, mak-
ing it highly relevant.
Flaws:  None.

Complexity:  2/5
Strengths:  The question is simple and easy to understand.
Flaws:  The simplicity of the question does not challenge the respondent
to think deeply about the topic.  It could be more complex.

Novelty:  2.5/5
Strengths:  The question is clear and directly related to the topic.
Flaws:  The question lacks originality and distinctiveness.  It is a basic
question that many people might ask when first learning about PPP.

Strengths in the Question Based on the Evaluation Scores:

Grammatical Correctness:  The question is well-formed and free from gram-
matical errors.

Relevance:  The question directly addresses the topic of the context, mak-
ing it highly relevant.

Appropriateness:  The question is semantically correct and understandable.

Flaws in the Question Based on the Evaluation Scores:

Complexity:  The question is quite simple and does not challenge the re-
spondent to think deeply about the topic.

Novelty:  The question is basic and lacks originality.  It could be im-
proved by asking something more specific or unique about PPP.
```

Figure 7: Output provided by $LLM_1$ for a generated question based on Economics from the EduProbe dataset.

The context involved in the question generation process is: *"Purchasing power parity (PPP) is an economic indicator that signifies the purchasing power of the currencies of various nations of the world against each other. It helps in comparing living standards between different countries and estimating economic productivity.".* The question generated by GPT-3.5 Turbo is: *"What does purchasing power parity do?"*

## 7   Conclusion

In this work, we propose a novel system, MIRROR (Multi-LLM Iterative Review and Response for Optimized Rating), which utilizes LLMs to automate the evaluation process for open-ended questions generated by AQG systems. Our experiments with state-of-the-art LLMs, including GPT-4, Gemini, and Llama2-70b, show that the scores for the human evaluation metrics namely *relevance*,

```
To evaluate the question "What does purchasing power parity do?" based on
the human evaluation metrics, we will consider the context and score the
question accordingly.  Here are the scores and the rationale for each:

Grammaticality:  5/5
Strengths:  The question is grammatically correct and straightforward.
Flaws:  None.

Appropriateness:  4.5/5
Strengths:  The question is semantically correct, asking about the function
of purchasing power parity (PPP).
Flaws:  The question could be more specific, such as inquiring about par-
ticular aspects or impacts of PPP.

Relevance:  5/5
Strengths:  The question is highly relevant to the context, directly asking
about purchasing power parity, which is the topic of the given context.
Flaws:  None.

Complexity:  2/5
Strengths:  The question is simple and easy to understand.
Flaws:  The simplicity of the question does not require much reasoning or
cognitive effort to answer.  It could be more complex to encourage deeper
thinking.

Novelty:  2.5/5
Strengths:  The question is clear and directly related to the topic.
Flaws:  The question lacks originality and distinctiveness, as it is a very
basic question that many people might ask when first learning about PPP.

Strengths in the Question Based on the Evaluation Scores:

Grammatical Correctness:  The question is well-formed and free from gram-
matical errors.

Relevance:  The question directly addresses the topic of the context, mak-
ing it highly relevant.

Appropriateness:  The question is semantically correct and understandable.

Flaws in the Question Based on the Evaluation Scores:

Complexity:  The question is quite simple and does not challenge the re-
spondent to think deeply about the topic.

Novelty:  The question is basic and lacks originality.  It could be im-
proved by asking something more specific or unique about PPP.
```

Figure 8: Output provided by $LLM_2$ for a generated question based on Economics from the EduProbe dataset.

*appropriateness*, *novelty*, *complexity*, and *grammaticality* improve when using the feedback-based approach called MIRROR and tend to be closer to the human baseline scores. We also observe that Pearson's correlation coefficient between GPT-4 and human experts improves when using MIRROR compared to using direct apparoach for evaluation. Our error analysis shows that *relevance* and *appropriateness* are two metrics where our proposed approach, MIRROR, significantly improves. The results demonstrate that LLMs have the potential to provide a scalable and effective alternative to human evaluation, offering a promising solution for assessing question quality in AQG systems. So far, we have focused on short and medium-sized contexts, whereas future work will focus on applying MIRROR to longer contexts.

# 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.

Said Al Faraby, Adiwijaya Adiwijaya, and Ade Romadhony. Review on neural question generation for education purposes. *International Journal of Artificial Intelligence in Education*, pages 1–38, 2023.

Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. The (un)suitability of automatic evaluation metrics for text simplification. *Computational Linguistics*, 47(4):861–889, December 2021. doi: 10.1162/coli_a_00418. URL https://aclanthology.org/2021.cl-4.28.

Manik Bhandari, Pranav Narayan Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. Re-evaluating evaluation in text summarization. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9347–9359, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.751. URL https://aclanthology.org/2020.emnlp-main.751.

Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik Narasimhan, Ameet Deshpande, and Bruno Castro da Silva. Rlhf deciphered: A critical analysis of reinforcement learning from human feedback for llms, 2024. URL https://arxiv.org/abs/2404.08555.

Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.870. URL https://aclanthology.org/2023.acl-long.870.

David Freedman, Robert Pisani, and Roger Purves. Statistics: Fourth international student edition. *WW Nort Co Httpswww Amaz ComStatistics-Fourth-Int-Stud-Free Accessed*, 22, 2020.

Huanli Gong, Liangming Pan, and Hengchang Hu. KHANQ: A dataset for generating deep questions in education. In Nicoletta Calzolari, Chu-Ren Huang, Hansaem Kim, James Pustejovsky, Leo Wanner, Key-Sun Choi, Pum-Mo Ryu, Hsin-Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi, Patrizia Paggio, Nianwen Xue, Seokhwan Kim, Younggyun Hahm, Zhong He, Tony Kyungil Lee, Enrico Santus, Francis Bond, and Seung-Hoon Na, editors, *Proceedings of the 29th International Conference on Computational Linguistics*, pages 5925–5938, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics. URL https://aclanthology.org/2022.coling-1.518.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/91edff07232fb1b55a505a9e9f6c0ff3-Paper-Conference.pdf.

Subhankar Maity, Aniket Deroy, and Sudeshna Sarkar. A novel multi-stage prompting approach for language agnostic mcq generation using gpt. In Nazli Goharian, Nicola Tonellotto, Yulan He, Aldo Lipani, Graham McDonald, Craig Macdonald, and Iadh Ounis, editors, *Advances in Information Retrieval*, pages 268–277, Cham, 2024a. Springer Nature Switzerland. ISBN 978-3-031-56063-7.

Subhankar Maity, Aniket Deroy, and Sudeshna Sarkar. Harnessing the power of prompt-based techniques for generating school-level questions using large language models. In *Proceedings of the 15th Annual Meeting of the Forum for Information Retrieval Evaluation*, FIRE '23, page 30–39, New York, NY, USA, 2024b. Association for Computing Machinery. ISBN 9798400716324. doi: 10.1145/3632754.3632755. URL https://doi.org/10.1145/3632754.3632755.

Nikahat Mulla and Prachi Gharpure. Automatic question generation: a review of methodologies, datasets, evaluation metrics, and applications. *Progress in Artificial Intelligence*, 12(1):1–32, 2023.

Preksha Nema and Mitesh M. Khapra. Towards a better metric for evaluating question generation systems. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3950–3959, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1429. URL https://aclanthology.org/D18-1429.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf.

Ehud Reiter. A structured review of the validity of BLEU. *Computational Linguistics*, 44(3):393–401, September 2018. doi: 10.1162/coli_a_00322. URL https://aclanthology.org/J18-3002.

Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. Generative language models for paragraph-level question generation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 670–688, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.42. URL https://aclanthology.org/2022.emnlp-main.42.

Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin, editors, *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology.org/W17-4413.

Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 59008–59033. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/b8c90b65739ae8417e61eadb521f63d5-Paper-Conference.pdf.

Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 10935–10950. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/23e6f78bdec844a9f7b6c957de2aae91-Paper-Conference.pdf.

Ruqing Zhang, Jiafeng Guo, Lu Chen, Yixing Fan, and Xueqi Cheng. A review on question generation from natural language text. *ACM Trans. Inf. Syst.*, 40(1), sep 2021. ISSN 1046-8188. doi: 10.1145/3468889. URL `https://doi.org/10.1145/3468889`.

# A Appendix

## A.1 Results

Table 4 shows the human evaluation metric scores for the EduProbe dataset. Table 5 shows the human evaluation metric scores for the SciQ dataset. Table 6 shows the Pearson's correlation coefficient scores for the EduProbe and SciQ datasets between GPT-4 and human experts. We observe that the human baseline scores the highest in *grammaticality*, *relevance*, and *appropriateness*, but lower in *novelty* and *complexity* for both the EduProbe and SciQ datasets. GPT-4 performs closest to the human baseline, especially in the feedback-based approach, with high scores in *grammaticality* and *appropriateness*. Gemini consistently scores slightly lower than GPT-4 across all metrics and evaluation methods. Llama2-70b has the lowest scores among the LLMs explored but still performs relatively well, with its highest score in *grammaticality*. Overall, LLMs tend to score higher in the feedback-based approach compared to the direct approach, suggesting that feedback helps improve perceived performance. In summary, GPT-4 outperforms Gemini and Llama2-70b in both evaluation methods, with scores closest to human performance, especially in *grammaticality* and *appropriateness*. Gemini performs moderately well, while Llama2-70b shows relatively lower scores across the metrics. For the EduProbe and SciQ datasets, GPT-4 shows higher correlation scores with human experts when using the feedback-based approach compared to the direct approach across all metrics. The highest correlation score is for *grammaticality* with a feedback-based approach. For the feedback-based approach, *appropriateness* and *relevance* show moderate correlations, suggesting a reasonable alignment with human expert evaluations in these areas. We observe lower correlations in *novelty* and *complexity* because both approaches (i.e., feedback-based and direct) show lower correlations in these areas, with the lowest being *complexity* in the direct approach. The direct approach generally yields lower correlation scores across all metrics, highlighting the importance of feedback in improving GPT-4's alignment with human expert evaluations. In summary, the correlation scores suggest that GPT-4's evaluations are more aligned with human experts when using the feedback-based approach, particularly in *grammaticality*. The alignment is moderate in *appropriateness* and *relevance*, while it is weaker in *novelty* and *complexity*. Direct approach results in overall lower correlation scores across all metrics.

For the SciQ and EduProbe datasets, the feedback-based approach achieves a higher correlation with human evaluations in terms of *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity*, suggesting that iteratively refining the questions using LLM feedback better aligns with human judgments on grammatical correctness. For *appropriateness*, the feedback-based approach shows a modest improvement over the direct approach. This indicates that feedback helps in evaluating whether the questions are suitable and contextually appropriate more closely to how humans would rate them. The feedback-based approach significantly outperforms direct prompting in relevance, suggesting that iterative LLM feedback effectively enhances the model's ability to judge how relevant the questions are to the given context. The feedback-based method also achieves a higher correlation in assessing *novelty*. This indicates that feedback helps the model better understand and evaluate the uniqueness and originality of the questions. While the improvement in complexity evaluation is smaller, the feedback-based approach still surpasses direct prompting. This suggests that feedback helps the model slightly better assess the intricacy and difficulty of the questions. Overall, the higher correlation scores for the feedback-based approach across all metrics suggest that iteratively using LLM feedback improves the alignment of GPT-4's evaluations with human expert judgments, making it a more effective method for assessing the quality of generated questions.

## A.2 Pearson's Correlation Coefficient between LLMs and Human Experts

Table 7 and Table 8 show the results of the Pearson correlation coefficient between LLMs and human experts for the EduProbe and SciQ datasets, respectively. For the EduProbe dataset, all three models show improved correlation scores across all metrics when the feedback-based approach is used compared to the direct approach. GPT-4 shows the most significant improvement in correlation scores with the feedback-based approach, indicating that it benefits the most from receiving feedback. GPT-4 consistently outperforms Gemini and Llama2-70b across all metrics and both approaches (i.e., direct and feedback-based). Gemini performs better than Llama2-70b in all metrics, both with the direct approach and feedback-based approach in terms of correlation scores. GPT-4 with the feedback-based method achieves the highest correlation scores across all metrics, making it the best-performing model in this comparison. The significant increase in scores with the feedback-based approach in all

Table 4: Human evaluation results on the EduProbe dataset for grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue indicates the highest metric values for the corresponding methods.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| Human Baseline | 4.95 | 4.97 | 4.48 | 3.98 | 3.10 |
| *Eduprobe (Direct Approach)* | | | | | |
| GPT-4 | 4.81 | 4.73 | 4.20 | 4.12 | 3.92 |
| Gemini | 4.61 | 4.51 | 4.02 | 4.03 | 3.88 |
| Llama2-70b | 4.38 | 4.20 | 3.84 | 4.01 | 3.88 |
| *EduProbe (Feedback-based Approach)* | | | | | |
| GPT-4 | 4.87 | 4.82 | 4.30 | 4.25 | 4.05 |
| Gemini | 4.72 | 4.64 | 4.14 | 4.10 | 4.00 |
| Llama2-70b | 4.60 | 4.62 | 4.08 | 4.06 | 3.83 |

Table 5: Human evaluation results on the SciQ dataset for grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue indicates the highest metric values for the corresponding methods.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| Human Baseline | 4.90 | 4.93 | 4.38 | 3.99 | 3.20 |
| *SciQ (Direct Approach)* | | | | | |
| GPT-4 | 4.70 | 4.44 | 4.03 | 4.01 | 3.74 |
| Gemini | 4.42 | 4.34 | 3.92 | 3.84 | 3.65 |
| Llama2-70b | 4.23 | 4.10 | 3.73 | 3.67 | 3.28 |
| *SciQ (Feedback-based Approach)* | | | | | |
| GPT-4 | 4.77 | 4.74 | 4.24 | 4.20 | 4.01 |
| Gemini | 4.64 | 4.58 | 4.08 | 4.04 | 3.93 |
| Llama2-70b | 4.58 | 4.55 | 3.94 | 3.91 | 3.80 |

models indicates that providing feedback improves the correlation between human experts and the LLMs.

For the SciQ dataset, starting with GPT-4, the model performs moderately well in terms of correlation scores with the direct approach. However, when the feedback-based approach is applied, there is a noticeable improvement in performance, especially in terms of *grammaticality*, *appropriateness*, and *relevance*. The Gemini model shows a similar trend: its performance under direct prompting is moderate, but it improves with the feedback-based approach. The enhancements are particularly evident in *grammaticality*, *appropriateness*, and *relevance*, though not as pronounced as those seen in GPT-4. Nevertheless, Gemini performs better with the feedback-based approach than with the direct approach in terms of correlation scores. Llama2-70b performs moderately with the direct approach, but feedback-based prompting leads to better correlation scores across all metrics, with the most significant improvements observed in *grammaticality*, *relevance*, and *appropriateness*.

### A.3 Error Analysis

Table 9 shows the percentage of questions with exact matches in scores for different human evaluation metrics using both the direct approach and the feedback-based approach. We conducted a human study based on 100 questions from the EduProbe and SciQ datasets, respectively evaluated by each LLM using both approaches. We observed that the scores provided by human experts matched with the direct approach in 54%, 39%, 36%, 42%, and 46% of cases for *grammaticality*, *appropriateness*, *relevance*, *novelty*, and *complexity*, respectively, for GPT-4. In contrast, the scores provided by human

Table 6: Pearson's correlation coefficient scores for the EduProbe and SciQ datasets between GPT-4 and human experts on grammaticality (Gram), appropriateness (App), relevance (Rel), novelty (Nov), and complexity (Com). Blue denotes the highest correlation values for a particular dataset and approach.

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| *EduProbe* | | | | | |
| GPT-4 (Direct Approach) | 0.41 | 0.32 | 0.26 | 0.25 | 0.28 |
| GPT-4 (Feedback-based Approach) | 0.62 | 0.48 | 0.51 | 0.42 | 0.38 |
| *SciQ* | | | | | |
| GPT-4 (Direct Approach) | 0.40 | 0.36 | 0.30 | 0.32 | 0.30 |
| GPT-4 (Feedback-based Approach) | 0.65 | 0.44 | 0.55 | 0.44 | 0.34 |

Table 7: Pearson's correlation coefficient scores for the EduProbe dataset between LLMs and human experts on *grammaticality* (Gram), *appropriateness* (App), *relevance* (Rel), *novelty* (Nov), and *complexity* (Com). Blue denotes the highest correlation values among all approaches and LLMs. Underline denotes the better correlation values for a particular LLM among the two approaches (i.e., direct and feedback-based).

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| GPT-4 (Direct Approach) | 0.41 | 0.32 | 0.26 | 0.25 | 0.28 |
| GPT-4 (Feedback-based Approach) | 0.62 | 0.48 | 0.51 | 0.42 | 0.38 |
| Gemini (Direct Approach) | 0.40 | 0.30 | 0.23 | 0.22 | 0.24 |
| Gemini (Feedback-based Approach) | 0.58 | 0.45 | 0.48 | 0.39 | 0.35 |
| Llama2-70b (Direct Approach) | 0.38 | 0.27 | 0.22 | 0.20 | 0.21 |
| Llama2-70b (Feedback-based Approach) | 0.54 | 0.44 | 0.44 | 0.35 | 0.31 |

Table 8: Pearson's correlation coefficient scores for the SciQ dataset between LLMs and human experts on *grammaticality* (Gram), *appropriateness* (App), *relevance* (Rel), *novelty* (Nov), and *complexity* (Com). Blue denotes the highest correlation values among all approaches and LLMs. Underline denotes the better correlation values for a particular LLM among the two approaches (i.e., direct and feedback-based).

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| GPT-4 (Direct Approach) | 0.40 | 0.36 | 0.30 | 0.32 | 0.30 |
| GPT-4 (Feedback-based Approach) | 0.65 | 0.44 | 0.55 | 0.44 | 0.34 |
| Gemini (Direct Approach) | 0.38 | 0.33 | 0.28 | 0.30 | 0.27 |
| Gemini (Feedback-based Approach) | 0.62 | 0.40 | 0.51 | 0.40 | 0.30 |
| Llama2-70b (Direct Approach) | 0.36 | 0.30 | 0.25 | 0.27 | 0.24 |
| Llama2-70b (Feedback-based Approach) | 0.60 | 0.37 | 0.48 | 0.36 | 0.28 |

experts matched with the feedback-based approach in 67%, 64%, 62%, 55%, and 61% of cases for the same metrics. All three LLMs (i.e., GPT-4, Gemini, and Llama2-70b) show a noticeable improvement in their scores when using the feedback-based approach compared to the direct approach, indicating that iterative feedback helps the models perform better. GPT-4 shows the highest scores in both direct and feedback-based methods, making it the most effective model among the three. Gemini and Llama2-70b exhibit similar trends, but Gemini generally scores higher than Llama2-70b in both approaches. For all three LLMs, we observed that *appropriateness* and *relevance* showed the highest improvements with the feedback-based approach compared to the direct approach. These results indicate that providing feedback is particularly beneficial for appropriateness and relevance, significantly enhancing the performance of the LLMs.

## A.4   Examples of Outputs Generated by LLMs Using the Feedback-based Approach

The feedback approach consists of feedback occurring between two LLMs, namely $LLM_1$ and $LLM_2$. Figure 9 shows the output provided by $LLM_1$ on a generated question based on Economics framed from the EduProbe dataset. Figure 10 shows the output provided by $LLM_2$ on a generated question based on Economics framed from the EduProbe dataset. We observe that the scores provided by $LLM_1$ and $LLM_2$ for different human evaluation metrics have become the same after completion of the Feedback Approach. The context involved in the question generation process is: *"Purchasing power parity (PPP) is an economic indicator that signifies the purchasing power of the currencies of various nations of the world against each other. It helps in comparing living standards between different countries and estimating economic productivity.".* The question generated by GPT-3.5 Turbo is: *"What does purchasing power parity do?"*

Table 9: Percentage of questions showing the exact match in scores for different human evaluation metrics for the direct approach and feedback-based approach. Blue denotes the highest percentage values among all approaches and LLMs. Underline denotes the better percentage values for a particular LLM among the two approaches (i.e., direct and feedback-based).

| Model | Gram | App | Rel | Nov | Com |
|---|---|---|---|---|---|
| GPT-4 (Direct Approach) | 54 | 39 | 36 | 42 | 46 |
| GPT-4 (Feedback-based Approach) | 67 | 64 | 62 | 55 | 61 |
| Gemini (Direct Approach) | 51 | 36 | 33 | 38 | 41 |
| Gemini (Feedback-based Approach) | 65 | 60 | 57 | 50 | 56 |
| Llama2-70b (Direct Approach) | 47 | 32 | 30 | 34 | 35 |
| Llama2-70b (Feedback-based Approach) | 62 | 57 | 53 | 45 | 53 |

Figure 11 shows the output provided by $LLM_1$ on a generated question based on History framed from the EduProbe dataset. Figure 12 shows the output provided by $LLM_2$ on a generated question based on History framed from the EduProbe dataset. We observe that the scores provided by $LLM_1$ and $LLM_2$ for different human evaluation metrics have become same after completion of the Feedback Approach. The context involved in the question generation process is: *"During the medieval period in India, Islamic rulers held significant power, leading to the blending of Indian and Islamic cultures, which can still be observed in the architecture and artwork created at that time. The country was governed and administered by notable rulers who made exceptional contributions in the fields of art, culture, literature, and architecture during this period."* The question generated by GPT-3.5 Turbo is: *"How did medieval period contribute to India?"*

Figure 13 shows the output provided by $LLM_1$ on a generated question based on Biology framed from the SciQ dataset. Figure 14 shows the output provided by $LLM_2$ on a generated question based on Biology framed from the SciQ dataset. We observe that the scores provided by $LLM_1$ and $LLM_2$ for different human evaluation metrics have become same after completion of the Feedback Approach. The context is: "*Mesophiles grow best in moderate temperature, typically between 25'C and 40'C. Mesophiles are often found living in or on the bodies of humans or other animals. The optimal growth temperature of many pathogenic mesophiles is 37".* The question generated by GPT-3.5 Turbo is: "*Where is Mesophiles found?*"

Figure 15 shows the output provided by $LLM_1$ on a generated question based on Earth Science framed from the SciQ dataset. Figure 16 shows the output provided by $LLM_2$ on a generated question based on Earth science framed from the SciQ dataset. We observe that the scores provided by $LLM_1$ and $LLM_2$ for different human evaluation metrics have become same after completion of the Feedback Approach. The context is: "*Without Coriolis Effect the global winds would blow north to south or south to north. But Coriolis makes them blow north-east to southwest or the reverse in the Northern Hemisphere. The winds blow northwest to southeast or the reverse in the southern hemisphere.*" The question generated by GPT-3.5 Turbo is: "*How does Coriolis effect impact global winds?*"

## A.5 Implementation Details

For GPT-3.5 Turbo we use a `temperature` value of 0.7 and `max_tokens` size of 2048. For GPT-4 we use a `temperature` value of 0.7 and `max_tokens` size of 2048. For Llama2-70b we use a `temperature` value of 0.8 and `max_tokens` size of 2048. For Gemini we use a `temperature` value of 0.7 and `max_tokens` size of 2048.

```
Scores and Rationale for the Question:  "What does purchasing power
parity do?"

Grammaticality:  5/5
Strengths:  The question is well-formed and free from grammatical
errors.
Flaws:  None.

Appropriateness:  4.5/5
Strengths:  The question is semantically correct and understandable.
Flaws:  The question is very general and could be more specific
about PPP.

Relevance:  5/5
Strengths:  The question directly addresses the topic of the con-
text, making it highly relevant.
Flaws:  None.

Complexity:  2/5
Strengths:  The question is simple and easy to understand.
Flaws:  The simplicity of the question does not challenge the re-
spondent to think deeply about the topic.  It could be more complex.

Novelty:  2.5/5
Strengths:  The question is clear and directly related to the topic.
Flaws:  The question lacks originality and distinctiveness.  It is a
basic question that many people might ask when first learning about
PPP.

Strengths in the Question Based on the Evaluation Scores:

Grammatical Correctness:  The question is well-formed and free from
grammatical errors.

Relevance:  The question directly addresses the topic of the con-
text, making it highly relevant.

Appropriateness:  The question is semantically correct and under-
standable.

Flaws in the Question Based on the Evaluation Scores:

Complexity:  The question is quite simple and does not challenge the
respondent to think deeply about the topic.

Novelty:  The question is basic and lacks originality.  It could be
improved by asking something more specific or unique about PPP.
```

Figure 9: Output provided by $LLM_1$ for a generated question based on Economics from the EduProbe dataset.

```
To evaluate the question "What does purchasing power parity do?" based
on the human evaluation metrics, we will consider the context and
score the question accordingly.  Here are the scores and the rationale
for each:


Grammaticality:  5/5
Strengths:  The question is grammatically correct and straightforward.
Flaws:  None.


Appropriateness:  4.5/5
Strengths:  The question is semantically correct, asking about the
function of purchasing power parity (PPP).
Flaws:  The question could be more specific, such as inquiring about
particular aspects or impacts of PPP.


Relevance:  5/5
Strengths:  The question is highly relevant to the context, directly
asking about purchasing power parity, which is the topic of the given
context.
Flaws:  None.


Complexity:  2/5
Strengths:  The question is simple and easy to understand.
Flaws:  The simplicity of the question does not require much reasoning
or cognitive effort to answer.  It could be more complex to encourage
deeper thinking.


Novelty:  2.5/5
Strengths:  The question is clear and directly related to the topic.
Flaws:  The question lacks originality and distinctiveness, as it is
a very basic question that many people might ask when first learning
about PPP.


Strengths in the Question Based on the Evaluation Scores:


Grammatical Correctness:  The question is well-formed and free from
grammatical errors.


Relevance:  The question directly addresses the topic of the context,
making it highly relevant.


Appropriateness:  The question is semantically correct and understand-
able.


Flaws in the Question Based on the Evaluation Scores:


Complexity:  The question is quite simple and does not challenge the
respondent to think deeply about the topic.


Novelty:  The question is basic and lacks originality.  It could be
improved by asking something more specific or unique about PPP.
```

Figure 10: Output provided by $LLM_2$ for a generated question based on Economics from the EduProbe dataset.

```
To evaluate the question "How did medieval period contribute to India?"
based on the human evaluation metrics, we will consider the context and
score the question accordingly.  Here are the scores and the rationale
for each:

Grammaticality:  3.5/5
Strengths:  The question is mostly grammatically correct.
Flaws:  The article "the" is missing before "medieval period," which
makes the question slightly awkward.  It should read:  "How did the me-
dieval period contribute to India?"

Appropriateness:  4.5/5
Strengths:  The question is semantically correct, asking about the con-
tributions of the medieval period.
Flaws:  The wording could be more precise, such as specifying which as-
pects of India (e.g., culture, art, architecture).

Relevance:  5/5
Strengths:  The question is highly relevant to the context, directly in-
quiring about the contributions of the medieval period, which includes
significant cultural, artistic, and administrative aspects.
Flaws:  None.

Complexity:  4/5
Strengths:  The question requires a detailed and thoughtful response,
considering various contributions in multiple fields (art, culture, lit-
erature, architecture).
Flaws:  The question could be more complex by asking about specific con-
tributions or comparing them to those of other periods.

Novelty:  3.5/5
Strengths:  The question is clear and prompts a discussion on an impor-
tant historical period.
Flaws:  While the question is pertinent, it is somewhat general and could
be seen as a common inquiry in the study of history.

Strengths in the Question Based on the Evaluation Scores:

Relevance:  The question directly addresses the topic of the context,
making it highly relevant.

Appropriateness:  The question is semantically correct and understand-
able.

Complexity:  The question prompts a detailed response, requiring consid-
eration of multiple aspects of the medieval period.

Flaws in the Question Based on the Evaluation Scores:

Grammaticality:  The question is slightly awkward due to the missing ar-
ticle "the" before "medieval period."

Novelty:  The question, while clear and pertinent, is somewhat general
and lacks a unique angle.  It could be improved by focusing on specific
contributions or comparing the medieval period's contributions to those
of other periods.

Complexity:  Although reasonably complex, it could be further detailed to
enhance the level of reasoning or cognitive effort required to answer.
```

Figure 11: Output provided by $LLM_1$ for a generated question based on History from the EduProbe dataset.

```
To evaluate the question "How did medieval period contribute to India?"
based on the human evaluation metrics, we will consider the context and
score the question accordingly.  Here are the scores and rationale for
each metric:

Grammaticality:  3.5/5
Strengths:  The question is mostly grammatically correct.
Flaws:  The article "the" is missing before "medieval period," making
the question slightly awkward.  It should read:  "How did the medieval
period contribute to India?"

Appropriateness:  4.5/5
Strengths:  The question is semantically correct, asking about the con-
tributions of the medieval period.
Flaws:  The wording could be more precise, such as specifying which as-
pects of India (e.g., culture, art, architecture).

Relevance:  5/5
Strengths:  The question is highly relevant to the context, directly
inquiring about the contributions of the medieval period, which includes
significant cultural, artistic, and administrative aspects.
Flaws:  None.

Complexity:  4/5
Strengths:  The question requires a detailed and thoughtful response,
considering various contributions in multiple fields (art, culture, lit-
erature, architecture).
Flaws:  The question could be more complex by asking about specific con-
tributions or comparing them to those of other periods.

Novelty:  3.5/5
Strengths:  The question is clear and prompts a discussion on an impor-
tant historical period.
Flaws:  While the question is pertinent, it is somewhat general and
could be seen as a common inquiry in the study of history.

Strengths in the Question Based on the Evaluation Scores:

Relevance:  The question is highly pertinent to the context, directly
addressing the main topic.

Appropriateness:  The question is clear and semantically accurate.

Complexity:  The question encourages a detailed response, requiring con-
sideration of various aspects of the medieval period.

Flaws in the Question Based on the Evaluation Scores:

Grammaticality:  The question is slightly awkward due to the missing
article "the" before "medieval period."

Novelty:  While clear and relevant, the question is somewhat general
and lacks a unique angle.  It could be improved by focusing on specific
contributions or comparing the medieval period's contributions to those
of other periods.

Complexity:  Although reasonably complex, the question could be more
detailed to enhance the level of reasoning or cognitive effort required
to answer.
```

Figure 12: Output provided by $LLM_2$ for a generated question based on History from the EduProbe dataset.

```
Evaluation Scores:

Grammaticality:  5/5
Appropriateness:  5/5
Relevance:  5/5
Complexity:  3/5
Novelty:  4/5

Strengths in the question based on the evaluation scores provided above:

Grammaticality:  The question is grammatically correct, with proper sentence
structure and punctuation.

Appropriateness:  The question is semantically correct, making sense in the
context of the given information.

Relevance:  The question is highly relevant to the context provided, di-
rectly addressing the key point about the optimal growth temperature of
mesophiles.

Novelty:  The question is somewhat original and distinct, asking both for
the optimal growth temperature and the reason for its significance, which
adds depth to the inquiry.

Flaws in the question based on the evaluation scores provided above:

Complexity:  The question is relatively straightforward and does not require
a high level of reasoning or cognitive effort to answer.  A more complex
question might delve deeper into the implications of the optimal growth tem-
perature or its effects on human health.
```

Figure 13: Output provided by $LLM_1$ for a generated question based on Biology from the SciQ
dataset.

```
Evaluation Scores:

Grammaticality:  5/5
Appropriateness:  5/5
Relevance:  5/5
Complexity:  3/5
Novelty:  4/5

Strengths in the question based on the evaluation scores provided above:

Grammaticality:  The question is grammatically accurate, with correct sentence
structure and punctuation.

Appropriateness:  The question is semantically accurate and fits well with the
given context.

Relevance:  The question is highly relevant to the provided context, directly
addressing the key point about the optimal growth temperature of mesophiles.

Novelty:  The question is somewhat original and unique, as it asks not only
about the optimal growth temperature but also why it is significant, adding
depth to the inquiry.

Flaws in the question based on the evaluation scores provided above:


Complexity:  The question is relatively simple and does not require a high
level of reasoning or cognitive effort to answer.  A more complex question
could explore the implications of the optimal growth temperature or its ef-
fects on human health.
```

Figure 14: Output provided by $LLM_2$ for a generated question based on Biology from the SciQ dataset.

```
Evaluation Scores:

Grammaticality:  5/5
Appropriateness:  5/5
Relevance:  4/5
Complexity:  2/5
Novelty:  2/5

Strengths in the question based on the evaluation scores provided above:

Grammaticality:  The question is free of grammatical errors, with proper
sentence structure and punctuation.

Appropriateness:  The question is semantically accurate and fits the con-
text.

Relevance:  The question is pertinent to the context, as it asks about the
Coriolis effect mentioned in the provided information.

Flaws in the question based on the evaluation scores provided above:

Complexity:  The question is very straightforward and does not require
significant reasoning or cognitive effort to answer.  It could be more
challenging by asking for a detailed explanation or the implications of the
Coriolis effect.

Novelty:  The question lacks originality and distinctiveness, as it is a
basic inquiry that doesn't explore unique aspects or implications of the
Coriolis effect mentioned in the context.
```

Figure 15: Output provided by $LLM_1$ for a generated question based on Earth Science from the SciQ dataset.

```
Evaluation Scores:

Grammaticality:  5/5
Appropriateness:  5/5
Relevance:  4/5
Complexity:  2/5
Novelty:  2/5


Strengths in the question based on the evaluation scores provided above:

Grammaticality:  The question is grammatically correct, with no errors
in sentence structure or punctuation.

Appropriateness:  The question is semantically correct and makes sense
given the context.

Relevance:  The question is relevant to the context, as it directly asks
about the Coriolis effect mentioned in the provided information.

Flaws in the question based on the evaluation scores provided above:

Complexity:  The question is very simple and does not require much rea-
soning or cognitive effort to answer.  It could be improved by asking
for more detailed information or explanation.

Novelty:  The question lacks originality and distinctiveness, as it is a
very basic inquiry that does not delve into the unique aspects or impli-
cations of the Coriolis effect described in the context.
```

Figure 16: Output provided by $LLM_2$ for a generated question based on Earth Science from the SciQ dataset.