

# “Did my figure do justice to the answer?”: Towards Multimodal Short Answer grading with Feedback (MMSAF)

Anonymous ACL submission

## Abstract

Assessments play a vital role in a student’s learning process. This is because they provide valuable feedback crucial to a student’s growth. Such assessments contain questions with open-ended responses, which are difficult to grade at scale. These responses often require students to express their understanding through textual and visual elements together as a unit. In order to develop scalable assessment tools for such questions, one needs multimodal LLMs having strong comparative reasoning capabilities across multiple modalities. Thus, to facilitate research in this area, we propose the Multimodal Short Answer grading with Feedback (MMSAF) problem along with a dataset of 2,197 data points. Additionally, we provide an automated framework for generating such datasets. As per our evaluations, existing Multimodal Large Language Models (MLLMs) could predict whether an answer is correct, incorrect or partially correct with an accuracy of 55%. Similarly, they could predict whether the image provided in the student’s answer is relevant or not with an accuracy of 75%. As per human experts, Pixtral was more aligned towards human judgement and values for biology and ChatGPT for physics and chemistry and achieved a score of 4 or more out of 5 in most parameters.

## 1 Introduction

Assessments play a vital role in a student’s learning process as they provide valuable feedback (Deeva et al., 2021), crucial to a student’s growth. Moreover, corrective, motivational and informative feedback can drastically speed up a student’s learning process and help the student develop an innate curiosity.

Grading such assessments and providing individual feedback to students is often difficult, especially in classrooms with a low teacher-to-student ratio (Burrows et al., 2015). Such assessments

contain open-ended responses, which are often difficult to grade at scale. Responses to such questions require students to express their understanding through textual and visual elements, leading to deeper levels of cognitive engagement.

This leads to the question: Can we develop scalable assessment tools that can assist teachers in evaluating such questions with open-ended responses while providing quality feedback? Developing such tools requires multimodal LLMs (MLLMs) capable of reasoning over different modalities such as text and images. These MLLMs should be capable of identifying key concepts across both text and images in the reference answer and comparing them to those in the student answer. Moreover, to develop such MLLMs, researchers need access to reliable datasets that are representative of student answers in examinations.

To facilitate research towards developing such systems, we propose the problem of Multimodal Short Answer grading with Feedback (MMSAF) along with a dataset of 2,197 data points. The primary motivation of this work is to provide effective feedback to the student on the shortcomings of their response and ways to mitigate them if needed.

Our contributions are as follows -

1. Introduction of the MMSAF problem, along with a dataset of 2,197 instances. (Section 3, Section 4)
2. An automated framework to generate an MMSAF dataset for any set of questions and reference answers. (Section 4)
3. A rubric-based approach to evaluate the quality of feedback coupled with extensive zero-shot evaluation on existing Large Language Models (LLMs). As per our evaluations, existing Multimodal Large Language Models (MLLMs), we achieve the following (Section 6)-

- (a) MLLMs could predict whether an answer is correct, incorrect or partially correct with an accuracy of 55%.
- (b) MLLMs could predict whether the image provided in the student’s answer is relevant or not with an accuracy of 75%.
- (c) Per human experts, Pixtral was more aligned towards human judgement and values for biology and ChatGPT for physics and chemistry and achieved a score of 4 or more out of 5 in most parameters.

## 2 Related Work

In recent years, there has been growing interest from both the natural language processing (NLP) and education research communities in the task of Automatic Short Answer Grading (ASAG) with feedback. A notable milestone in this direction came in 2022, when [Filighera et al. \(2022\)](#) introduced the first dataset for ASAG with feedback problem. This bilingual dataset focused on short textual responses to questions across various topics, primarily in computer science. However, the dataset was limited to only about 2,000 responses, and lacked diversity across different engineering disciplines.

To address these shortcomings, [Aggarwal et al. \(2024\)](#) introduced the EngSAF dataset, which contained about 5,800 student answers drawn from 25 different engineering courses spanning multiple subfields. This dataset laid the groundwork for more robust benchmarking and model development in the ASAG with feedback task.

Following this, research began to shift towards more advanced methods for the ASAG with feedback problem. While earlier approaches primarily leveraged prompt engineering, [Fateen et al. \(2024\)](#) proposed a retrieval-augmented generation (RAG) based approach to enhance response quality and contextual relevance.

The feedback so generated by the model serves as a way to explain its grading rationale, which adds a layer of explainability to the grading task. [Li et al. \(2023\)](#) introduced the Automated Explainable Student Response Assessment (AERA) framework, which generates scoring rationales using ChatGPT. AERA demonstrated rationales comparable in quality to human explanations and achieved a Quadratic Weighted Kappa (QWK) score of 11%

on benchmark datasets. Similarly, [Tornqvist et al. \(2023\)](#) presented ExASAG, an explainable grading framework that integrates SHAP (Shapley Additive exPlanations) with SciBERT to introduce interpretability to the process of grading such assignments.

While these efforts significantly advanced the field of text-based ASAG, they fall short in handling multimodal responses, especially answers that combine textual explanations with supporting visual models or diagrams. Visual models, as part of student answers, play an important role in demonstrating the student’s proficiency level on a particular concept. To address this, [Leong et al. \(2018\)](#) and later [Sagherian et al. \(2022\)](#) proposed automated grading systems for scientific visual models, evaluated on proprietary datasets from Educational Testing Services (ETS). However, these visual models were created using predefined shapes such as boxes, arrows, and fish, limiting student expressiveness and creativity. Similarly, [Lee and Zhai \(2023\)](#) showed that GPT-4V is capable of evaluating such visual models. However, none of these works deal with the scenario of evaluating hand-drawn diagrams or diagrams that are far more complex than simple visual models.

To bridge this gap and handle real-world scenarios where students provide multimodal responses that contain textual answers coupled with diagrams, we introduce the MMSAF problem in the next section.

## 3 The Multimodal Short Answer Grading with Feedback (MMSAF) Problem

We introduce the Multimodal Short Answer Grading with Feedback (MMSAF) problem, which focuses on evaluating student responses containing both textual and visual content. Given a Question (Q), a Reference Answer (RA), and a Student Answer (SA), the objective is to assign a Level of Correctness (LC) label, an Image Relevance (IR) label, and a Feedback (F) which provides rationale behind the LC and IR labels.

Figure 1 provides an illustrative example. In the given question, the student is asked to explain the flow of blood in the human heart using a diagram. While a textual response alone can convey the explanation, a supporting diagram adds clarity and depth to the response. Note that it is difficult to create such diagrams using only simple visual models. This highlights the motivation for introducing this

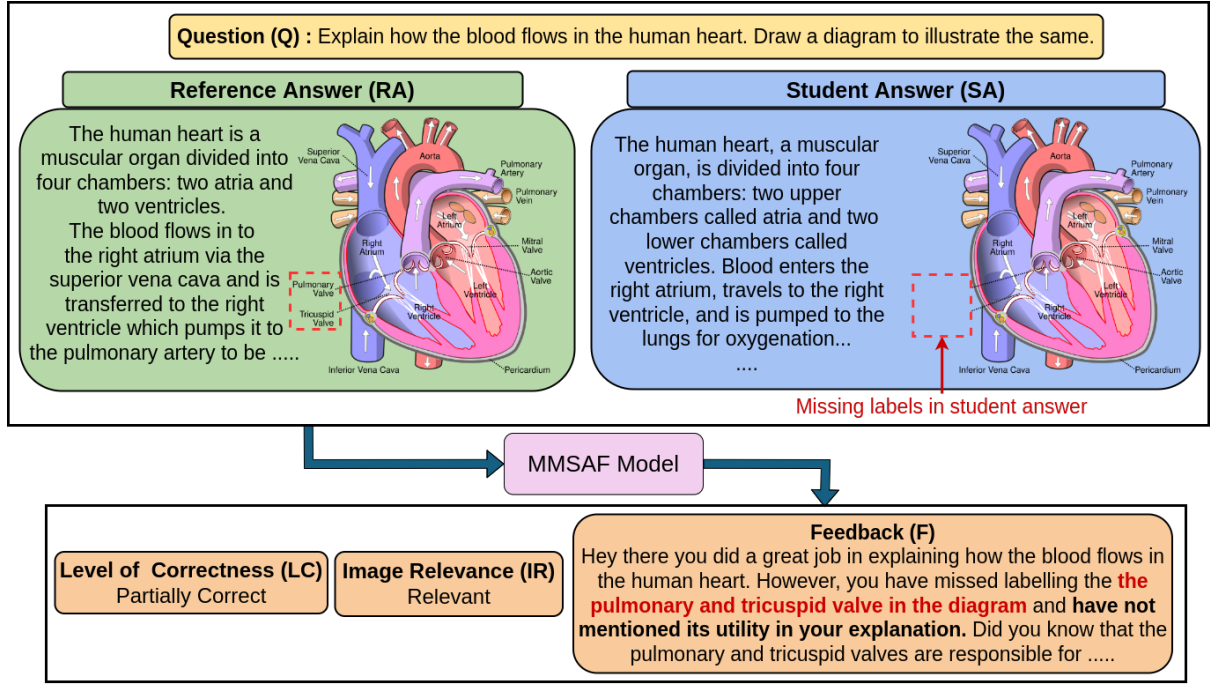


Figure 1: Illustration of the MMSAF problem with an example. (Image source for heart diagram: <https://edurev.in/t/131714/STRUCTURE-OF-HUMAN-HEART>)

problem.

We decompose the MMSAF problem into two sub-problems: a classification task for LC and IR labels and a reasoning task for feedback generation.

### 3.1 Classification of Level of Correctness and Image Relevance

The problem of determining LC is formulated as a multi-class classification problem where given  $(Q, SA, RA)$ , the model  $M$  must assign one of three correctness levels:  $(Q, SA, RA) \xrightarrow{M} \{\text{Correct, Partially Correct, Incorrect}\}$ , where the LC reflects how accurately the student's response aligns with the reference.

Similarly, the IR problem is framed as a binary classification problem where given  $(Q, SA, RA)$ , the model determines whether the image in the student response is relevant or not:  $(Q, SA, RA) \xrightarrow{M} \{\text{Relevant, Irrelevant}\}$ . This sub-problem relies heavily on the model's ability to perform multimodal reasoning.

### 3.2 Feedback Generation

The problem of feedback generation in MMSAF requires comparative reasoning (Yu et al., 2023), wherein the model compares the student and reference answers to identify conceptual matches and deviations. This involves verifying whether key concepts present in the textual and visual part of he

RA have been captured in the SA or not.

Thus, given  $(Q, RA, SA)$ , a model  $M$  must generate feedback:  $(Q, SA, RA) \xrightarrow{M} \text{Feedback}$ .

The feedback should identify the errors in the SA and, where necessary, suggest methods to correct them. Since both SA and RA contain multiple modalities, the model has to perform comparative reasoning across different modalities. In the next section, we introduce the MMSAF dataset, developed to support and benchmark progress on this task.

## 4 Multimodal Short Answer Grading with Feedback (MMSAF) Dataset

The MMSAF dataset serves as a benchmark for the MMSAF problem. It consists of 181 high school-level questions across physics, chemistry, and biology. Synthetic student responses were generated for each, resulting in 2,197 total data points. Each data point is a tuple of seven elements: Question (Q), Reference Answer (RA), Student Answer (SA), Level of correctness (LC), Image Relevance (IR), Sample Feedback (F) and Rubrics for error detection in feedback (FR).

As per the problem, given a  $(Q, RA, SA)$  tuple, the task is to generate feedback that evaluates both the correctness and image relevance of the student's answer. Since the student answers are synthetically

generated, we record the errors present in them as rubrics under the FR column. These can then be used to evaluate the generated feedback along with standard metrics like BLEU (Papineni et al., 2002) or ROUGE (Lin, 2004).

Figure 1 provides an example from the dataset<sup>1</sup>. The corresponding FR column will contain rubrics such as “Did it detect that the pulmonary and tricuspid valves have not been labelled in the diagram?” while F will have the following value “The answer is partially correct. The student does not label the pulmonary and tricuspid valves. The student fails to mention the utility of the pulmonary and the tricuspid valve....”.

Next, we describe how the textual and image components of student answers were generated, followed by the process of assigning correctness and relevance levels, feedback, and rubrics. Figure 2 illustrates the overall data generation pipeline.

#### 4.1 Generation of Textual and Image Segments of Student Answers

We first extract 160 question-reference answer pairs from high school textbooks and generate 21 additional ones using Gemini (Team, 2024), verified by a Subject Matter Expert (SME) (Step 1 in Figure 2). Each student’s response contains a textual answer and a supporting image, generated separately and then combined using a correctness matrix (Step 2 in Figure 2).

In consultation with SMEs and based on Marsh and Eliseev (2019), student mistakes are categorized as:

- *Errors made with confidence*: Answers which appear as a confident attempt but are actually incorrect answers derived from misunderstood or fabricated facts.
- *Misunderstanding*: Incorrect answers due to misinterpreting the question or confusing related concepts.
- *Conceptual Change*: Misapplication of a known concept in a new context.

More examples are provided in Appendix M. These errors align with hallucinations commonly seen in LLMs, particularly factual fabrication and inconsistency (Huang et al., 2025). Hence, we simulate such student answers by making use of hallucinations.

All textual answers are generated using Gemini, for which we do the following -

- *Correct Answers*: For each question, we prompt Gemini (see Appendix A) with the question and reference answer, except for 26 numerical questions, where the correct answer has been reused.
- *Incorrect/Partially Correct Answers*: Using the termite strategy (Saxena, 2024), we introduce hallucinations into the reference answer to simulate common student errors. Prompts are in Appendices B and C. For numerical questions, calculation errors in the textual part and minor perturbations, such as changes in the direction of force in free body diagrams, have been added manually.

For images, SMEs noted common issues like missing labels or use of incorrect domain-specific objects (E.g., using concave instead of a convex lens), which have been simulated as follows -

- *Correct Images*: Directly taken from the reference answer.
- *Partially Correct Images*: Created by removing parts, swapping labels, or replacing similar objects.
- *Incorrect Images*: Randomly assigned from within the subject’s image pool.

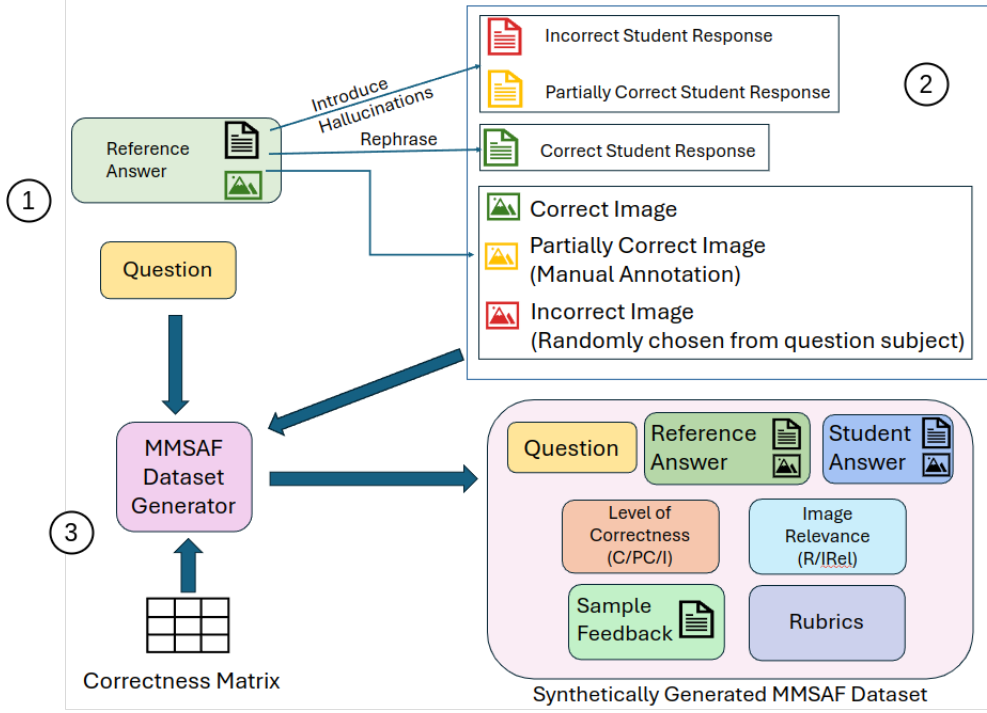
#### 4.2 Generation of Level of Correctness, Image Relevance and Rubrics

Now, given the correct/partially correct/incorrect textual answer and images, we generate the final set of student answers using the correctness matrix (Table 4 in Appendix D). This resulted in the construction of the final dataset, which is of 2,197 data points. Note that the level of correctness for each student’s answer is defined by the output label in Table 4. Similarly, any student answer using an image that is correct or partially correct is assigned an Image Relevance level of “Relevant”, and an incorrect image as “Irrelevant”.

However, to generate the sample feedback and rubrics for error detection in feedback, we use a simple templating strategy. The errors introduced were recorded while constructing the incorrect and partially correct responses and used to construct the sample feedback and rubrics. Hence, the rubrics are beneficial in detecting whether the feedback has successfully detected the errors present in the student’s answer and, if possible, ways to mitigate the error. This ensures that the feedback is beneficial for the student and is helpful in correcting the

<sup>1</sup>The dataset will be released upon publication of this work.





C - Correct PC - Partially Correct I - Incorrect R - Relevant IRel - Irrelevant

Figure 2: An automatic framework to generate the MMSAF dataset

student’s mistakes. The dataset split statistics are in Appendix K.

## 5 LLMs in Consideration

The MMSAF problem involves evaluating multiple images and text as a whole, requiring MLLMs capable of complex reasoning. We select four such models: ChatGPT, Gemini, Pixtral, and Molmo. Open-source models (Molmo and Pixtral) were accessed via the Huggingface library<sup>2</sup>, while APIs were used for ChatGPT<sup>3</sup> and Gemini<sup>4</sup>.

**ChatGPT:** ChatGPT (GPT-4o mini) from OpenAI is a well-established multilingual and multimodal model. It has demonstrated strong performance in grading and educational tasks, including diagram-based scoring (Lee and Zhai, 2023), making it a suitable candidate for MMSAF.

**Gemini:** Gemini, developed by Google, is at par with ChatGPT on multimodal reasoning benchmarks like MMMU. The latest Gemini Ultra has reportedly outperformed GPT-4V<sup>5</sup>. We use the freely available *gemini-1.5-flash* for our experiments.

**Pixtral:** Pixtral 12B (Agrawal et al., 2024), by Mistral.ai, is an open-source multimodal model that outperforms models like LLaVA, Claude-3, and Gemini-1.5 Flash on benchmarks like MMMU<sup>6</sup>. Built on Mistral Nemo 12B and a custom vision encoder (Pixtral-ViT), it combines efficiency with strong performance, making it a lightweight but robust MLLM. Additionally, it is trained to interpret diagrams while providing detailed and structured interpretations.

**Molmo:** Compared to other MLLMs, Molmo (Deitke et al., 2024) by AllenAI has been specifically trained on academic datasets, which makes it a possible candidate for solving the MMSAF problem. As per their website<sup>7</sup>, Molmo has beaten existing leading models such as GPT-4o, Gemini 1.5 and Claude-3 on 11 different academic benchmarks and human evaluation. In particular, we choose the *Molmo-7B-D-0924* model variant, suitable for our dataset’s academic and multimodal nature. The next section details how these models were evaluated on our dataset.

<sup>2</sup><https://huggingface.co/>

<sup>3</sup><https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

<sup>4</sup><https://ai.google.dev/gemini-api/docs/api-key>

<sup>5</sup><https://blog.google/technology/ai/google-gemini-ai/#performance> (Last Accessed: 05-12-2024)

<sup>6</sup><https://mistral.ai/news/pixtral-12b/> (Last Accessed: 05-12-2024)

<sup>7</sup><https://molmo.allenai.org/blog> (Last Accessed: 05-12-2024)

## 6 Evaluation of LLM Generated Feedback

### Prompt 1: Feedback Prompt

...

Task : You have to generate the level of correctness, the image relevance and the feedback.

The feedback should point out any errors in the text as well as the image.

It should also provide the reason for the level of correctness and image relevance.

It can contain some additional information and facts to complement the student’s understanding as well.

It should be a conversation between you as a teacher and a student.

It should be of 500 words.

...

The goal of this evaluation was to quantify the zero-shot performance of existing LLMs on this dataset and grade their capabilities. To do so, 221 data points (130 from biology, 56 from chemistry and 35 for physics) were sampled randomly and fed to different LLMs using Prompt 1. The LLMs used were ChatGPT, Gemini, Pixtral and Molmo and their respective details can be found in Section 5.

Note that the complete prompt can be found in Appendix E and the experimental setup in Appendix L. The corresponding LC, IR and feedback values generated were collected and then analysed in the following subsections.

### 6.1 Analysis of Correctness and Relevance levels

To evaluate the performance of the models on predicting Level of Correctness and Image Relevance labels, we report macro-averaged accuracy, precision, recall, and F1-score values.

Model	Accuracy	Precision	Recall	F1
ChatGPT	0.50	0.32	0.31	0.30
Molmo	0.42	0.36	0.34	0.21
Pixtral	0.52	0.32	0.32	0.32
<b>Gemini</b>	<b>0.55</b>	<b>0.44</b>	<b>0.68</b>	<b>0.48</b>

Table 1: Metrics for generated Level of Correctness labels

As shown in Table 1, Gemini outperforms all other models across the evaluated metrics, suggesting higher reliability in predicting LC labels and

fewer false positives. ChatGPT had a tendency to label most answers as “Partially Correct” class, impacting its overall performance. Molmo, in contrast, exhibited a strong bias towards labelling answers as “Incorrect”. Pixtral, while more lenient, frequently confused “Incorrect” and “Partially Correct” responses, which led to reduced precision. However, this behaviour indicates potential for performance improvement through fine-tuning.

Model	Accuracy	Precision	Recall	F1
<b>ChatGPT</b>	<b>0.75</b>	<b>0.76</b>	<b>0.81</b>	<b>0.74</b>
Molmo	0.29	0.15	0.49	0.23
Pixtral	0.66	0.59	0.59	0.59
Gemini	0.58	0.71	0.70	0.58

Table 2: Metrics for generated Image Relevance levels

Table 2 presents the results for IR prediction. ChatGPT achieves the highest performance across all metrics, indicating strong multimodal reasoning in assessing image relevance. Molmo frequently predicted most images as relevant, which led to poor precision. Pixtral also suffered from false positives, while Gemini often misclassified relevant images as irrelevant.

Further breakdowns and confusion matrices can be found in Appendix G and Appendix H.

### 6.2 Evaluation Task for Experts

The 221 data points mentioned earlier were provided to the LLMs along with a prompt (as in Appendix E), and their feedback, level of correctness and image relevance values were recorded and presented to six Subject Matter Experts (SMEs), where we have 3 experts for each domain namely, physics, chemistry and biology. Relevant details about SMEs are mentioned in Appendix F.

The SMEs were instructed to score each feedback on a scale of 1 to 5 based on the following parameters -

1. *Fluency and Grammatical Correctness (FGE)*: This metric denotes the fluency and grammatical correctness of the LLM-generated feedback. The idea is to check if the LLM-generated sentences are grammatically correct or not. A score of 1 denotes that the FGE level of the feedback is extremely poor while a score of 5 indicates that it is excellent.
2. *Emotional Impact (EI)*: This metric is to check whether the LLM-generated feedback will have a positive impact on the student or not,

that is, whether the feedback is more encouraging and assistive for the student or not. A score of 1 denotes negative impact, while a score of 5 denotes a positive impact.

3. *Level of Feedback Correctness (LC)*: This metric is to determine whether the feedback has properly captured all the errors present in the student’s answer. A score of 1 denotes that no error has been captured in the feedback, while a score of 5 denotes that all the errors have been captured in the feedback.
4. *Error Mitigation in Feedback (EM)*: This metric evaluates whether the feedback has properly addressed each and every error present in the student’s answer and suggested ways to correct them. A score of 1 denotes no such error mitigation has been done, while a score of 5 denotes that all the ways necessary to correct all errors are present.
5. *Rubrics for error detection (FR)*: While traditional NLP metrics like ROUGE-2 (Ganesan, 2018), SCAReBLEU (Post, 2018), METEOR (Banerjee and Lavie, 2005), and BERTScore (Zhang et al., 2020) rely on n-gram overlap, they often miss the semantic accuracy of feedback, particularly in identifying and addressing student errors. To address this, we propose a rubric-based evaluation that checks whether the feedback captures all relevant errors. Annotators assess each rubric as a True/False question based on the LLM-generated feedback.

### 6.3 Analysis of Expert Evaluation

Once the *three* SMEs had completed their respective evaluation tasks, all the scores for each metric were collected and averaged out to present the final data. Any disagreement was resolved using majority voting among three raters. Table 3 summarises the average ratings assigned by annotators to the feedback generated by different LLMs over various criteria mentioned in Section 6.2 for each and every subject.

*Physics*: Physics questions involve interpreting abstract diagrams containing objects such as arrows, circles and rectangles. It also involves interpreting certain domain-specific objects such as lenses and mirrors. Apart from diagrams, such questions involve calculations and physics concepts based on reasoning. ChatGPT outperformed others in all areas, with SMEs highlighting its strength in

identifying calculation errors and providing step-by-step explanations. This is reflective of its strong reasoning and math skills, also validated by benchmarks like MATH-500 and MMMU<sup>8</sup>. Molmo often hallucinated, mislabeled correct answers, and struggled with reasoning, though it handled concrete diagrams reasonably well. Pixtral delivered structured feedback with clearly structured explanations, aligning with its training data where it had to interpret graphs and diagrams and provide structured analysis.

*Chemistry*: Chemistry questions involve domain-specific chemical formulas as part of their questions and diagrams. They also include interpreting graphs. Some diagrams involve interpreting domain-specific objects such as a beaker, a scientific fork and others. When compared to other models, ChatGPT performs the best on all parameters. SMEs point out that ChatGPT provides detailed explanations. The reason is the same as for physics questions, as these questions test the reasoning skills of such MLLMs and how well they can interpret abstract diagrams. However, they also point out that Molmo’s tone was too direct for students and was not good at providing proper explanations for its assigned labels for the level of correctness and image relevance. This is because Molmo is not trained for chemistry questions, nor for interpreting such domain-specific concepts from chemistry.

*Biology*: Biology questions are more factual, involving real-life diagrams such as a heart. MLLMs trained on such data will perform better in identifying key concepts. Pixtral excelled, providing structured, concept-based feedback and maintaining a polite tone, leading to higher emotional impact scores. Though Molmo showed decent visual understanding, it failed to link errors to relevant concepts, highlighting a gap in reasoning. ChatGPT and Gemini performed well, but not as effectively as Pixtral.

To summarise, ChatGPT is best suited for physics and chemistry due to its reasoning and mathematical strengths. Pixtral is more effective in biology, owing to its detailed, empathetic feedback and structured analysis. SMEs noted Molmo’s feedback often included Chinese characters, was overly direct, and sometimes mislabeled correct answers, suggesting a need for fine-tuning. More details are included in Appendices I and J and ablation

<sup>8</sup><https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

Subject	Model	Evaluation Parameters				
		FGE	EI	LC	EM	Rubrics
Physics	ChatGPT	<b>5.00</b>	<b>4.59</b>	<b>4.53</b>	<b>4.56</b>	<b>0.78</b>
	Gemini	5.00	4.56	4.12	4.09	0.68
	Molmo	4.75	2.66	2.69	2.25	0.53
	Pixtral	5.00	4.15	4.24	4.24	0.55
Chemistry	ChatGPT	<b>4.98</b>	<b>4.98</b>	<b>4.42</b>	<b>4.51</b>	0.63
	Gemini	4.95	4.67	4.37	4.42	<b>0.65</b>
	Molmo	4.70	3.25	2.93	2.98	0.48
	Pixtral	4.95	4.71	4.10	4.05	0.49
Biology	ChatGPT	5.00	3.07	3.24	3.15	0.50
	Gemini	5.00	3.59	<b>3.43</b>	3.07	0.53
	Molmo	4.92	3.11	3.04	2.69	0.54
	Pixtral	<b>5.00</b>	<b>3.87</b>	3.40	<b>3.48</b>	<b>0.58</b>

Table 3: Average expert evaluation scores of different metrics on each subject

studies in Appendix N.

## 7 Conclusion

This paper introduces the MMSAF problem along with a dataset of 2,197 data points. The MMSAF dataset contains physics, chemistry and biology questions from high school textbooks. Additionally, we provide an automated framework to generate similar datasets, given a set of questions, reference answers, and annotated images. We also establish a baseline using 4 models, namely ChatGPT, Gemini, Pixtral and Molmo, across all three subjects.

Our evaluations show that while Gemini performed the best in generating the correctness labels and ChatGPT excelled in generating the image relevance labels, human evaluations proved that Pixtral was more aligned towards human judgement and values for biology and ChatGPT for physics and chemistry. Future work also involves exploring other solutions, including Retrieval augmented generation (RAG) based approaches to add more insight and conceptual depth to the feedback.

## Limitations

This work introduces the Multimodal Short Answer grading with Feedback (MSMAF) problem along with a dataset of 2,197 data points. Note that all the student answers are synthetically generated. They are only representative of a subset of mistakes made by students in real-life examinations. However, collecting real-life data from examinations

involves legal considerations and barriers. Thus, such synthetic datasets often serve as an alternative to real-life data to aid progress in such scenarios.

While the dataset is currently restricted to physics, chemistry, and biology questions at the high school level, adding data points from other subjects and university-level courses can increase its complexity and richness. Another limitation of the automated pipeline pertaining to partially correct images is that it needs to be manually annotated, which can lead to a scalability issue. This problem can be automated using simple text manipulation operations performed via OpenCV<sup>9</sup> or by generating diagrams using solutions such as DiagrammerGPT (Zala et al., 2024) once their code is released.

Since such systems are still in their nascent stages, it will be a better practice to use them in case of practice exercises or examinations having low weightage rather than examinations, which have a higher weightage. This will partly reduce the grading load from teachers and also provide feedback for the students’ improvement.

## Ethical Considerations

The questions and reference answers have been extracted from the National Council of Educational Research and Training (NCERT) textbooks for 10<sup>th</sup> standard science, 11<sup>th</sup> standard, and 12<sup>th</sup> standard physics, chemistry and biology. NCERT is an autonomous organisation established by the Govern-

<sup>9</sup><https://opencv.org/>



ment of India that provides guidance and recommendations to both central and state governments on policies and initiatives to improve the quality of school education. The textbooks have been downloaded from <https://www.ncrtsolutions.in/>. We adhere to NCERT guidelines, which state that "NCERT books can also be downloaded free of cost from our website for non-commercial purposes."

## References

Dishank Aggarwal, Pushpak Bhattacharyya, and Bhaskaran Raman. 2024. "i understand why i got this grade": Automatic short answer grading with feedback. *Preprint*, arXiv:2407.12818.

Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, Teven Le Scao, Andy Lo, William Marshall, Louis Martin, Arthur Mensch, Pavankumar Muddireddy, Valera Nemychnikova, Marie Pellat, Patrick Von Platen, Nikhil Raghuraman, Baptiste Rozière, Alexandre Sablayrolles, Lucile Saulnier, Romain Sauvestre, Wendy Shang, Roman Soletskyi, Lawrence Stewart, Pierre Stock, Joachim Studnia, Sandeep Subramanian, Sagar Vaze, Thomas Wang, and Sophia Yang. 2024. *Pixtral 12b*. *Preprint*, arXiv:2410.07073.

Satanjeev Banerjee and Alon Lavie. 2005. *METEOR: An automatic metric for MT evaluation with improved correlation with human judgments*. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Steven Burrows, Iryna Gurevych, and Benno Stein. 2015. *The eras and trends of automatic short answer grading*. *International Journal of Artificial Intelligence in Education*, 25:60–117.

Galina Deeva, Daria Bogdanova, Estefanía Serral, Monique Snoeck, and Jochen De Weerd. 2021. *A review of automated feedback systems for learners: Classification framework, challenges and opportunities*. *Computers & Education*, 162:104094.

Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz,

Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wiltliff, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. 2024. *Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models*. *Preprint*, arXiv:2409.17146.

Menna Fateen, Bo Wang, and Tsunenori Mine. 2024. *Beyond scores: A modular rag-based system for automatic short answer scoring with feedback*. *Preprint*, arXiv:2409.20042.

Anna Filighera, Siddharth Parihar, Tim Steuer, Tobias Meuser, and Sebastian Ochs. 2022. *Your answer is incorrect... would you like to know why? introducing a bilingual short answer feedback dataset*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8577–8591, Dublin, Ireland. Association for Computational Linguistics.

Kavita Ganesan. 2018. *Rouge 2.0: Updated and improved measures for evaluation of summarization tasks*. *Preprint*, arXiv:1803.01937.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. *A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions*. *ACM Trans. Inf. Syst.*, 43(2).

Gyeong-Geon Lee and Xiaoming Zhai. 2023. *Nerif: Gpt-4v for automatic scoring of drawn models*. *Preprint*, arXiv:2311.12990.

Chee Wee Leong, Lei Liu, Rutuja Ubale, and Lei Chen. 2018. *Toward large-scale automated scoring of scientific visual models*. In *Proceedings of the Fifth Annual ACM Conference on Learning at Scale, L@S '18*, New York, NY, USA. Association for Computing Machinery.

Jiazheng Li, Lin Gui, Yuxiang Zhou, David West, Cesare Aloisi, and Yulan He. 2023. *Distilling ChatGPT for explainable automated student answer assessment*. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6007–6026, Singapore. Association for Computational Linguistics.

Chin-Yew Lin. 2004. *Rouge: A package for automatic evaluation of summaries*. In *Annual Meeting of the Association for Computational Linguistics*.

Elizabeth J. Marsh and Emmaline Drew Eliseev. 2019. *Correcting Student Errors and Misconceptions*, page 437–459. Cambridge Handbooks in Psychology. Cambridge University Press.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. *Bleu: a method for automatic evaluation of machine translation*. In *Proceedings of the*

40th Annual Meeting on Association for Computational Linguistics, ACL '02, page 311–318, USA. Association for Computational Linguistics.

Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ari Sagherian, Suhasini Lingaiah, Mohamed Abouelenien, Chee Wee Leong, Lei Liu, Mengxuan Zhao, Blake Lafuente, Shu-Kang Chen, and Yi Qi. 2022. [Learning progression-based automated scoring of visual models](#). In *Proceedings of the 15th International Conference on PErvasive Technologies Related to Assistive Environments*, PETRA '22, page 213–222, New York, NY, USA. Association for Computing Machinery.

Ashita Saxena. 2024. Hallucination detection in machine generated text. Master’s thesis, Indian Institute of Technology Bombay.

Gemini Team. 2024. [Gemini: A family of highly capable multimodal models](#). *Preprint*, arXiv:2312.11805.

Maximilian Tornqvist, Mosleh Mahamud, Erick Mendez Guzman, and Alexandra Farazouli. 2023. [ExASAG: Explainable framework for automatic short answer grading](#). In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 361–371, Toronto, Canada. Association for Computational Linguistics.

Mengxia Yu, Zhihan Zhang, Wenhao Yu, and Meng Jiang. 2023. [Pre-training language models for comparative reasoning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12421–12433, Singapore. Association for Computational Linguistics.

Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2024. [Diagrammergpt: Generating open-domain, open-platform diagrams via llm planning](#). *Preprint*, arXiv:2310.12128.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

## A Prompt used for synthetically generating correct responses

As mentioned in Section 4.1, the prompt used to generate the correct responses is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER> tags.

### Prompt 2: Prompt for generating correct responses

Prompt: You are a student who is attempting an examination.

Task: Given a question and its original answer, rewrite the original answer as a paragraph in a different tone so that it is correct and captures all the necessary facts in the original answer.

<QUESTION></QUESTION>

<ANSWER></ANSWER>

Output Requirements:

<HANSWER>The correct answer</HANSWER>

## B Prompt used for synthetically generating partially correct responses

As mentioned in Section 4.1, the prompt used to generate the partially correct responses by introducing factual fabrication is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER> tags. It is similar to the one used by Ashita et al in (Saxena, 2024).

### Prompt 3: Prompt for generating partially correct responses (Factual Fabrication)

Factual fabrications in text can be created by introducing contextually relevant information or facts that are not verifiable by established real-world evidence. This technique involves adding elements that fit logically within the context but are factually incorrect or unverifiable.

Task: Given a question and its corresponding original factual answer, you are to create a factual fabrication in the answer by introducing a new, contextually relevant information or fact that cannot be verified against real-world knowledge. You should rewrite the answer (the hallucinated answer) by adding the fabricated information in the original answer as well as the information/-fact you introduced in the original answer to generate the hallucinated answer.

<QUESTION></QUESTION>

<ANSWER></ANSWER>

Output Requirements:

<HANSWER>The hallucinated answer</HANSWER>

<RDETAILS>The additional information/-fact introduced to generate the hallucinated answer in a single sentence .<RDETAILS>

The prompt used to generate the partially correct responses by introducing factual inconsistency is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER> tags. It is similar to the one used by [Saxena \(2024\)](#).

Prompt 4: Prompt for generating partially correct responses (Factual Inconsistency)

Prompt: In the field of dependency parsing, modifiers are defined as words or phrases that provide additional information about other elements in a sentence. One technique to generate deliberate factual inconsistencies in text, termed the "Termite Strategy" targets these modifiers. This strategy involves replacing modifiers with alternative words or phrases that are factually inconsistent yet still maintain the overall coherence of the sentence.

Task:

Given a question and its original answer, apply the Termite Strategy to introduce a factual inconsistency in the original answer. Replace a modifier in the original answer with an alternative that contradicts the factual information in the answer, but still retains sentence coherence. You must rewrite the "complete" answer with the modifications (the "hallucinated answer") and also provide the replacement details.

<QUESTION></QUESTION>

<ANSWER></ANSWER>

Output Requirements:

<HANSWER>The hallucinated answer</HANSWER>

<RDETAILS>Describe the original modifier and the replacement word or phrase used to create the inconsistency in a single sentence<RDETAILS>

## C Prompt used for synthetically generating incorrect responses

As mentioned in Section 4.1, the prompt used to generate the incorrect responses by introducing factual fabrication is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER> tags. It is similar to the one used by [Saxena \(2024\)](#).

Prompt 5: Prompt for generating incorrect responses (Factual Fabrication)

Prompt: Factual Fabrication refers to instances where the LLM's output contains facts that are unverifiable against established real-world knowledge.

Task:

Given a question, the task is to generate an incorrect answer using the techniques of factual fabrication.

<QUESTION></QUESTION>

You should follow the following output format:

Output Requirements:

<HANSWER>The hallucinated answer</HANSWER>

<RDETAILS>Mention why the answer is incorrect in a single sentence<RDETAILS>

The prompt used to generate the incorrect responses by introducing factual inconsistency is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER> tags. It is similar to the one used by [Saxena \(2024\)](#).

#### Prompt 6: Prompt for generating incorrect responses (Factual Inconsistency)

Prompt: Factual Inconsistency refers to situations where the LLM's output contains facts that can be grounded in real-world information, but present contradictions.

Task:

Given a question , the task is to generate an incorrect answer using the techniques of factual inconsistency .

<QUESTION></QUESTION>

You should follow the following output format:

Output Requirements:

<HANSWER>The hallucinated answer</HANSWER>

<RDETAILS>Mention why the answer is incorrect in a single sentence</RDETAILS>

tags and the student answer within the <STUDENT></STUDENT> tags. Additionally, we place the reference answer image and the student answer image one after the other, just after the prompt and pass it to the corresponding LLM API.

#### Prompt 7: Feedback Prompt

Act as a teacher and grade the student answer given the question, reference answer and the student answer.

Task : You have to generate the level of correctness, the image relevance and the feedback. The feedback should point out any errors in the text as well as the image. It should also provide the reason for the level of correctness and image relevance. It can contain some additional information and facts to complement the student's understanding as well. It should be a conversation between you as a teacher and a student. It should be of 500 words.

Input format -

<QUESTION>The question</QUESTION>

<ANSWER>The reference answer. Note that the first image corresponds to the image in the reference answer</ANSWER>

<STUDENT>The student answer. Note that the second image corresponds to the image in the reference answer</STUDENT>

You have to strictly follow this output format-

<CORRECTNESS>Predict whether the answer is Correct, Partially Correct or Incorrect. Note you should evaluate both the text and image of the student answer as a whole. </CORRECTNESS>

<RELEVANCE>Predict whether the second image is relevant or irrelevant to the question</RELEVANCE>

<REASON>The feedback</REASON>

Here is the input :

<QUESTION></QUESTION>

## D Correctness Matrix

The correctness mentioned in Section 4.2 has been explained here. Given the correct/partially correct/incorrect textual answer and images, the final set of student answers is generated using the correctness matrix as below.

Student Response Text	Student Response Image	Overall Correctness Label
C	C	C
C	PC/I	PC
PC	C/PC	PC
PC	I	I
I	C	PC
I	PC/I	I

C - Correct PC - Partially Correct I - Incorrect

Table 4: Matrix for determining the Level of Correctness

## E Prompt to Generate Feedback

As mentioned in Section 6.2, the prompt used to generate the feedback given the question, reference answer, and the student answer is as follows. Note we place the question within the <QUESTION> and the </QUESTION> tags. Similarly, we place the reference answer within the <ANSWER> </ANSWER>



<ANSWER></ANSWER>  
<STUDENT></STUDENT>

## F Subject Matter Expert Details

The expert evaluation of LLMs involved the help of six Subject Matter Experts (SMEs) (as per Section 6.2), where two SMEs were present each for biology, physics, and chemistry.

*Biology:* SME 1 has been a high school educator with teaching experience for over 25 years. SME 2 is a research scholar with more than 2 years of experience. SME 3 is a researcher with more than 5 years of teaching experience.

*Physics:* SME 1 is currently a professor at a reputed institute with over 13 years of teaching experience. SME 2 is a research scholar with more than 2 years of experience. SME 3 is a research scholar with over 2 years of teaching experience.

*Chemistry:* Both our SMEs are research scholars with more than 3 and 4 years of teaching experience, respectively. SME 3 is an educator with over 20 years of teaching experience.

The LLM-generated feedback had an average of 1758 words and a maximum of 3538 words. Each reviewer had to go through 880 instances of LLM-generated feedback and rate them based on the evaluation parameters as mentioned in Section 6.2. As pointed out by SMEs, this was a very tiring and time-consuming task as they had to meticulously go through a single feedback multiple times for each parameter. Hence, each of them has been compensated appropriately as per the norms for their evaluations and comments.

## G Additional Information for Level of Correctness Labels

This appendix analyses the confusion matrix obtained for all the LLMs over the generated level of correctness labels. The major observations have already been mentioned in Section 6.1.

Figure 3 displays the confusion matrix obtained for Gemini over the Level of Correctness labels. While it has correctly predicted the "Correct" labels, it has often labelled "Partially Correct" answers as "Incorrect" and vice versa. However, the generated trend indicates that the model was being lenient in labelling most "Incorrect" answers as "Partially Correct". However, for some answers, it failed to

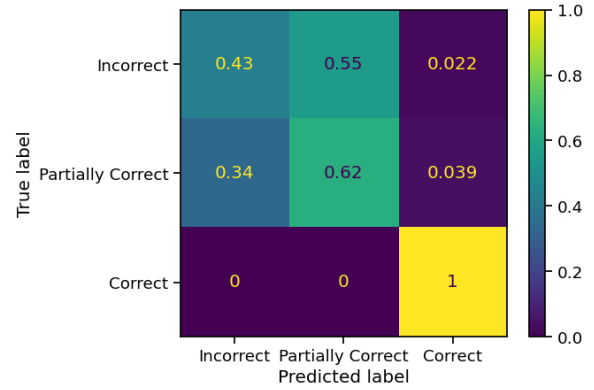


Figure 3: Confusion Matrix for Gemini after True Class Normalization

evaluate the answer properly and marked "Partially Correct" answers as "Incorrect".

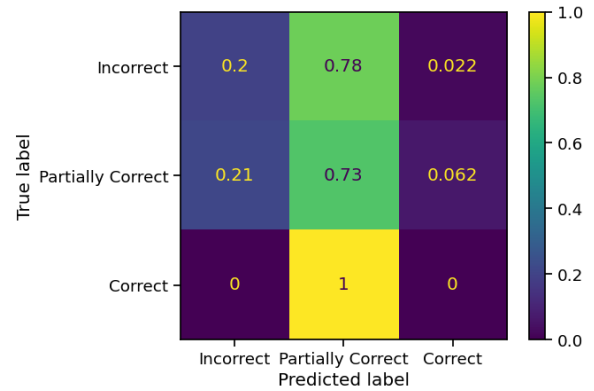


Figure 4: Confusion Matrix for ChatGPT after True Class Normalization

Figure 4 displays the confusion matrix obtained for ChatGPT over the Level of Correctness labels. It can be seen that ChatGPT was lenient in grading the student's answer as it labelled most "Incorrect" answers as "Partially Correct" while it made a mistake of labelling "Correct" answers as "Partially Correct" as well.

Figure 5 displays the confusion matrix obtained for Pixtral over the Level of Correctness labels. Similar to ChatGPT, Pixtral was lenient and labelled "Incorrect" answers as "Partially Correct" which also led to the reduced metrics.

Figure 5 displays the confusion matrix obtained for Molmo over the Level of Correctness labels. As it can be seen, it labelled almost all answers as "Incorrect" indicating that it was unable to show considerable performance while performing complex comparative reasoning over multiple modalities.

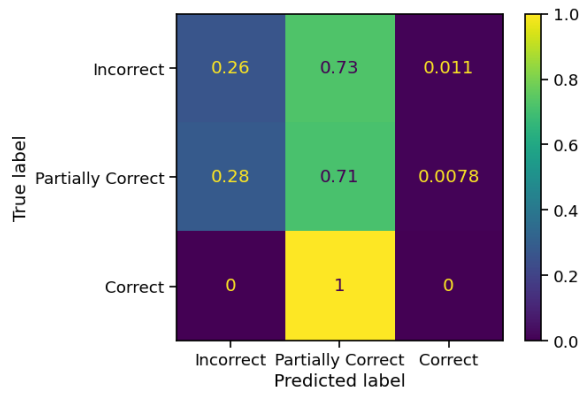


Figure 5: Confusion Matrix for Pixtral after True Class Normalization

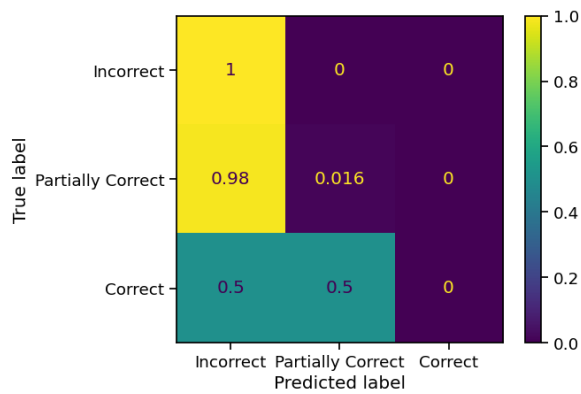


Figure 6: Confusion Matrix for Molmo after True Class Normalization

## H Additional Information for Image Relevance Labels

This appendix analyses the confusion matrix obtained for all the LLMs over the generated image relevance labels. The major observations have already been mentioned in Section 6.1.

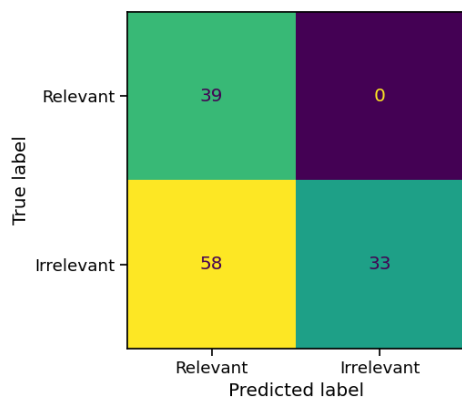


Figure 7: Confusion Matrix for Gemini

Figure 7 displays the confusion matrix obtained for Gemini over the Image Relevance labels. Gemini often predicted "Irrelevant" images as "Relevant", which led to reduced metrics.

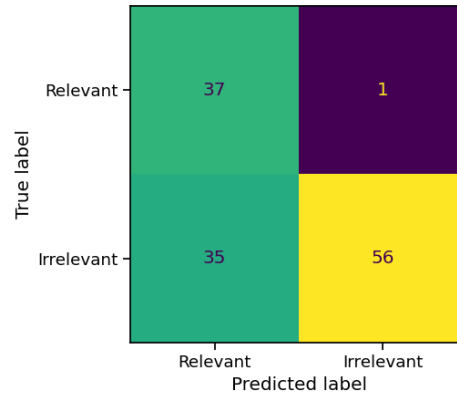


Figure 8: Confusion Matrix for ChatGPT

Figure 8 displays the confusion matrix obtained for ChatGPT over the Image Relevance labels. ChatGPT had the highest accuracy among all the LLMs. However, it still failed to classify a small part of "Relevant" images and "Irrelevant". This is also consistent with the fact that ChatGPT has superior multimodal reasoning capabilities compared to Gemini, Pixtral and Molmo, and this fact has also been verified by different metrics on standard benchmarks such as MMMU.

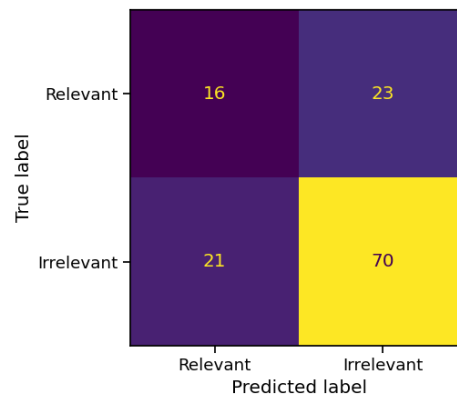


Figure 9: Confusion Matrix for Pixtral

Figure 9 displays the confusion matrix obtained for Pixtral over the Image Relevance labels. When it comes to Pixtral, it has often incorrectly classified "Relevant" images as "Irrelevant" and vice versa.

Figure 10 displays the confusion matrix obtained for Molmo over the Image Relevance labels. Molmo has classified all images as "Relevant", bar-

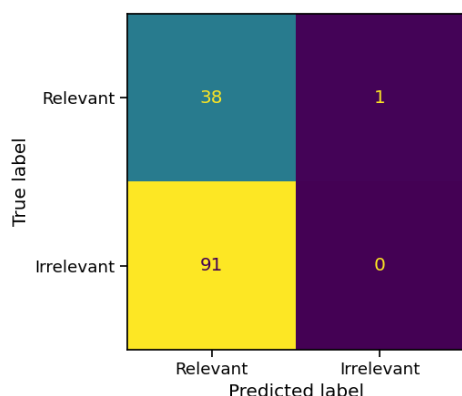


Figure 10: Confusion Matrix for Molmo

ring one. This indicates that there is a scope to improve Molmo’s multimodal reasoning capabilities.

## I Analysis of Rubrics for error detection

This is an addition to Section 6.3. The synthetically generated student answers contain both text and image parts, which might contain errors. Thus, the generated rubrics can either be used to evaluate errors present in the textual part or in the image part. We first analysed how well the LLMs performed in detecting the textual errors, followed by their performance in the case of image errors. Note that this data was again annotated by the SMEs as stated in section 6.2. The error detection accuracy ( $Ac_{ED}$ ) for textual or image errors is computed as

$$Ac_{ED} = \frac{\text{Number of rubrics detected}}{\text{Total number of rubrics}}$$

Note that number of rubrics detected is equivalent to number of True values obtained for rubrics divided by total number of rubrics. Also, we use majority voting to resolve any disagreements.

Figure 11 contains the  $Ac_{ED}$  pertaining to the text errors of each LLM on each subject and on all subjects as a whole, while Figure 12 contains the same but for image errors of each LLM on each subject and on all subjects as a whole. From both Figure 11 and Figure 12, it can be concluded that ChatGPT performed the best across all subjects with scores of 0.73 and 0.9, respectively. However, if we analyse across individual subjects, ChatGPT was able to beat all except for Gemini in biology for text error-based rubrics and was comparable to other models in biology for image-based rubrics.

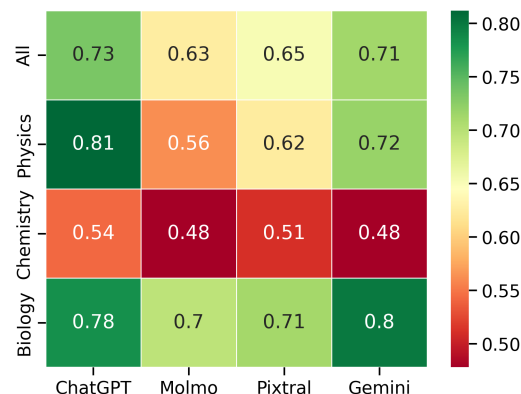


Figure 11:  $Ac_{ED}$  values for Text based Rubrics

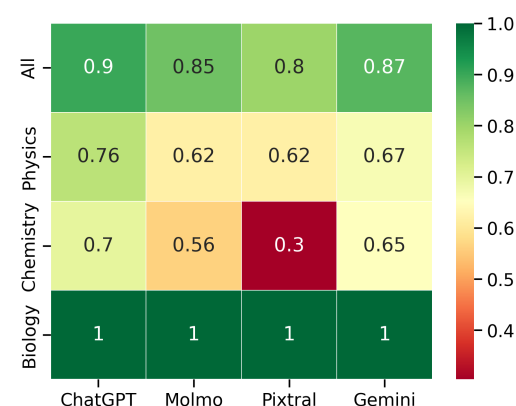


Figure 12:  $Ac_{ED}$  values for Image based Rubrics

## J Extended Results

This is an addition to Section 6.3. In this appendix, we calculate the difference in performance that arises due to one prompt working better for a given model than others. For the level of correctness, Gemini displayed a 6% improvement over ChatGPT, a 23.64% improvement over Molmo and a 5.43% improvement over Pixtral. For Image relevance labels, ChatGPT displayed a 20.74% improvement over Gemini, 56.14% over Molmo, and 12.2% over Pixtral.

Earlier in Section 6.3, we had shown that ChatGPT performs the best in physics and chemistry, while Pixtral performs the best in biology. Table 5 denotes the percentage improvement of the best model over the other models in each subject.

*Physics:* It can be seen that ChatGPT significantly beats other models in all parameters except for Gemini in FGE, EI, and Pixtral in FGE. Since LC and EM play a major role in the generated feedback, we can conclude that ChatGPT is the optimal solution for generating feedback for physics.

*Chemistry:* When it comes to chemistry, ChatGPT again significantly beats other models in all parameters, which is indicative of the fact that it is the best solution for generating feedback in Chemistry.

*Biology:* Here, Pixtral significantly beats other models in all parameters except for Gemini and ChatGPT in FGE. However, since LC and EM play a major role in the generated feedback, we can still conclude that Pixtral is the better solution for generating feedback in the case of biology.

## K Dataset Split Statistics

As mentioned in Section 4.2, the data has been split into the train, test and validation sets in the ratio of 3:1:1 and the statistics is as follows -

The mean and standard deviation of word lengths are as follows -

## L Experimental Setup

As mentioned in Section 6, the system specifications used for LLM Inference are as follows -

- CPU: Intel(R) Xeon(R) Gold 6348 CPU @ 2.60GHz with 112 CPUs
- RAM: 503 GB
- GPU: NVIDIA A100-SXM GPU with 80 GB of memory (Only 1 out of 4 used)

The relevant models used from Huggingface, OpenAI, and Gemini APIs have already been mentioned in Section 5. All inference operations use default hyperparameters.

## M Common errors made by students during exams

This is an elaboration on Section 4.1. To facilitate better development of such synthetic datasets, we list down the common mistakes made by students during examinations. These examples have also been suggested by SMEs and categorised as per the categories listed by Marsh and Eliseev (2019). Thus broadly, such errors are categorised as follows -

- *Errors made with confidence:* These errors are commonly made by students where the student confidently provides an answer that is derived from a correct fact, but is totally incorrect given the context or a completely made-up fact. This again contains scenarios such as -
  - *Incorrect Translations:* These errors are made by students when they generate a wrong summarization of a fact. One example from Marsh et al’s work includes that while the correct fact is “George Washington has dentures made of bone and other non-wood materials.”, the student mistakenly mentions “George Washington has wooden teeth”. This is a wrong summarisation of the earlier-mentioned fact.
  - *Factual Errors:* These errors are made by students when they use a correct fact but in the wrong context. An example provided by our SME is as follows: given the fill-in-the-blank type question “Haemoglobin consists of . . . and globin.” the student incorrectly mentions calcium instead of haem. While calcium is a known fact and is present in the human body, it is not present in haemoglobin.
  - *High confidence misconceptions or Erroneous guess:* These errors are made by students when they try to fill in the gaps in their knowledge with generic concepts. For example, a student mentioned the Sahara consists of sand, while, in reality, it is a rocky desert. Another example mentioned by our SME is that when a student



Subject	Model	Evaluation Parameters (% Improvement)				
		FGE	EI	LC	EM	Rubrics
Physics	Gemini	0.00	0.66	9.95	11.49	14.71
	Molmo	5.26	72.56	68.40	102.67	47.17
	Pixtral	0.00	10.60	6.84	7.55	41.82
Chemistry	Gemini	0.61	6.64	1.14	2.04	-3.08
	Molmo	5.96	53.23	50.85	51.34	31.25
	Pixtral	0.61	5.73	7.80	11.36	28.57
Biology	ChatGPT	0.00	26.06	4.94	10.48	16.00
	Gemini	0.00	7.80	-0.87	13.36	9.43
	Molmo	1.63	24.44	11.84	29.37	7.41

Table 5: Percentage improvement of best model over other models in each subject

	Train	Validation	Test
Correct	102	34	35
Partially Correct	692	232	232
Incorrect	522	174	174
<b>Total</b>	<b>1316</b>	<b>440</b>	<b>441</b>

Table 6: Dataset split statistics

	Mean	Standard Deviation
Question	118.61	84.20
Reference Answer	924.30	688.72
Student Answer	888.13	569.12

Table 7: Word statistics of dataset

is asked to label the parts of a leaf, they mistakenly label the petiole as the stem.

- *Calculations errors:* Such errors are made by students in the case of numerical questions where they can guess the correct answer but are unable to produce it due to failure to recall certain formulas or another calculation mistake. These errors are common in questions where the students have to prove a theorem. Since the theorem is a well-known fact, students often make up the reasoning process to account for it.

- *Misunderstanding:* These errors are commonly made by students when they fail to understand a question or concept and answer

the question incorrectly. Some such scenarios include cases -

- *Misunderstanding a concept:* Such errors are seen in students’ answers when they have a wrong understanding of a certain concept. One such example is when students are asked how echolocation works, they mention that bats shout at objects to judge their distance, and other mammals can do the same as well. However, the correct fact is that bats are able to echolocate because of specialised organs that might be absent in other mammals. As indicated by this example and also pointed out by our SME, the answer is wrong due to wrong reasoning, which is evident by the fact that the student’s concept is not clear.
- *Memory problem:* The student makes such mistakes when they are unable to recall facts relevant to a question due to a wrong conceptual clarity. One such example pointed out by our SME is that a student fails to provide the scientific name of Mango as they do not have concept clarity.
- *Misunderstanding the question:* Such errors happen when the student fails to interpret a question properly. These errors are commonly seen in passage-based questions where the student is unable to comprehend the passage properly or the question relevant to it, and provides wrong facts from the passage as their answer.

- *Conceptual Change:* These errors happen when a student understands a concept but fails to use it in a different context or determine how it interacts with other concepts. One such example from Marsh et al.'s work is that a student answers that diffusion happens because the oxygen and carbon dioxide molecules move willingly. In reality, the correct fact is that diffusion occurs due to different concentration levels of these molecules and their existing Brownian motion. Another example pointed out by our SME is when the student is asked how cross-pollination occurs, and the student fails to answer because they are unable to recall how pollination takes place.

## N Ablation Studies

This is an extension to Section 6.3. To study the variance in the feedback generated after running the same prompt 10 times on the same LLM and input, we do an empirical evaluation and show two such outputs below. The common observation is that the core content of the feedback, which detects inaccuracies or suggests ways to correct them in the student's answer, remains consistent. However, the presentation becomes different.

### Feedback 1

... The question asked where to find lactic acid bacteria and some of their uses. Instead, you wrote a detailed description of the human nervous system, including the brain, eye, and ear. The image provided is a cross-section of the cochlea, part of the inner ear and also irrelevant to the topic ... Lactic acid bacteria are found in various fermented foods. The most common example is curd (yogurt). They are also found in other fermented foods like sauerkraut, kimchi, pickles, and some cheeses. Additionally, they naturally exist in our gastrointestinal tracts... Fermentation provides preservation, adds flavor, and creates beneficial nutritional changes...Lactic acid bacteria can improve gut health by promoting the growth of beneficial bacteria in the intestines. They aid in digestion, and some strains produce vitamins like vitamin K and B vitamins....The image you included is a diagram of the cochlea, a part of the inner

ear. This is entirely unrelated to the question about lactic acid bacteria.

### Feedback 2

... The question asked about lactic acid bacteria and where they are found, including their applications. Instead, your response is a detailed description of the human nervous system, including the brain, eye, and ear. ... The image you provided is a diagram of the cochlea, a part of the inner ear. A picture of curd or a diagram illustrating fermentation would have been appropriate.... Lactic acid bacteria are found in various fermented foods, including yogurt (curd), kefir, sauerkraut, kimchi, pickles, and many other products. They are also found in some parts of the human gut ... Lactic acid bacteria are crucial in food preservation (fermentation), increasing the nutritional value of food (e.g., vitamin B12 in yogurt), and improving gut health (probiotics). They can also have industrial applications, such as producing lactic acid for various uses.

In the two feedbacks below, the blue regions comment on the question and response provided. Note that while the content is the same, it only varies in the natural language used. Similar to the regions marked using violet, the feedback corrects information pertaining to the question, serving the error mitigation task. However, as seen earlier, they vary in the natural language used.