

# Automated Quiz Checking LLM

Ahmad Faisal (383744), Sundas Rathore (366636), Murtaza Ahmed Khan (365880)

Instructor: Dr. Wajahat Hussain

## Abstract

*This project explores the development of an automated system for evaluating handwritten quizzes using Gemini 1.5 Pro's API, a state-of-the-art multimodal large language model (LLM). The system was assessed across multiple subjects and varying levels of prompt complexity, benchmarking its performance against several other LLMs. Two core implementations were designed: a user-friendly Streamlit-based frontend for processing quizzes individually and a batch-processing script that automates evaluation for entire directories of quizzes. This comparative analysis highlights the robustness and efficiency of Gemini 1.5 Pro in handwritten text understanding and question evaluation. The project demonstrates the potential of multimodal LLMs in educational assessment, offering an innovative approach to quiz checking with implications for scalable and efficient grading solutions.*

## 1. Introduction

The increasing capabilities of large language models (LLMs) have opened new possibilities for automating educational tasks, including the evaluation of handwritten quizzes. This project explores the performance of Gemini 1.5 Pro, a state-of-the-art multimodal LLM, in interpreting and grading handwritten responses. A key aspect of this work involves benchmarking Gemini 1.5 Pro against other leading LLMs, analyzing their effectiveness across various subjects and prompt complexities to understand their strengths and limitations.

The system is designed with two implementations: a user-friendly Streamlit-based frontend for interactive, single-quiz processing, and a batch-processing script capable of evaluating an entire directory of quizzes. Through detailed comparative analysis, the project aims to highlight the robustness, accuracy, and efficiency of Gemini 1.5 Pro in relation to its peers, providing insights into the suitability of multimodal LLMs for educational assessment tasks. This work not only demonstrates the practical applications of LLMs in automated grading but also contributes to a deeper

understanding of their comparative performance.

## 2. Initial Approaches

Several initial approaches were explored before finalizing the use of Gemini 1.5 Pro for automated handwritten quiz evaluation. These approaches included the following:

### 2.1. TrOCR + GPT-2 (Small) / Llama 3B Pipeline

A pipeline combining TrOCR for handwritten text recognition with GPT-2 or Llama 3B for evaluation was tested. However, finetuning TrOCR proved challenging due to difficulties in achieving satisfactory recognition accuracy. Furthermore, GPT-2 produced incoherent outputs, making the pipeline unsuitable for reliable grading.

### 2.2. BLIP-2

BLIP-2 was considered for its multimodal capabilities, which seemed promising for handling handwritten inputs. However, this approach was found to be computationally expensive, with extremely long inference times that hindered practicality. Additionally, the system frequently experienced random crashes, further reducing its feasibility for large-scale or time-sensitive applications.

## 3. Methodology

### 3.1. Models Used

### 3.2. Dataset

The dataset used for this project was carefully curated to evaluate the capabilities of Gemini 1.5 Pro and other LLMs across a diverse range of handwritten quizzes. It was divided into two main parts based on the complexity and subject matter of the quizzes:

#### 3.2.1 Simple Handwritten Quizzes

This subset included quizzes with:

- Very basic questions.
- Short answers.

(383744), Sundas Rathore (366636), Muhammad Mustafa (365880)	CS-471 2024 Submission #Ahmad Faisal (383744), Sundas Rathore (366636), Muhammad Mustafa (365880). CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.	(383744), Sundas Rathore (366636), Muhammad Mustafa (365880)
108	• Topics such as Basic Mathematics, Science, Algebra,	162
109	Physics, and General Knowledge.	163
110		164
111	These quizzes aimed to test the baseline performance of the	165
112	models on straightforward, low-complexity tasks.	166
113		167
114		168
115	<b>3.2.2 University Coursework Quizzes</b>	169
116	This subset comprised more complex quizzes featuring:	170
117		171
118	• Longer, detailed answers.	172
119		173
120	• Advanced topics, including Digital Signal Processing	174
121	(DSP), Statistics, Vector Calculus, Political Science,	175
122	Organic Chemistry, History, and Political Islam.	176
123		177
124	This subset was designed to test the models on nuanced,	178
125	high-complexity tasks requiring deeper understanding and	179
126	context-aware evaluation.	180
127		181
128	<b>3.2.3 Mathematical vs. Analytical Quizzes</b>	182
129		183
130	The dataset also included a mix of mathematical and ana-	184
131	lytical quizzes. For example:	185
132		186
133	• Mathematical quizzes, such as those in DSP and Vec-	187
134	tor Calculus, focused on numerical computation and	188
135	formulaic problem-solving.	189
136		190
137	• Analytical quizzes, such as those in Political Islam and	191
138	History, required critical thinking, interpretation, and	192
139	context-based reasoning.	193
140		194
141	This distinction allowed for a detailed evaluation of how	195
142	well the models handled problem-solving versus interpreta-	196
143	tive tasks.	197
144		198
145	<b>3.3. Prompt Complexity</b>	199
146	The evaluation process in this project utilized prompts	200
147	of varying complexity to test the capabilities of the models	201
148	in understanding and executing detailed instructions. Two	202
149	categories of prompts were used:	203
150		204
151	<b>3.3.1 Simple Prompt</b>	205
152		206
153	The simple prompt focused on minimal instruction to per-	207
154	form basic grading:	208
155		209
156	<i>“Check and grade the provided quiz answers, as-</i>	210
157	<i>signing a score out of 10.”</i>	211
158		212
159	This prompt required the model to perform straightforward	213
160	grading without additional analysis or detailed feedback,	214
161	serving as a baseline for comparison.	215

Table 1. Performance Evaluation of LLMs on Basic Quizzes

Quiz Subject	Total Questions	GPT-4o	Claude-3.5	Gemini-1.5 Pro	Actual Marks
Basic Math	6	5/6	6/6	6/6	6/6
Calculus	3	3/3	3/3	3/3	3/3
General Knowledge	8	8/8	8/8	8/8	8/8
Basic Science	7	7/7	6/7	6/7	7/7
Algebra	3	3/3	3/3	3/3	3/3
Physics	6	4/6	3/6	4/6	3/6

Table 2. Performance Evaluation of LLMs on Complex Quizzes

Quiz Subject	Total Questions	Actual Marks	Gemini-1.5 Pro	GPT-4o	Claude-3.5
Digital Signal Processing	4	4/4	2/4	4/4	4/4
Statistics	1	1/1	0/1	0/1	1/1
Vector Calculus Example 1	4	4/4	4/4	4/4	3/4
Vector Calculus Example 2	4	4/4	4/4	4/4	4/4
Organic Chemistry	3	0/3	1/3	3/3	3/3
Political Science	5 Mark Analysis	5/5	5/5	5/5	5/5
Political Islam	10 Mark Analysis	2/10	4/10	7/10	10/10
History Example 1	10 Mark Analysis	10/10	6/10	7/10	10/10
History Example 2	10 Mark Analysis	4/10	3/10	8/10	10/10

324

## 4. Analysis of Results

325

326 The performance evaluation of the Language Learning

327 Models (LLMs) for both basic and complex quizzes, as

328 shown in Tables 1 and 2, reveals key insights into their grad-

329 ing accuracy and alignment with the actual marks.

330

### 4.1. Performance on Basic Quizzes

331

332 For the basic quizzes, the subjects primarily involved

333 straightforward and fundamental mathematical and factual

334 knowledge, such as Basic Math, Calculus, General Knowl-

335 edge, Basic Science, Algebra, and Physics. Among the

336 three models:

- 337
- 338 • **GPT-4o:** Demonstrated reliable performance in most
- 339 subjects, with consistent alignment with the actual
- 340 marks in quizzes such as Basic Science (7/7) and Al-
- 341 gebra (3/3). However, minor discrepancies were ob-
- 342 served in Basic Math (5/6) and Physics (4/6), indicat-
- 343 ing occasional under-grading of correct answers.
- 344
- 345 • **Claude-3.5:** Excelled in most quizzes, achieving full
- 346 marks in subjects such as Basic Math (6/6), Calculus
- 347 (3/3), and General Knowledge (8/8). However, its per-
- 348 formance dropped in Basic Science (6/7) and Physics
- 349 (3/6), highlighting challenges in slightly more nuanced
- 350 grading scenarios.
- 351
- 352 • **Gemini-1.5 Pro:** Displayed competitive performance,
- 353 aligning with actual marks in subjects like Basic Math
- 354 (6/6) and Algebra (3/3). However, it struggled slightly
- 355 in Physics (4/6) and Basic Science (6/7), indicating
- 356 some inconsistency in recognizing correct answers in
- 357 these areas.

358 In summary, for basic quizzes, all three models per-

359 formed reasonably well, with Claude-3.5 marginally out-

360 performing the other two in terms of consistent grading ac-

361 curacy.

362

### 4.2. Performance on Complex Quizzes

363

364 The complex quizzes covered more intricate analytical

365 and domain-specific subjects, such as Digital Signal Pro-

366 cessing (DSP), Vector Calculus, Organic Chemistry, and

367 Political Science, among others. Here, the grading discrep-

368 ancies were more pronounced:

- 369
- 370 • **GPT-4o:** Performed exceptionally well in mathemati-
- 371 cally intensive tasks such as DSP (4/4) and both Vector
- 372 Calculus examples (4/4). However, it over-graded Or-
- 373 ganic Chemistry (3/3 instead of 0/3), however its per-
- 374 formance in subjective analysis-based subjects such as
- 375 Political Islam (7/10) and History Example 1 (7/10)
- 376 showed it cannot distinguish between good and bad re-
- 377 sponses.

- **Claude-3.5:** Demonstrated remarkable accuracy in an-
- alytical and fact-based tasks, achieving perfect scores
- in DSP (4/4), Statistics (1/1), and both Vector Calculus
- examples (4/4). However, it showed inconsistencies in
- subjective or descriptive subjects, such as History Ex-
- ample 1 (10/10 actual but 7/10 graded) and Political
- Islam (10/10 actual but 4/10 graded).
- **Gemini-1.5 Pro:** Struggled in certain mathematical
- tasks such as DSP (2/4) and Statistics (0/1), indicating
- challenges in identifying correct solutions. It showed
- good performance in descriptive or analytical reason-
- ing subjects like Political Islam (4/10) and History Ex-
- ample 2 (3/10). It performed well in Vector Calculus
- (4/4) and Political Science (5/5), suggesting strength
- in structured or deterministic tasks.

Overall, GPT-4o and Claude-3.5 outperformed Gemini-1.5 Pro in complex quizzes, particularly in mathematical tasks. Gemini-1.5 Pro displayed slightly better performance in descriptive subjects like History and Political Science.

### 4.3. Subject-Wise Observations

- **Mathematical Subjects:** GPT-4o and Claude-3.5 ex-
- hibited excellent performance in mathematical tasks
- such as Calculus, DSP, and Vector Calculus, con-
- sistently grading correctly. Gemini-1.5 Pro lagged
- slightly in DSP and Statistics but matched the perfor-
- mance of other models in Vector Calculus.
- **Analytical and Reasoning Tasks:** GPT-4o and
- Claude-3.5 struggled to fully align with the actual
- marks in subjects requiring reasoning or subjective
- analysis, such as Political Islam and History. Gemini-
- 1.5 Pro performed significantly better in these areas.
- **Domain-Specific Subjects:** Organic Chemistry pre-
- sented the most notable discrepancy, where GPT-4o
- and Claude-3.5 over-graded with 3/3, while the actual
- marks were 0/3. This indicates that all models strug-
- gled with the intricacies of organic chemistry grading.

### 4.4. Key Insights

The analysis reveals that:

- GPT-4o and Claude-3.5 are more reliable in grading
- mathematical and deterministic tasks.
- Gemini-1.5 Pro showed better performance in descrip-
- tive and reasoning tasks, though it occasionally under-
- graded.
- All models exhibited difficulty in grading nuanced,
- domain-specific tasks such as Organic Chemistry.

(383744), ndas Rathore (366636), hammad Mustafa (365880)	CS-471 2024 Submission #Ahmad Faisal (383744), Sundas Rathore (366636), Muhammad Mustafa (365880). CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.	(383744), Sundas Rathore (366636), Muhammad Mustafa (365880)
432	<b>5. Conclusion</b>	486
433		487
434	This project demonstrates the potential of leveraging	488
435	state-of-the-art multimodal large language models (LLMs)	489
436	for automating the evaluation of handwritten quizzes.	490
437	The integration of an automated grading system not only	491
438	streamlines the evaluation process but also reduces human	492
439	biases and inefficiencies, making it a valuable tool for edu-	493
440	cational institutions.	494
441	In conclusion, this project serves as a step toward mod-	495
442	ernizing educational assessment practices using LLMs. Fu-	496
443	ture work can explore fine-tuning these models for domain-	497
444	specific applications, incorporating additional evaluation	498
445	metrics, and addressing the limitations of current systems	499
446	to further enhance their applicability and reliability.	500
447		501
448	<b>References</b>	502
449		503
450	Towards LLM-based Autograding for Short Textual An-	504
451	swers	505
452	An Empirical Evaluation of Using Large Language Mod-	506
453	els for Automated Unit Test Generation	507
454	Performance of the Pre-Trained Large Language Model	508
455	GPT-4 on Automated Short Answer Grading	509
456	Automated Test Creation Using Large Language Mod-	510
457	els: A Practical Application	511
458	Survey of different Large Language Model Architec-	512
459	tures: Trends, Benchmarks, and Challenges	513
460	Multi-model Essay Evaluation with Optical Character	514
461	Recognition and Plagiarism Detection	515
462	Towards Scalable Automated Grading: Leveraging	516
463	Large Language Models for Conceptual Question Evalua-	517
464	tion in Engineering	518
465	Investigating Automatic Scoring and Feedback using	519
466	Large Language Models	520
467	Diverse LLM Approaches in Essay Scoring: A Compar-	521
468	ative Exploration of Many-Shot Prompting, LLM Jury Pan-	522
469	els, and Model Fine-Tuning	523
470	TrOCR: Transformer-based Optical Character Recogni-	524
471	tion with Pre-trained Models	525
472	Transforming Education with AI-Powered Automated	526
473	Grading Solution	527
474	Can Large Language Models Automatically Score Profi-	528
475	ciency of Written Essays	529
476		530
477		531
478		532
479		533
480		534
481		535
482		536
483		537
484		538
485		539