

MaPLe: Marker-guided Partial Labeling

Anonymous ACL submission

Abstract

In recent years, using large language models (LLMs) as evaluators has emerged as a new evaluation paradigm. However, when reasoning processes are complex, models often struggle to determine appropriate analysis directions, and unguided evaluation may lead to erroneous judgments. To address this, we propose a novel training strategy called MaPLe (Marker-guided Partial Labeling). This method explicitly triggers the model’s implicit reasoning paths by randomly masking prompt information, thereby guiding the reasoning direction and enhancing evaluation accuracy. To validate the method’s cross-lingual and multi-scenario adaptability, we constructed an automatic question-answering scoring chinese dataset for second language learners, Chinese-L2. Experimental results demonstrate that MaPLe achieves superior performance across multiple benchmarks and exhibits strong generalization capabilities in cross-lingual and multi-scenario data environments. Our method and related resources are released at <https://anonymous.4open.science/r/MaPLe1-60D6>

1 Introduction

In recent years, LLMs have experienced rapid development, demonstrating formidable capabilities across numerous natural language processing tasks (Bonthu et al., 2021; Mann et al.,

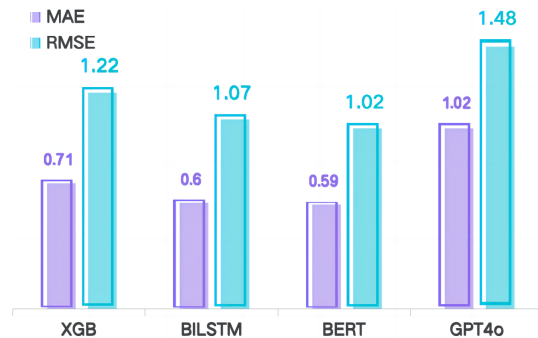


Figure 1: For the Mohler dataset, both MAE and RMSE metrics (where lower values are better) show that small models outperform large models. (Ferreira Mello et al., 2025).

2020; Ouyang et al., 2022; OpenAI, 2023). However, their performance as evaluators in automated scoring tasks—which demand fine-grained reasoning and precise numerical calibration—remains unsatisfactory, sometimes even underperforming models with smaller parameter scales (Ferreira Mello et al., 2025; Henkel et al., 2024a; Carpenter et al., 2024) (Figure 1).

Research using LLMs as evaluators can be broadly categorized into two types: the first involves directly invoking LLM for zero-shot or few-shot prompt scoring (Yuan et al., 2023); the second entails supervised fine-tuning of LLM for specific scoring tasks (Wei et al., 2022; Kojima et al., 2022; Rosset). However, both approaches exhibit inherent limi-

Dataset	Reference Answer	Scoring Criteria	Maximum Length	Average Length	Score Range	Number of Questions (essays)
Chinese-L2	×	✓	218	52	0-8	181
Mohler	×	×	1031	163	0-5	80
AES ENEM	✓	×	2004	4571	0-1000	2263

Table 1: Key information table for the Chinese-L2, Mohler and AES ENEM datasets

051	tations. Prompt-based approaches heavily rely	thinking through masking, introducing ana-	088
052	on closed-source and expensive commercial	logical reasoning chains without requiring ad-	089
053	models, posing significant challenges in cost	ditional annotations, significantly improving	090
054	control, assessment transparency, and data pri-	LLM grading accuracy.	091
055	privacy protection (OpenAI, 2023). Supervised	• We systematically evaluated the perfor-	092
056	fine-tuning methods, while achieving efficient	mance of mainstream LLMs in cross-language	093
057	and accurate scoring on specific datasets, suffer	and cross-scenario automatic grading tasks, re-	094
058	from severe limitations in generalizability	vealing both challenges and opportunities in	095
059	and flexibility. Consequently, neither paradigm	multilingual environments.	096
060	fully unleashes the inherent reasoning potential		
061	of LLMs to achieve robust and generalizable	2 Related Work	097
062	automated scoring.		
063	We contend that the core issue lies in the ab-	This paper focuses on automatic scoring tasks,	098
064	sence of explicit guidance for the model’s im-	bringing together two major mainstream sce-	099
065	PLICIT reasoning process during scoring (though	narios: Automated Essay Scoring (AES) and	100
066	methods like Chain of Thought can guide	Automated Short Answer Grading (ASAG) un-	101
067	LLMs toward deep reasoning, they cannot	der a unified perspective. It systematically	102
068	guarantee the correctness of the reasoning pro-	reviews and compares the latest research ad-	103
069	cess). To address this, we propose the MaPLE	vancements of large models in these areas.	104
070	training strategy. By masking key words in	LLMs for AES. The emergence of LLMs	105
071	the prompt template, we guide the LLM to	and their versatility across various downstream	106
072	perform implicit reasoning within “informa-	tasks (Wei et al., 2022; Kojima et al., 2022)	107
073	tion gaps” while specifying the direction of	has drawn attention to their potential for essay	108
074	reasoning. This approach transforms masking	grading. This has spurred the development of	109
075	into an internal reasoning chain, activating the	zero-shot and few-shot automated essay grad-	110
076	model’s latent scoring potential. Simultane-	ing techniques based on LLMs. Concurrently,	111
077	ously, to evaluate large models’ generalization	research has shown that incorporating context	112
078	capabilities across languages, we constructed	learning and the Chain of Thought (COT) into	113
079	Chinese-L2—the first multi-scenario question-	prompt design for automated grading can en-	114
080	answering scoring dataset tailored for Chinese	hance the effectiveness of automated essay	115
081	as a second language learners.	scoring (Lee et al., 2024a,b). Other studies	116
082	Our main contributions are listed below:	have explored prompt engineering, investigat-	117
083	• We independently constructed Chinese-	ing how different prompt templates influence	118
084	L2, the first multi-scenario question-answering	automatic scoring tasks (Zhou et al., 2022; Ye	119
085	grading dataset focused on Chinese second-	et al., 2024; Chen et al., 2024; Xue et al., 2025).	120
086	language learners.	However, LLMs employing these strategies	121
087	• We proposed MaPLE, which promotes	still perform poorly, with results comparable	122
		to classifiers relying solely on predicting text	123

length.

LLMs for ASAG. The application of LLMs for ASAG has become a current research focus. Recent studies have explored the potential of LLMs in grading short-answer questions, comparing their performance through zero-shot and few-shot settings (Henkel et al., 2024b,a). Additional research has focused on utilizing various prompt templates to reflect LLM performance in ASAG tasks (Yang et al., 2023; Kamesh, 2024; Duong and Meng, 2024; Zhao et al., 2025; Wang et al., 2025; Ferreira Mello et al., 2025). However, these studies have not demonstrated the advantages of LLMs.

3 Why Do LLMs Underperform on Scoring Tasks?

Many researchers have observed that current LLMs exhibit limitations in numerical reasoning (scoring tasks) (Lai et al., 2025). In this section, we will explore potential reasons for LLMs’ poor performance on scoring tasks and propose recommended approaches.

3.1 The Pre-training Objectives of LLMs

Common LLMs primarily optimize the “next-word prediction” objective during pretraining, granting them significant advantages in generating coherent long-form text. However, this objective predominantly learns superficial statistical correlations rather than interpretable discriminative criteria. Consequently, models often exhibit scale instability and bias in scoring tasks requiring numerical calibration and multidimensional trade-offs, making it challenging to achieve stable alignment with human evaluations.

3.2 Implicit Reasoning Capabilities Remain Underutilized

Although existing research has endowed LLMs with some automated scoring capabilities through techniques like prompt engineering,

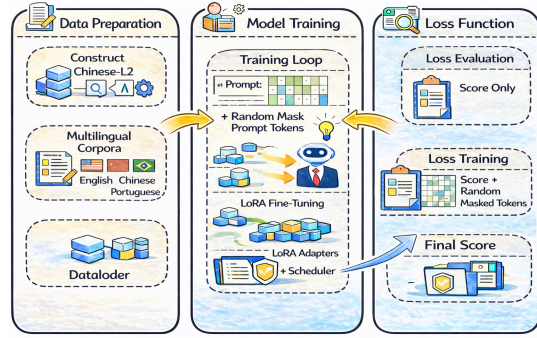


Figure 2: Data construction workflow and MaPLE method (we used LoRA Fine-Tuning with a rank of 64).

few-shot learning, and supervised fine-tuning (Zhao et al., 2023; Henkel et al., 2024a; Fan et al., 2024), these approaches often remain confined to outcome-level alignment. They fail to sufficiently activate and constrain the model’s implicit reasoning processes. Without explicit guidance on reasoning pathways, models tend to rely on superficial linguistic cues or heuristic associations during scoring. This leads to unstable scoring scales and drifting dimension weights, making reliable alignment with human evaluation criteria difficult. This issue is particularly pronounced in dimensions requiring deep semantic understanding, such as content quality and logical consistency (Seßler et al., 2024).

3.3 MaPLE

MaPLE is a straightforward plug-and-play training strategy designed to enhance the reasoning capabilities of LLMs. This approach processes input prompts (as definition in the theorem below), randomly selects a starting masking position, and subsequently masks all subsequent tokens (including the ground truth label) starting from that position. The model must then infer the content of the masked tokens. By converting the single score label into partial labeling for tokens and scores, and in-

192 incorporating masked tokens and scores into the
 193 loss function calculation, MaPLE provides ex-
 194 plicit direction for the model’s reasoning pro-
 195 cess. Masked tokens are treated as the model’s
 196 “thought trail” or “reasoning path,” enhanc-
 197 ing contextual understanding during inference.
 198 This activates latent reasoning capabilities, sig-
 199 nificantly improving inference accuracy.

Definition MaPLE Given a prompt tem-
 plate P , whose token sequence is repre-
 sented as $P = (t_1, t_2, \dots, t_n)$, a start-
 ing index q_i is randomly generated, where
 $q_i \in U \{i_1, i_2, \dots, i_n\}$, and U denotes a
 discrete uniform distribution. Each batch
 uses a different random q to better stimu-
 late LLMs’ implicit reasoning capabilities
 and generalization. For the original
 sequence $P = (t_1, t_2, \dots, t_n)$, applying
 masking yields: $M_q(P) = (t_1, \dots, t_{q-1},$
 $[Mask], [Mask], \dots, [Mask])$. The par-
 tial labeling consisting of $[Mask]$ is where
 we need to compute the loss, providing
 clear direction for model inference.

Corollary 1 (Progressive Masking Prop-
 erty) When $q = 1$, the entire prompt
 is masked, requiring the model to gener-
 ate content entirely autonomously. When
 $q = n$, the task approximates traditional
 language modeling. As q decreases from
 n to 1, reconstruction difficulty monoton-
 ically increases.

Corollary 2 (Structural Integrity Con-
 straint) Unlike traditional random sparse
 masking, this method forces the model to
 learn the structural integrity of the prompt
 template. Specifically, it must correctly
 infer the structure and semantics of the en-
 tire suffix (t_q, \dots, t_n) based on the pre-
 fix (t_1, \dots, t_{q-1}) .

3.4 Loss Function

To support our MaPLE, we modified the com-
 ponents involved in loss calculation. During
 training: we compute loss for both the random
 masking portion and the scores simultaneously.
 During validation: we mask only the scores,
 aiming to reduce interference from random se-
 lection and identify the model that most accu-
 rately outputs standard scores. Specifically, for
 the training phase, we first construct labels. We
 define tokens for which no loss is computed as
 γ . Then, for the input to the model $x_i = [x_{i1}, \dots,$
 $x_{iL}]$, we randomly determine the start position
 a_i of the masking. From a_i onwards, all subse-
 quent tokens are included in loss computation.
 We denote the interval where loss is computed
 as k_i . Thus, the labels for the training phase
 are:

$$y_{it} = \begin{cases} y, & t \in k_i \\ x_{it}, & t \notin k_i \end{cases} \quad (1)$$

For the validation phase, we design an iden-
 tifier g_i to separate the prompt content from
 the final output score. We do not compute loss
 for any content preceding g_i ; only the loss for
 the score is calculated. Thus, the labels for the
 validation phase are:

$$y_{it} = \begin{cases} \gamma, & t \leq g_i \\ x_{it}, & t > g_i \end{cases} \quad (2)$$

Therefore, the final loss function is:

$$\mathcal{L}_i = -\frac{1}{|\Omega_i|} \sum_{t \in \Omega_i} \log \frac{\exp(z_{i,t,y_{i,t}})}{\sum_{v=1}^V \exp(z_{i,t,v})} \quad (3)$$

where Ω_i denotes the valid set for loss cal-
 culation (equation 4), z_{itv} represents the logits
 for token t at position v in the vocabulary, and
 y_{it} is the target token.

$$\Omega_i = \{t \mid y_{i,t} \neq \gamma\} \quad (4)$$

Models	Mohler			AES ENEM			Chinese-L2		
	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$
Praetor(Leng et al., 2025)	1.84	2.16	0.41	-	-	-	1.08	1.93	0.74
MTS(Lee et al., 2024c)	-	-	-	193.71	215.94	0.58	2.25	2.79	0.73
Qwen3-8B(ours)	0.32	0.54	0.86	59.56	87.42	0.83	0.61	1.11	0.91
Qwen3-14B(ours)	0.30	0.55	0.86	54.86	81.25	0.85	0.58	1.05	0.92
Llama-3.1-8B-Instruct(ours)	0.36	0.59	0.83	59.31	88.15	0.83	0.62	1.10	0.92

Table 2: Comparison results with existing methods across three datasets. Δ represents MAE, $|\Delta|$ represents RMSE, ρ represents Pearson.

4 Experiment

Currently, most research on LLMs heavily relies on massive computational resources (see Appendix A.1 for comparison), posing significant challenges for universities with limited resources. Due to insufficient computational resources, we introduced LoRA fine-tuning technology into model training and designed a universal prompt template (see Appendix A.2). Figure 2 illustrates our data construction workflow and training methodology.

4.1 Evaluation Indicators

To assess the performance of our proposed method, we employ Pearson’s correlation coefficient (Pearson) to measure the trend correlation between the scores generated by LLMs and the true labels. Additionally, we calculate the Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) to quantify the gap between the LLMs’ output scores and the true labels.

4.2 Baseline

We selected the state-of-the-art open-source LLMs as baseline models: Qwen3-8B, Qwen3-14B (Yang et al., 2025), and Llama-3.1-8B-Instruct (Dubey et al., 2024). (Detailed experimental settings are provided in A.3.) Additionally, we compared with recent scoring task studies: MTS (Lee et al., 2024c) and Praetor

(Leng et al., 2025). (Detailed experimental settings for Praetor and MTS are provided in Appendix A.4. Furthermore, we selected additional recent relevant studies for comparison, as detailed in Appendix A.5.)

4.3 Datasets

Currently, automatic scoring research primarily focuses on English-language scenarios, while data resources for other languages remain relatively scarce. To address this imbalance and validate the effectiveness of the proposed method in multilingual settings, we constructed the Chinese-L2 dataset, with its construction details outlined in Appendix A.6. Additionally, we incorporate the English Mohler corpus (Mohler et al., 2011) and the Portuguese AES ENEM dataset (Silveira et al., 2024). Key information for the three datasets is summarized in Table 1, with a more detailed introduction provided in Appendix A.6.

4.4 Results: Comparison of Open-Source LLMs

Through systematic comparisons of various LLMs, we found Qwen3-14B to deliver the most competitive overall performance on automated scoring tasks. Its MAE on the Mohler, AES ENEM, and Chinese-L2 benchmark datasets was 0.30, 54.86, and 0.58, respectively. Notably, under comparable parameter counts, Qwen3-8B consistently outper-

Models	Method	Mohler			AES ENEM			Chinese-L2		
		$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$
Qwen3-8B	Base	0.34	0.62	0.82	60.92	91.06	0.82	0.69	1.21	0.90
	NoMask	0.34	0.60	0.84	60.18	89.30	0.82	0.65	1.18	0.90
	Fixed_Random	0.35	0.59	0.84	78.01	103.86	0.78	0.66	1.20	0.90
	MaPLe	0.32	0.54	0.86	59.56	87.42	0.83	0.61	1.11	0.91
Qwen3-14B	Base	0.32	0.56	0.86	63.52	88.04	0.82	0.62	1.14	0.90
	NoMask	0.35	0.62	0.81	59.93	87.00	0.83	0.59	1.10	0.91
	Fixed_Random	0.35	0.59	0.84	58.69	85.04	0.83	0.60	1.12	0.91
	MaPLe	0.30	0.55	0.86	54.86	81.25	0.85	0.58	1.05	0.92
QLlama-3.1-8B-Instruct	Base	0.39	0.70	0.81	59.46	95.88	0.79	0.66	1.15	0.91
	NoMask	0.37	0.60	0.82	60.55	93.29	0.80	0.68	1.23	0.89
	Fixed_Random	0.39	0.62	0.82	59.31	89.97	0.83	0.69	1.16	0.90
	MaPLe	0.36	0.59	0.83	59.31	88.15	0.83	0.62	1.10	0.92

Table 3: Results of ablation experiments across three datasets. Δ represents MAE, $|\Delta|$ represents RMSE, ρ represents Pearson.

formed Llama3.1-8B-Instruct. Furthermore, while larger models within the same series typically exhibit superior performance across most metrics, exceptions were observed: for instance, on the Mohler dataset, Qwen3-8B’s RMSE correlation coefficient was marginally lower than Qwen3-14B’s. Performance gains from model scale were less pronounced than anticipated—e.g., a mere 0.02 difference in MAE on Mohler and 0.03 difference on Chinese-L2—likely due to our limited dataset size. Therefore, in practical deployment scenarios, selecting models with smaller parameter counts remains viable when balancing inference efficiency and resource costs.

4.5 Results: Method Comparison

As shown in Table 2, our method demonstrates superior competitiveness compared to existing approaches such as Praetor and MTS, particularly on the AES ENEM dataset. This indicates that existing methods are often constrained by insufficient transferability, making it challenging to effectively generalize to similar tasks. In contrast, our method can be directly applied to common rating tasks without requiring task-

specific tuning, while still achieving excellent performance.

5 Analysis

This section provides an in-depth analysis of each module within MaPLe, exploring the factors behind its outstanding performance and their impact on the model. To this end, we conducted three additional comparative experiments: (1) Base: Removing the random masking mechanism and using standard cross-entropy loss; (2) NoMask: Retaining the proposed loss function but not applying random masking to the prompt; (3) Fixed_Random: Using a fixed masking region for each sample during training. The experimental results are shown in Table 3.

5.1 Synergy Between MaPLe and Loss Functions

To investigate the synergistic effect of MaPLe and the loss function in enhancing model performance, we designed ablation experiments by simultaneously removing both components for evaluation. The experimental results are shown in Table 3. After removing both MaPLe

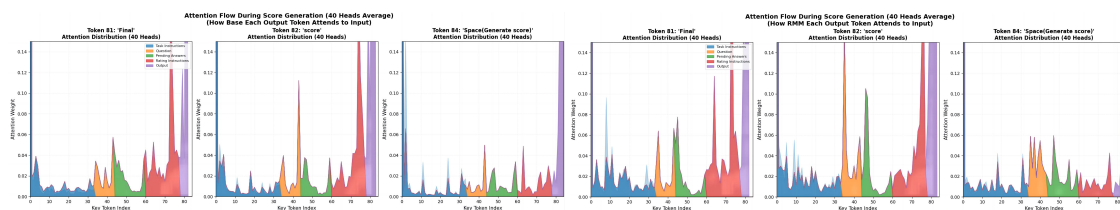


Figure 3: Comparison between Base and MaPLE. We divided the prompt into distinct regions to observe the distribution of attention weights during generation scoring, with colors representing different regions.

344 and the loss function, evaluation metrics significantly
 345 decreased across all three datasets. For instance, on the Mohler dataset, the RMSE
 346 increased by 12.62%. This result indicates a synergistic effect between MaPLE and the loss
 347 function, not only validating the advantages of MaPLE but also emphasizing the critical role
 348 of selecting an appropriate loss function during training and validation for model stability.

349 5.2 Role of MaPLE

350 NoMask and Fixed_Random are two supplementary control experiments we designed to
 351 further validate the effectiveness of the MaPLE module. As shown in Table 3, removing the
 352 MaPLE module resulted in a 6.61% increase in root mean square error (RMSE) on the AES
 353 ENEM dataset. These experiments clearly demonstrate the role of introducing random
 354 masking into prompts: the model no longer blindly memorizes prompt content but instead
 355 learns to “answer the structure.” This forces the model to engage in internal reasoning: to
 356 correctly predict masked tokens and fill the information gaps we created, the model must per-
 357 form implicit inference internally. This drives the model to more deeply understand and inte-
 358 grate unmasked contextual information, building stronger logical connections and semantic
 359 associations within the neural network. Simultaneously, masked tokens function as a “chain
 360 of reasoning” during training, simulating intermediate thought steps in human inference.
 361 This enhances the model’s ability to reason and

377 evaluate complex tasks.

378 5.3 Implicit Reasoning Analysis Capability

379 To verify whether our method genuinely enhances the implicit reasoning capabilities of
 380 large language models, we visualized the final layer attention weights of the Qwen3-14B
 381 model on the Mohler dataset, with results shown in Figure 3. By comparing the performance
 382 of the baseline model and the MaPLE model, we observe that our approach enables the
 383 model to focus not only on the output segment but also on multiple regions within the
 384 prompt. Specifically, by comparing the first three subfigures (Base model) with the last
 385 three subfigures (MaPLE model) in Figure 3, we can clearly observe this shift. This indicates
 386 that MaPLE enables the model to synthesize information from multiple parts during reason-
 387 ing to arrive at the final inference result, which requires robust reasoning capabilities. For a
 388 more detailed comparative analysis, please refer to Appendix A.7.

399 5.4 Score Distribution Analysis

400 To further validate that our method genuinely activates the model’s implicit reasoning capa-
 401 bilities rather than merely falling within the inherent scoring preference range of large lan-
 402 guage models (Zheng et al., 2023), we compared the scoring distributions generated by
 403 Qwen3-14B across three datasets with manually annotated distributions. Figure 4 demon-
 404 strates the scoring distributions for the Base and MaPLE models on the AES ENEM dataset.
 405 Figure 4 shows that the MaPLE model’s scoring distribution is more aligned with the manually
 406 annotated distribution compared to the Base model, indicating that MaPLE effectively
 407 activates the model’s implicit reasoning capabilities. This suggests that MaPLE not only
 408 improves the model’s performance but also enhances its ability to reason and generate more
 409 accurate and meaningful outputs.

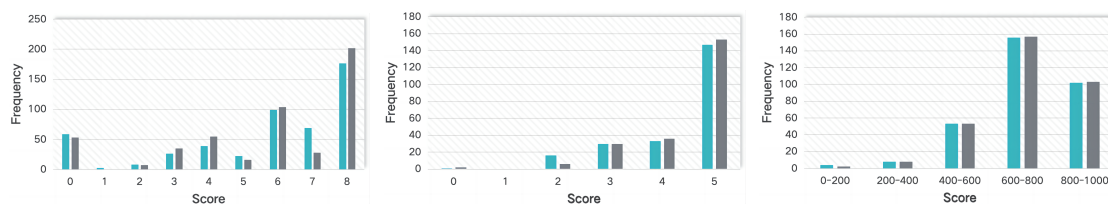


Figure 4: Distribution figures of actual versus predicted scores across the AES ENEM, Mohler, and Chinese-L2 datasets. In each figure, the left side shows human-labeled scores, while the right side displays model-predicted scores.

409 strates high consistency in distribution shapes
 410 between the two, indicating that the model
 411 is no longer constrained by its inherent scor-
 412 ing preferences and can produce results more
 413 aligned with human judgments. Furthermore,
 414 this result demonstrates that MaPLe provides
 415 the model with explicit “alignment direction”
 416 during scoring inference, prompting its scor-
 417 ing behavior to more stably align with human
 418 standards. This highlights its robustness and
 419 effectiveness in automated scoring tasks.

420 5.5 Evaluation Feedback Analysis

421 Current research also focuses on enabling
 422 LLMs to generate scoring feedback to further
 423 elucidate their interpretability. While this is not
 424 the primary focus of this paper, we addition-
 425 ally verified whether our method possesses the
 426 potential to generate scoring feedback. We per-
 427 formed zero-shot reasoning using our frozen
 428 MaPLe checkpoint by incorporating relevant
 429 prompt words to guide the model in generating
 430 scoring rationale. Appendix A.8 presents our
 431 generated prompt templates and feedback ex-
 432 amples, demonstrating our method’s capability
 433 for feedback generation.

434 6 Conclusion

435 The MaPLe method we propose is simple
 436 and efficient. It does not require designing
 437 task-specific prompts or introducing additional
 438 multi-dimensional labels, and is compatible
 439 with existing large language models. This

method randomly masks tokens in the input
 prompt, with the masked tokens forming par-
 tial labeling, providing clear direction for rea-
 soning. It also simulates "thinking chains"
 and "reasoning paths," activating the implicit
 reasoning abilities of LLMs and significantly
 improving evaluation accuracy. Furthermore,
 MaPLe is task-agnostic and can be widely ap-
 plied to various scoring tasks, from essay scor-
 ing to short-answer question evaluation, show-
 ing outstanding performance across multiple
 languages and scenarios.

452 Limitations

453 **Domain Extension Potential.** Although this
 454 study has validated the effectiveness of the
 455 MaPLe method across multiple scoring task
 456 datasets, its potential for extension to broader
 457 domains remains to be further explored. As
 458 LLMs increasingly penetrate scientific do-
 459 mains, MaPLe holds promise for migration to
 460 complex tasks such as drug discovery (Zheng
 461 et al., 2024), disease diagnosis (Zhou et al.,
 462 2024), and weather forecasting (Wang and
 463 Karimi, 2024). These directions constitute re-
 464 search topics worthy of exploration in future
 465 work.

466 **Capability to Generate Detailed Feedback.**
 467 While our method demonstrated preliminary
 468 feedback generation on the Chinese-L2 dataset,
 469 its feedback generation mechanism requires
 470 further training and validation on more com-
 471 plex tasks. Generating high-quality, inter-

472	pretable feedback remains a key research focus (Nguyen et al., 2023; Ouyang et al., 2022; Leng et al., 2025), and we are committed to advancing the practical application of LLMs in this domain.		
473			
474			
475			
476			
477	Prompt Design and Model Coverage. This study employed uniformly designed prompts for the evaluation task and has not yet systematically validated the method’s dependence on different prompts. Furthermore, constrained by computational resources, experiments were conducted solely on two open-source models Qwen3 and Llama3, without extending to larger-parameter models or closed-source models. The original intent of this work was to propose an efficient, generalizable, and resource-conserving research approach. Future efforts will further validate the method’s universality across a broader range of model scales and architectures.		
478			
479			
480			
481			
482			
483			
484			
485			
486			
487			
488			
489			
490			
491			
492	Ethics Statement		
493	We confirm that all authors of this study have adhered to the ACL ethical guidelines and recommended code of conduct. All code and datasets used in this research are publicly accessible, and we have fully cited the sources of all datasets. We believe this study poses no potential risks. It should also be noted that all research content, experimental design, data results, conclusion analysis, and academic viewpoints in this paper were independently completed and are the sole responsibility of the author. GPT-5.2 (OpenAI, 2025) was used solely to enhance the quality of textual expression and did not participate in any creative research processes nor automatically generate any substantive academic content.		
494			
495			
496			
497			
498			
499			
500			
501			
502			
503			
504			
505			
506			
507			
508			
509	References		
510	Sridevi Bonthu, S Rama Sree, and MHM Krishna Prasad. 2021. Automated short answer	grading using deep learning: A survey. In <i>International cross-domain conference for machine learning and knowledge extraction</i> , pages 61–78. Springer.	512 513 514 515
511			
		Dan Carpenter, Wookhee Min, Seung Lee, Gamze Ozogul, Xiaoying Zheng, and James Lester. 2024. Assessing student explanations with large language models using fine-tuning and few-shot learning. <i>Proceedings of the 19th Workshop on Innovative Use of NLP for Building . . .</i>	516 517 518 519 520 521
		Tzu-Lin Chang, Keng-Pei Lin, and Zong-Shun Chen. 2024. Automatic short-answer grading with a pseudo-siamese neural network.	522 523 524
		Weizhe Chen, Sven Koenig, and Bistra Dilkina. 2024. Reprompt: Planning by automatic prompt engineering for large language models agents. <i>arXiv preprint arXiv:2406.11132</i> .	525 526 527 528
		Niharika Dadu, Harsh Vardhan Singh, and Romi Banerjee. 2025. Grade guard: A smart system for short answer automated grading. <i>arXiv preprint arXiv:2504.01253</i> .	529 530 531 532
		Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv e-prints</i> , pages arXiv–2407.	533 534 535 536 537 538
		Ta Nguyen Binh Duong and Chai Yi Meng. 2024. Automatic grading of short answers using large language models in software engineering courses. In <i>2024 IEEE Global Engineering Education Conference (EDUCON)</i> , pages 1–10. IEEE.	539 540 541 542 543 544
		Zhiyuan Fan, Weinong Wang, Debing Zhang, and 1 others. 2024. Sedareval: Automated evaluation using self-adaptive rubrics. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 16916–16930.	545 546 547 548 549
		Rafael Ferreira Mello, Cleon Pereira Junior, Luiz Rodrigues, Filipe Dwan Pereira, Luciano Cabral, Newarney Costa, Geber Ramalho, and Dragan Gasevic. 2025. Automatic short answer grading in the llm era: Does gpt-4 with prompt engineering beat traditional models? In <i>Proceedings of the 15th international learning analytics and knowledge conference</i> , pages 93–103.	550 551 552 553 554 555 556 557

653	Igor Cataneo Silveira, André Barbosa, and Denis Deratani Mauá. 2024. A new benchmark for automatic essay scoring in portuguese. In <i>Proceedings of the 16th International Conference on Computational Processing of Portuguese-Vol. 1</i> , pages 228–237.	Haifeng Zhao, Yuguang Jin, and Leilei Ma. 2025. Dynamic prompt adjustment formulti-label class-incremental learning. In <i>International Conference on Brain Inspired Cognitive Systems</i> .	697 698 699 700 701
659	Hanling Wang, Banghao Chi, Yufei Wu, Kexin Chen, Di Wu, Songning Liu, Yiwei Li, Hanyan Niu, and Xiaohui Zhu. 2025. Llmaking: Adaptive automatic short-answer grading using large language models. In <i>Proceedings of the Twelfth ACM Conference on Learning@ Scale</i> , pages 105–115.	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> , 1(2).	702 703 704 705 706
666	Yang Wang and Hassan A. Karimi. 2024. Exploring large language models for climate forecasting . <i>CoRR</i> , abs/2411.13724.	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.	707 708 709 710 711 712 713
669	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	Yizhen Zheng, Huan Yee Koh, Maddie Yang, Li Li, Lauren T. May, Geoffrey I. Webb, Shirui Pan, and George Church. 2024. Large language models in drug discovery and development: From disease mechanisms to clinical trials . <i>CoRR</i> , abs/2409.04481.	714 715 716 717 718 719
675	Mingfeng Xue, Yunting Liu, Xingyao Xiao, and Mark Wilson. 2025. Automatic prompt engineering for automatic scoring. <i>Journal of Educational Measurement</i> .	Shuang Zhou, Zidu Xu, Mian Zhang, Chunpu Xu, Yawen Guo, Zaifu Zhan, Sirui Ding, Jiashuo Wang, Kaishuai Xu, Yi Fang, Liqiao Xia, Jeremy Yeung, Daochen Zha, Genevieve B. Melton, Mingquan Lin, and Rui Zhang. 2024. Large language models for disease diagnosis: A scoping review . <i>CoRR</i> , abs/2409.00097.	720 721 722 723 724 725 726
679	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. In <i>The eleventh international conference on learning representations</i> .	727 728 729 730 731 732
684	Xianjun Yang, Wei Cheng, Xujiang Zhao, Linda Petzold, and Haifeng Chen. 2023. Dynamic prompting: A unified framework for prompt tuning. <i>ArXiv</i> , abs/2303.02909.	A Appendix	733
688	Qinyuan Ye, Mohamed Ahmed, Reid Pryzant, and Fereshte Khani. 2024. Prompt engineering a prompt engineer. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 355–385.	A.1 Comparison of Computing Resources	734
693	Chenhan Yuan, Qianqian Xie, and Sophia Ananiadou. 2023. Zero-shot temporal relation extraction with chatgpt. <i>arXiv preprint arXiv:2304.05454</i> .	This study systematically reviews the computational resources required for current large language model-based automatic scoring tasks and compares our proposed MaPLe method with existing representative approaches. The specific results are summarized in Table 4. The data indicates that mainstream approaches	735 736 737 738 739 740 741

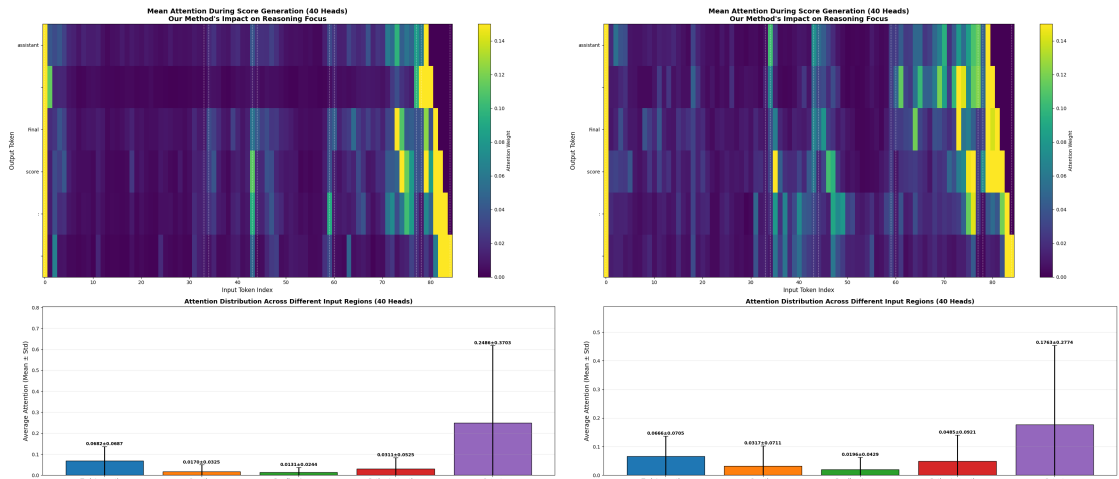


Figure 5: Comparison between Base and MaPLE. The left figure shows the Base method, while the right figure illustrates the MaPLE method. The upper section displays which tokens were focused on during inference-based score generation, while the lower section presents the mean and standard deviation of attention weights across regions (where a smaller standard deviation is preferable).

Method	computing resources
Zheng et al., 2023	8 A100 (80GB) GPUs
Duong and Meng, 2024	GPT-3.5-Turbo API
Fan et al., 2024	128 H100 (80GB) GPUs
Lee et al., 2024c	GPT-4 API
Lee et al., 2024c	1 A40 (48GB) GPUs
Lai et al., 2025	4 A800 (80GB) GPUs
Leng et al., 2025	8 A100 (80GB) GPUs
RMM(ours)	2 A40 (48GB) GPUs

Table 4: Comparison of computing resources

reducing hardware configuration and VRAM occupancy requirements. This design enables the method to maintain comparable or even superior performance to existing approaches while substantially enhancing feasibility and practical value in resource-constrained environments, providing a viable technical pathway for broader research teams to pursue related explorations.

A.2 Prompt Template

This study proposes and designs a prompt template with strong versatility. To ensure the method’s reusability and transferability, we directly apply the same template across different tasks and datasets without additional task-specific rewriting or customization. Structurally, the template integrates both comprehensive scoring criteria and non-scoring criteria formats. This enables the model to adhere to unified evaluation standards while maintaining stable comprehension of task-critical information and output constraints. Experimental results demonstrate consistent and significant

typically rely on large-scale computational infrastructure, such as multiple high-end GPUs (e.g., A100/H100) for parallel training, or high-capacity GPU memory to support full-parameter fine-tuning and large-scale inference. This imposes significant barriers for universities and research institutions with generally constrained computational resources. In contrast, the proposed MaPLE method integrates parameter-efficient fine-tuning techniques (e.g., LoRA) with lightweight prompt design. It completes all training and evaluation tasks using only two A40 GPUs, substantially

778	applicability and robustness across three exper-		
779	imental datasets, reflecting strong cross-task		
780	generalization capabilities. Detailed design		
781	specifications and template composition are		
782	presented in Tables 5 and 6.		
783		A.4.2 Detailed Experimental Setup for	817
		MTS	818
784	LLMs are deployed on 2 NVIDIA A40 GPUs	MTS A zero-shot prompting framework that	819
785	(48GB VRAM) to achieve high-performance	automatically decomposes questions/essay into	820
786	real-time processing. GPU acceleration is im-	distinct dimensions and generates scoring crite-	821
787	plemented using Python 3.12+, PyTorch 2.6,	ria for each dimension. Subsequently, through	822
788	and CUDA 12.4. Asynchronous inference is	multi-round dialogues, it guides LLMs to ex-	823
789	enabled via vLLMs to enhance throughput	tract scores for each dimension. Each dialogue	824
790	(Kwon et al., 2023). Transformers is lever-	round completes the scoring for its correspond-	825
791	aged for optimized integration of pre-trained	ing dimension based on its specific criteria,	826
792	models, facilitating model management. For	ultimately yielding a total score.	827
793	more experimental details from our study in	This paper utilizes the prompt provided by	828
794	Table 7.	MTS to invoke GPT-5.2 (OpenAI, 2025), auto-	829
		matically generating scoring criteria for each	830
795	A.4 Detailed Experimental Setup for	dimension of the AES ENEM and Chinese-L2	831
796	Praetor and MTS	datasets (experiments were not conducted on	832
		the Mohler dataset as scoring dimensions were	833
797	A.4.1 Detailed Experimental Setup for	not identified). The dimension-specific scoring	834
798	Praetor	criteria for the AES ENEM dataset are pre-	835
799	Praetor A fine-grained generative LLMs eval-	sented in Tables 10, 11, 12, 13, and 14 (these	836
800	uator trained using multi-task learning, sup-	dimensions are provided in the AES ENEM	837
801	porting instance-level customizable evaluation	paper), while the scoring criteria for each di-	838
802	criteria. It can evaluate LLMs through single-	dimension of the Chinese-L2 dataset are shown	839
803	sentence scoring or pairwise comparisons, sup-	in Tables 15 and 16 (these dimensions are ex-	840
804	porting both English and Chinese languages,	tracted from the overall scoring criteria).	841
805	while offering high flexibility in setting evalua-		
806	tion standards.	A.5 Comparison of MaPLe with Existing	842
807	In this paper, we applied the prompt design	Methods	843
808	methodology from Praetor to construct corre-	To comprehensively evaluate MaPLe’s com-	844
809	sponding prompts on the Mohler and Chinese-	petitiveness, we systematically compared it	845
810	L2 datasets, with specific designs detailed in	against multiple representative methods pro-	846
811	Table 8,9. Based on the code provided by Prae-	posed in recent years. As shown in Table 17	847
812	tor, we conducted comparative experiments.	(Mohler dataset) and Table 18 (AES ENEM	848
813	The reason for selecting only the Mohler and	dataset), we included not only state-of-the-art	849
814	Chinese-L2 datasets for experimentation is that	approaches based on large language models but	850
815	this method currently supports only English	also traditional methods leveraging lightweight	851
816	and Chinese.	models (e.g., BERT, XGBoost) to conduct	852
		a cross-paradigm, cross-model-scale perfor-	853
		mance assessment.	854
		Experimental results demonstrate that	855
		MaPLe exhibits significant advantages in au-	856
		tomated short-answer grading tasks. On the	857
		Mohler dataset, our method substantially out-	858

Prompt Template

###Task Description###

You are a professional {subject}. You will assign a reasonable score to answers regarding {subject}-related questions/essay, with the score not exceeding {score_range} points.

###Evaluation Materials###

Question content: {question/essay}

Rated answers: {answer}

###Output format###

Directly provide the final score. The output format is: Final score: {score}

###End marker and Target output###

Final score:

Table 5: Prompt template without scoring criteria

Prompt Template

###Task Description###

You are a professional {subject}. You will assign a reasonable score to answers for {subject}-related questions/essay based on the scoring criteria, with the score not exceeding {score_range} points.

###Evaluation Materials###

Scoring criteria: {scoring criteria}

Question content: {question/essay}

Rated answers: {answer}

###Output format###

Directly provide the final score. The output format is: Final score: {score}

###End marker and Target output###

Final score:

Table 6: Prompt template with scoring criteria

Parameters	Details
Batch_size	24
Epoch_num	8
Learning_rate	3e-04
Lora_rank	64
Lora_alpha	16
Lora_dropout	0.05
Temperature	0
Top_k	1
Top_p	1

Table 7: Experimental details across all datasets

performs all comparison baselines—including high-performing models like PSNN and large-model approaches such as GPT-4o—on key metrics like RMSE. This confirms that MaPLE’s effectiveness in enhancing scoring accuracy stems not only from the powerful capabilities of large models but also from its unique training design. Furthermore, on the AES ENEM dataset, MaPLE achieved the best results, even surpassing the inter-rater reliability levels reported in the original paper for this dataset, further highlighting the method’s robustness and practicality in complex scoring tasks.

A.6 Detailed Implicit Reasoning Analysis

In Section 5.3, we conducted a preliminary analysis of the information regions the model focuses on when generating scores. Here, we will elaborate on our analytical process in detail. First, we extracted the attention weights from the final layer of both the Base and MaPLE methods (using the Qwen3-14B model as the foundation) on the Mohler dataset. This layer comprises 40 attention heads. We computed the mean and standard deviation of attention across these 40 heads at their respective positions. Subsequently, we segmented the model input into five distinct regions: Task

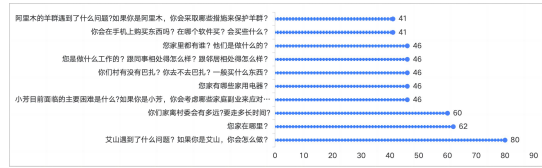


Figure 6: Question content distribution figure (showing only the top 10 most frequent questions).

Instructions, Question, Pending Answers, Rating Instructions, and Output. By quantifying the model’s attention distribution across these regions during final score generation, we examined whether it relies solely on the immediate context tokens or can effectively gather information across regions.

To delve deeper into the analysis, we further plotted heatmaps of the model’s reasoning process during score generation and calculated the standard deviation of the final attention layer (as shown in Figure 5). The lower half of Figure 5 demonstrates that our method consistently exhibits a stable low standard deviation across generated score regions. This indicates that the 40 attention heads maintain consistent reasoning paths, enabling the model to clearly identify key areas. The heatmap in the upper part of Figure 5 indicates that during scoring inference, the model simultaneously focuses on multiple information regions rather than being confined to local context, thereby supporting a more comprehensive reasoning process. These findings further validate that our approach effectively activates the latent reasoning capabilities of LLMs.

A.7 Dataset Overview

A.7.1 Construction Process of the Chinese-L2 Dataset

We selected 5,000 data entries from our research group’s projects. These entries primarily aim to assess second-language learners’ proficiency through question-answering tasks, covering multiple scenarios: “single-question

Prompt Template

Task Description ### You are a lecturer for a data structures course (computer science major). The system will present you with questions related to the data structures course and students' responses. Please evaluate the accuracy of the responses and grade them based on their correctness.

1. Write a detailed commentary on the answer. This commentary will be used to assess the answer's quality.
2. After completing the review, provide a decimal score between score_range representing your assessment of the answer. A higher score indicates better quality.
3. Output format: Your review of the response Overall score: Your rating for the answer
4. Generate content strictly adhering to this format. Do not add any additional opening, closing, or explanatory text.

User question ### question ### Answer to be evaluated ### answer

Table 8: Using the Praetor method on the prompt template for Mohler dataset

Prompt Template

Task Description ### You are a Chinese linguistics instructor. The system will present you with an open-ended question and a second language learner's response to it. Evaluate the response according to the grading criteria and assign a score based on its accuracy.

1. Write a detailed commentary on the answer. This commentary will be used to assess the answer's quality.
2. After completing the review, provide a decimal score between score_range representing your assessment of the answer. A higher score indicates better quality.
3. Output format: Your review of the response Overall score: Your rating for the answer
4. Generate content strictly adhering to this format. Do not add any additional opening, closing, or explanatory text.

User question ### question ### Answer to be evaluated ### answer ### Scoring Criteria ### scoring criteria

Table 9: Using the Praetor method on the prompt template for Chinese-L2 dataset

Formal Language

0: The language does not comply with written language standards at all, contains serious grammatical errors, frequent spelling mistakes, is extremely casual and informal, and does not meet academic or formal writing requirements in any way.

1-40: The language contains quite a few grammatical or spelling errors, some word choices are inaccurate or unsuitable for a formal style, the expression is somewhat chaotic or irregular, giving an overall unprofessional impression and lacking clear logic.

41-80: The language meets basic written language standards, with occasional grammatical or spelling errors that have little impact on overall understanding; word choices are relatively formal but not precise, and some expressions may appear unclear or lack professionalism.

81-120: The language is fairly standard, with occasional minor grammatical or spelling errors; overall expression is clear and basically meets the requirements of formal written language, with accurate, concise, and relatively professional word choices; the language style is consistent, though some parts may lack greater depth or refinement.

121-160: The language is very standard, with almost no grammatical, spelling, or logical errors; word choices are appropriate and accurate, fully meeting the requirements of written and academic language; expression is concise and clear, with a consistent and professional style, demonstrating a strong level of formality.

161-200: The language perfectly meets written language standards, with no grammatical or spelling errors; every word choice and sentence structure is extremely precise, showing a high degree of formality and academic quality; logic is tight, expression is clear and fluent, fully complying with written cultural standards, demonstrating exceptional language ability.

Table 10: Scoring criteria for the formal language dimension using the MTS method on the AES ENEM dataset

Understanding the Task

0: Completely does not understand the task requirements, and the response content is entirely off-topic. No relevant knowledge is applied, the content is incoherent, and it lacks a reasonable logical structure.

1-40: Has a shallow understanding of the task; part of the response is relevant, but there is a lack of deep understanding of key concepts. Knowledge application is insufficient, reasoning is confused, and the structure is loose or unclear, failing to effectively support the main argument.

41-80: Has some understanding of the task, but still not deep enough. Able to apply some related knowledge and develop some reasonable arguments, but reasoning is inadequate, the structure is somewhat messy or not tight, and the organization is not optimal within the constraints of an essay.

81-120: Has a good understanding of the task and can effectively analyze and argue using relevant knowledge. The theme is fairly well developed, reasoning is strong, the structure is reasonable, and content is organized well within the given framework, though there is room for improvement, such as in depth or comprehensiveness.

121-160: Has a very deep understanding of the task and can accurately apply knowledge from various disciplines for analysis, fully developing arguments. Reasoning is clear and strong, structure is tight, showing strong logic and coherence. Completes the task within the essay structure while also deepening the discussion of the topic.

161-200: Has an extremely thorough understanding of the task and can integrate concepts from multiple disciplines for profound analysis, developing arguments comprehensively from multiple angles. Reasoning is comprehensive and precise, logic is rigorous and deep, and the theme is perfectly presented within structural constraints. The text structure is flawless, content is rich and highly persuasive, fully demonstrating a high level of academic understanding and critical thinking skills.

Table 11: Scoring criteria for the understanding the task dimension using the MTS method on the AES ENEM dataset

Organization of Information

0: The selection of information is completely irrelevant or misleading, unable to effectively connect different pieces of information, and the content lacks structure. There is no clear argument or viewpoint, and the reasoning cannot support any perspective.

1-40: The selection of information is relatively random, lacking focus, poorly connected, and some information is irrelevant or repetitive. The organization is loose, the content lacks clear hierarchy or logic, and some arguments are not strongly supported.

41-80: The selection of information is basically reasonable, but there are still some irrelevant or redundant contents. The connection between information is preliminary, the structure is somewhat loose, some arguments are basically supported, but the overall reasoning is insufficient and there is room for improvement.

81-120: The selection of information is clear and mostly relevant, connections are reasonable, and key information can be effectively organized and presented. The overall structure is clear, logic is coherent, and the arguments are fairly well supported, though some reasoning may still be insufficiently in-depth or have minor gaps.

121-160: The selection of information is precise and relevant, connections are reasonable and tight, effectively supporting the arguments. The structure is clear, content is well-organized, reasoning is strong, fully supporting the position with information, facts, and viewpoints, and the expression is very clear.

161-200: The selection of information is perfect, with all information closely connected and precisely supporting the arguments. The structure is rigorous, layered, with in-depth and persuasive reasoning. The integration and organization of information demonstrate strong logic and critical thinking. The arguments are strongly supported, showing high-level writing skills.

Table 12: Scoring criteria for the organization of information dimension using the MTS method on the AES ENEM dataset

Knowing Argumentation

0: Shows no understanding of the language mechanisms of argumentation; the argument lacks clear structure, logic is chaotic, reasoning is unsupported or contains obvious logical errors. Language mechanisms are used improperly, failing to effectively convey the argument.

1-40: The argument structure is confusing or incomplete, reasoning lacks logic, and evidence is insufficient. The use of language mechanisms is somewhat stiff, reasoning or refutation is ineffective, and the argument lacks persuasiveness.

41-80: The argument structure is basically clear, but there are still some logical gaps or imperfect reasoning. Able to use basic language mechanisms, such as causality or analogy, but not fully developed; the refutation part is still insufficient.

81-120: The argument structure is relatively clear, and reasoning is fairly rigorous. Language mechanisms are used appropriately, effectively employing causality, comparison, etc., to support points. Able to respond to counterarguments, enhancing the depth and persuasiveness of the argument, though some parts may still lack depth.

121-160: The argument structure is clear and rigorous, logical reasoning is solid, and evidence is strong. Able to skillfully use language mechanisms such as causality, analogy, comparison, and refutation, enhancing the persuasiveness and depth of the argument. Handles counterarguments very reasonably, demonstrating strong argumentation ability.

161-200: The argument structure is extremely clear and rational, reasoning is flawless, points are fully supported by evidence, and logic is rigorous and perfect. Language mechanisms are used with ease, flexibly employing multiple mechanisms (such as refutation, causality, analogy, etc.) to effectively support the argument. Responses to opposing views are very strong, showing an exceptionally high level of argumentation.

Table 13: Scoring criteria for the knowing argumentation dimension using the MTS method on the AES ENEM dataset

Solution Proposal

0: The solution lacks innovation or is completely infeasible, ignores human values, has negative impacts on individuals or groups, completely disregards social and cultural diversity, and has a negligible effect on the problem.

1-40: The solution lacks effectiveness or practical feasibility, has certain negative impacts or biases, fails to adequately respect human values or the diversity of social cultures, provides insufficient details, and is difficult to implement.

41-80: The solution has some innovation and feasibility but may face challenges during implementation. It respects basic human values but does not fully consider social and cultural diversity. Implementation details are somewhat lacking and it may have adverse effects on certain groups.

81-120: The solution is fairly innovative and considers practical feasibility while providing some implementation details. It adequately respects basic human values and makes efforts to consider social and cultural diversity. It can be effectively implemented in diverse cultural contexts, but may still require further optimization or improvement.

121-160: The solution is innovative and highly feasible, with thorough details and consideration of potential implementation obstacles. It respects and embodies basic human values and has good adaptability to cultural diversity, allowing effective application across different cultural contexts.

161-200: The solution is extremely innovative and practically feasible, with detailed and comprehensive consideration of all possible obstacles, and proposes concrete, workable countermeasures. It fully respects core human values and has high adaptability in global or multicultural contexts, effectively reflecting social and cultural diversity, and has a very significant impact on solving the problem.

Table 14: Scoring criteria for the solution proposal dimension using the MTS method on the AES ENEM dataset

Completeness of Response

0: Unable to answer or provides completely irrelevant answers.

1-3: Barely answers some questions; answers are incomplete, sentences are disjointed or lack key information; frequent pauses and repetitions prevent clear communication of all key points.

4-6: Answers all questions but with simple structure; answers lack detail or fail to fully cover some points. Expression is relatively concise, with occasional pauses, repetitions, or grammatical errors.

7-8: Effectively answers all questions using one or more complete sentences, with substantial and detailed content.

Table 15: Scoring criteria for the completeness of response dimension using the MTS method on the Chinese-L2 dataset

Answer Completeness

0: The answer is completely off-topic or does not address the question at all.

1-3: The answer has a weak connection to the question, contains some irrelevant information, and features significant pauses and repetition.

4-6: The answer is generally on topic, conveying the main information, but may be somewhat off-topic or not fully aligned with the question.

7-8: The answer is entirely on topic, clearly and accurately expressed, and closely aligned with the core of the question.

Table 16: Scoring criteria for the answer completeness dimension using the MTS method on the Chinese-L2 dataset

Models	Mohler		
	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$
GPT4o (Ferreira Mello et al., 2025)	0.69	1.01	-
Upstage Solar Minin (Dadu et al., 2025)	-	1.08	-
Gemini 1.5 Flash (Dadu et al., 2025)	-	1.12	-
GPT4o Minin (Dadu et al., 2025)	-	1.02	-
Upstage Solar Minin (Dadu et al., 2025)	-	0.93	-
PSNN (Chang et al., 2024)	-	0.705	0.774
BERT (Ferreira Mello et al., 2025)	0.64	1.11	-
XGB (Ferreira Mello et al., 2025)	0.69	1.18	-
Qwen3-8B (ours)	0.32	0.54	0.86
Qwen3-14B (ours)	0.30	0.55	0.86
Llama-3.1-8B-Instruct(ours)	0.36	0.59	0.83

Table 17: Results of comparing other methods with the Mohler dataset

Models	AES ENEM		
	$\Delta \downarrow$	$ \Delta \downarrow$	$\rho \uparrow$
Grader A (Silveira et al., 2024)	-	264.40	-
Grader B (Silveira et al., 2024)	-	237.55	-
Qwen3-8B(ours)	59.56	87.42	0.83
Qwen3-14B(ours)	54.86	81.25	0.85
Llama-3.1-8B-Instruct(ours)	59.31	88.15	0.83

Table 18: Results of comparing other methods with the AES ENEM dataset

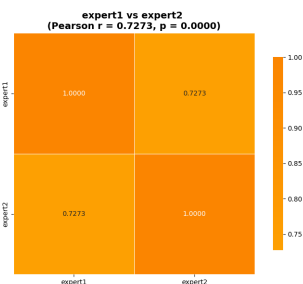


Figure 7: Heatmap of correlation among figure.

single-answer,” “multiple-question multiple-answer,” and “contextual Q&A.” The ratio of these scenarios is 3:5:2. Each question is scored out of 8 points. The dataset comprises 181 distinct, fully open-ended questions spanning daily life, social culture, and other domains. Figure 6 displays the top 10 most frequently occurring questions.

This dataset was constructed through manual annotation. First, 5,000 candidate samples were selected from within the research team. Subsequently, the samples were randomly and independently assigned to two groups of domain experts for double-blind annotation based on a unified scoring standard. After annotation, the Pearson correlation coefficient between the two groups’ scores was calculated and visualized in a heatmap (Figure 7). Results demonstrated a high positive correlation between expert evaluations. For samples with score discrepancies of ≥ 2 points, third-party experts were engaged for review and refinement, ulti-

mately forming the final dataset. This dataset encompasses diverse scenarios and broad problem coverage, enabling comprehensive and nuanced assessment of examinees’ overall capabilities.

Additionally, it should be noted that our dataset complies with the relevant regulations of the ethics review committee, does not involve any personal information, and paid data annotation experts \$2,800.

A.7.2 Mohler Dataset Overview

The Mohler Dataset is an English-language automated short-answer grading dataset comprising 80 questions sourced from assignments in a data structures course. These questions span ten assignments and two exams, encompassing 2,273 student responses. Two graduate students in related fields independently graded these answers using a scale from 0 (completely incorrect) to 5.

A.7.3 AES ENEM Dataset Overview

The AES ENEM dataset is a Portuguese-language automated essay grading dataset comprising 3,586 essays sourced from the websites Educação UOL and Brasil Escola. Two experienced annotators independently scored these essays across five dimensions, each with a maximum score of 200 points, resulting in a total possible score of 1,000 points.

A.8 Feedback Prompt Template and Examples

To validate the potential of our method in generating feedback, we directly invoked the trained parameters for zero-shot reasoning. Tables 19 and 20 present the prompt templates employed, while Tables 21 showcase feedback examples generated across two datasets. The results demonstrate that the model outputs high-quality feedback without additional fine-tuning, fully illustrating how our approach

983 effectively taps into the implicit reasoning ca-
984 pabilities of LLMs.

Prompt Template

###Task Description###

You are a professional {subject}. You will assign a reasonable score to answers regarding {subject}-related questions/essay, with the score not exceeding {score_range} points. Not only should you provide the final score, but you should also explain the rationale behind it.

###Evaluation Materials###

Question content: {question/essay}

Rated answers: {answer}

###Output format###

Provide the final score and the rationale for the scoring. The output format is: Final score: {score}. Reasons for the scoring: {Reasons}

###End marker and Target output###

Final score: Reasons for the scoring:

Table 19: Feedback Prompt template without scoring criteria

Prompt Template

###Task Description###

You are a professional {subject}. You will assign a reasonable score to answers for {subject}-related questions/essay based on the scoring criteria, with the score not exceeding {score_range} points. Not only should you provide the final score, but you should also explain the rationale behind it.

###Evaluation Materials###

Scoring criteria: {scoring criteria}

Question content: {question/essay}

Rated answers: {answer}

###Output format###

Provide the final score and the rationale for the scoring. The output format is: Final score: {score}. Reasons for the scoring: {Reasons}

###End marker and Target output###

Final score: {score}. Reasons for the scoring: {Reasons}

Table 20: Feedback Prompt template with scoring criteria

Feedback Examples

Question: Have you been to Tiananmen Square?

Answer: I went.

Final score: 6

###Reasoning for Scoring: ###

The response is on topic, clearly expressed, with minimal pauses, repetitions, and linguistic errors. However, the content is incomplete, lacking subjects and objects. Simply stating “I went” is insufficient. If the question asked how one went to Tiananmen Square, the answer would be even less appropriate. Since the question only asks whether the person has been there, the content is fundamentally on topic and clearly expressed, warranting a score of 6.

Table 21: Feedback examples on the Chinese-L2 dataset