

# Scaling Equitable Reflection Assessment in Education via Large Language Models and Role-Based Feedback Agents

Chenyu Zhang<sup>1\*</sup>, Xiaohang Luo<sup>2</sup>

<sup>1</sup>Harvard University

<sup>2</sup>University of Pennsylvania

chenyu\_zhang@alumni.harvard.edu, xiaohl@upenn.edu

## Abstract

Formative feedback is widely recognized as one of the most effective drivers of student learning, yet it remains difficult to implement equitably at scale. In large or low-resource courses, instructors often lack the time, staffing, and bandwidth required to review and respond to every student reflection, creating gaps in support precisely where learners would benefit most. This paper presents a theory-grounded system that uses five coordinated role-based LLM agents (Evaluator, Equity Monitor, Metacognitive Coach, Aggregator, and Reflexion Reviewer) to score learner reflections with a shared rubric and to generate short, bias-aware, learner-facing comments. The agents first produce structured rubric scores, then check for potentially biased or exclusionary language, add metacognitive prompts that invite students to think about their own thinking, and finally compose a concise feedback message of at most 120 words. The system includes simple fairness checks that compare scoring error across lower and higher scoring learners, enabling instructors to monitor and bound disparities in accuracy. We evaluate the pipeline in a 12-session AI literacy program with adult learners. In this setting, the system produces rubric scores that approach expert-level agreement, and trained graders rate the AI-generated comments as helpful, empathetic, and well aligned with instructional goals. Taken together, these results show that multi-agent LLM systems can deliver equitable, high-quality formative feedback at a scale and speed that would be impossible for human graders alone. The approach demonstrates how structured agent roles, fairness checks, and learning-science principles can work together to support instructors while preserving pedagogical intent. More broadly, the work points toward a future where feedback-rich learning becomes feasible for any course size or context, advancing long-standing goals of equity, access, and instructional capacity in education.

## Code, Prompts, and Anonymized Data —

<https://github.com/CharlieChenyuZhang/equitable-reflection-assessment>

## Introduction

Timely, high-quality formative feedback (Juwah et al. 2004; Pishchukhina and Allen 2021; Morris, Perry, and Wardle

\*Corresponding author.  
Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

2021) is one of the most powerful levers for closing achievement gaps in education, yet it remains out of reach for large classes and low-resource programs. Decades of scholarship show that rich, dialogic feedback can raise learning gains by up to 0.7 effect-size points (Black and Wiliam 1998; Hattie and Timperley 2007). Unfortunately, instructors cannot read and respond to every learner reflection when enrolments scale into the hundreds or when teaching assistants are scarce. The resulting feedback gap disproportionately undermines the academic growth and motivation of historically marginalized students, further exacerbating existing disparities in educational attainment (Nicolai, Koenka, and Braxton 2023).

Large language models (LLMs) offer a tantalising alternative: they can read and comment on text at super-human speed without task-specific fine-tuning. Early studies already suggest that GPT-4 can achieve human-level agreement in automated essay scoring when supplied with appropriate rubrics and prompt constraints (García-Varela et al. 2025). However, when an explicit, teacher-designed rubric is absent, LLMs have been shown to over-emphasize superficial presentation cues or implicitly apply value-laden constructs that fall outside the course’s intended learning outcomes. These risks highlight the necessity of grounding any AI-assisted grading workflow in a shared reference framework, whether through detailed rubrics, constrained prompting, or post-hoc human calibration, to ensure that automated evaluations reinforce rather than undermine pedagogical objectives (Ouyang, Dinh, and Xu 2023; Brown 2022).

We address these gaps with a theory-grounded, multi-agent pipeline that (1) produces *robust, equitable* rubric scores and (2) delivers *dialogic* feedback aligned with learning-science theory. By integrating advances in natural language processing with principles from the learning sciences, our approach provides formative feedback that is accurate and scalable while remaining pedagogically grounded and attentive to equity.

## Background & Related Work

*“Learning without reflection is a waste. Reflection without learning is dangerous.”*

Confucius

**Formative Feedback** Formative assessment serves as an ongoing gauge of student learning, supplying timely feedback rather than waiting for end-of-course summative judgments. High-quality formative feedback promotes self-regulated learning and narrows achievement gaps when it is specific, actionable, and dialogic (Juwah et al. 2004; Pishchukhina and Allen 2021; Morris, Perry, and Wardle 2021). Many studies have examined the use of formative feedback during the writing process. Providing students with formative feedback while they write is a key instructional practice that helps them improve as writers (Graham, Harris, and Hebert 2011; MacArthur 2016). By clearly communicating what high-quality performance looks like and how to achieve it, formative feedback directs students toward productive action and improvement in specific writing skills (Graham et al. 2012; Panadero et al. 2023). To alleviate the time-intensive nature of essay evaluation, researchers have applied natural language processing and artificial intelligence techniques (Grimes and Warschauer 2010; Graham, Hebert, and Harris 2015; Roscoe et al. 2017). However, these systems remain less reliable than expert human raters, and developing them demands extensive technical expertise and large, prompt-specific training corpora (Chen, Hebert, and Wilson 2022; Moore and MacArthur 2016). Recent work has started exploring the use of GPT-4 to generate written feedback for English language learners, highlighting its potential to provide timely, detailed, and structured responses (Carlson, Pack, and Escalante 2023).

**Automated Grading Systems.** Automated grading has a long history in educational technology, originating in the 1960s with formative work on programming assignment evaluation (Hollingsworth 1960) and automatic essay scoring (Page 1966). Early systems relied on deterministic pipelines, including unit-testing harnesses for code (Messer et al. 2024), rule-based heuristics (Liu et al. 2019; Uree II and Wallace 2019), and ensemble strategies that used stacking with domain adaptation to transfer across tasks (Heilman and Madnani 2013). The rapid progress of deep learning has fundamentally reshaped this landscape: neural models now routinely outperform traditional pipelines in both reliability and generalization (Riordan et al. 2017).

The rapid expansion of large language models (LLMs) in natural language processing has sparked their integration into automated grading pipelines, where prompt-based evaluation is aligned with explicit marking schemes. Researchers have successfully applied few-shot fine-tuning or in-context demonstration to a range of state-of-the-art models, including BERT (Sung et al. 2019), and GPT-4 (Chiang et al. 2024; Impey et al. 2025; Xiao et al. 2025). Across a range of essay and short-answer benchmarks, prior work has demonstrated grading accuracy that is comparable to, and in some cases better than, traditional baselines, with expert educators validating many of these gains. In comparison, fully *zero-shot* prompting is still under-examined and has not yet shown the same level of reliability. Early investigations using simple prompts on university-level courses show promise for low-stakes formative tasks yet fall short on comprehensive summative assessments such as final ex-

aminations (Kortemeyer 2023; Yeung et al. 2025).

**Gap and Contribution.** Most LLM-based grading systems concentrate on overall score agreement and leave broader social-impact goals unaddressed. These goals include equitable treatment of learners at different proficiency levels, feedback that supports motivation and self-regulation, and workflows instructors can audit, interpret, and adapt. Prior work on multi-agent or equity-aware LLM pipelines has examined cultural alignment, feedback quality, or bias detection in isolation, rather than as part of an end-to-end reflection-assessment system suitable for real classrooms. To our knowledge, no existing approach integrates stable rubric scoring, bias-aware dialogic feedback, and explicit fairness evaluation for reflective writing into a single, scalable workflow. **This work contributes:** (i) a scalable, self-consistent scoring pipeline that produces auditable rubric scores with minimal human supervision; and (ii) a role-based agent ensemble that generates bias-aware, conversational feedback aligned with formative-assessment and metacognitive principles.

## Problem Formulation

High-quality, *dialogic* formative feedback is one of the most effective ways to narrow achievement gaps, yet instructors in large or low-resource courses cannot realistically review and respond to every learner reflection. Addressing this challenge requires an approach that produces rubric-aligned scores that are accurate and fair, along with narrative comments that are pedagogically useful for learners, all without increasing the workload placed on instructors or creating new barriers for students who are already underserved. To formalize this challenge, we describe the reflection-assessment task and introduce notation that supports precise evaluation of scoring accuracy, equity, and feedback quality.

**Task Definition.** We cast reflection assessment as a **two-output prediction task** with explicit equity constraints:

- (a) a *rubric scorer*  $f : \mathbb{X} \rightarrow [0, 3]$  that predicts an ordinal score for each of the four dimensions in Table 1; and
- (b) a *feedback generator*  $g : (x_i, f(x_i)) \rightarrow \mathbb{T}$  that emits a concise, learner-facing comment  $t_i \in \mathbb{T}$  of no more than 120 words, aligned with the five quality criteria in Table 2.

**Notation.** Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  be a corpus of  $N$  reflections where:

- $x_i \in \mathbb{X}$  represents the reflection text, and
- $y_i \in \{0, \dots, 3\}^4$  represents expert rubric scores on four dimensions.

A disjoint evaluation set  $\mathcal{E} = \{(x_j, y_j, q_j)\}_{j=1}^M$  additionally contains  $q_j \in [1, 5]^5$ , the mean Likert ratings (Correctness, Alignment, Actionability, Depth, Empathy) that three trained human graders assign to the AI-generated comment for reflection  $x_j$ .

**Objective 1: Scoring Accuracy.** We quantify how closely the system’s predictions align with expert labels using two complementary statistics: absolute deviation and ordinal agreement.

1. **Mean Absolute Error (MAE).** Let the rubric comprise  $D = 4$  scored dimensions. For reflection  $x_i$ , denote the model-predicted score on dimension  $d$  by  $f_d(x_i)$  and the corresponding majority-vote human score by  $y_{i,d}$ . The per-dimension MAE is

$$\text{MAE}_d = \frac{1}{N} \sum_{i=1}^N |f_d(x_i) - y_{i,d}|, \quad (1)$$

$$\text{MAE}_{\text{overall}} = \frac{1}{D} \sum_{d=1}^D \text{MAE}_d, \quad (2)$$

where Eq. (1) averages the absolute error for a *single* rubric dimension and Eq. (2) averages those  $D$  values to yield one scalar. Lower scores indicate tighter model–human alignment;  $\text{MAE} = 0$  denotes perfect scoring fidelity.

2. **Quadratic weighted kappa (QWK).** Let  $C = \{0, \dots, K\}$  be the set of ordinal score categories,  $O_{ij}$  the observed co-occurrence matrix,  $E_{ij}$  the expected matrix under independence, and  $w_{ij} = \frac{(i-j)^2}{(K-1)^2}$  the quadratic weight. The kappa statistic is

$$\text{QWK} = 1 - \frac{\sum_{i,j} w_{ij} O_{ij}}{\sum_{i,j} w_{ij} E_{ij}}, \quad (3)$$

ranging from  $-1$  (complete disagreement) to  $1$  (perfect agreement); higher values signify stronger ordinal concordance.

**Objective 2: Equity.** Even though we do not collect sensitive demographic attributes, we can still operationalize *fair treatment* by requiring the scorer to be equally accurate for learners at different proficiency levels. Concretely, we follow established grading-fairness research that compares performance across *score bands* (Liu and Dorans 2016). Let  $\mathcal{B}$  partition  $\mathcal{D}$  into two bands based on the *human* rubric score  $y_i$ : *low* ability (0–1 points) and *high* ability (2–3 points). Fairness is quantified by the worst-band error gap

$$\Delta_{\text{MAE}} = \max_{b \in \mathcal{B}} |\text{MAE}_b(f) - \text{MAE}_{-b}(f)|. \quad (4)$$

where  $\text{MAE}_b$  is computed over reflections whose true score falls in band  $b$ . Minimizing  $\Delta_{\text{MAE}}$  ensures that lower-scoring learners, who are most in need of formative support, are not systematically under- or over-scored by the automated system.

**Objective 3: Feedback Usefulness.** For each generated comment  $t_j = g(x_j, f(x_j))$ , three trained graders rate its quality on the five Likert-scaled dimensions in Table 2,

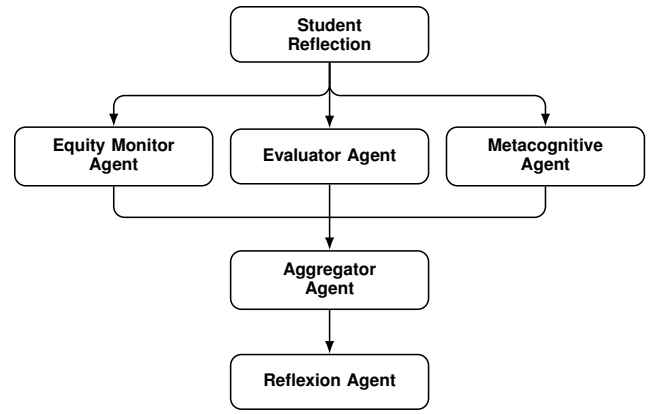


Figure 1: Role-based, multi-agent workflow for equitable reflection assessment.

yielding a vector  $q_j \in [1, 5]^5$ . Define the aggregate quality score

$$Q(g) = \frac{1}{M} \sum_{j=1}^M \frac{1}{5} \sum_{d=1}^5 q_{j,d}. \quad (5)$$

where  $M = |\mathcal{E}|$  and  $d$  indexes the rubric dimensions (*Correctness, Alignment, Actionability, Depth, Empathy*). The objective is to **maximize**  $Q(g)$  while (i) honoring the 120-word length cap and (ii) ensuring that neither the overall error  $\text{MAE}(f)$  nor the equity gap  $\Delta_{\text{MAE}}$  degrades as defined in Objectives 1–2.

**Research Questions.** Guided by the goals of accuracy, equity, and practicality, we investigate three core research questions:

**RQ1: Scoring Accuracy and Fairness** Can a zero-shot LLM scorer match expert graders on both the overall MAE and the equity gap  $\Delta_{\text{MAE}}$ ?

**RQ2: Feedback Usefulness** Does the role-based agent ensemble achieve a mean quality rating of  $Q(g) \geq 4.0$  on the five-point scale?

**RQ3: Practicality** Are the time and cost requirements of the pipeline low enough to support deployment in large or resource-constrained courses?

No reinforcement learning or fine-tuning is performed. All automation relies on *prompt-based* GPT-4o calls whose outputs are evaluated against *human-expert rubrics* for both scoring and feedback quality, aligning the evaluation with the social-impact emphasis on equity, transparency, and pedagogical value.

## Methodology

Our pipeline is implemented as a role-based, multi-agent workflow orchestrated with AutoGen (Wu et al. 2024) (Figure 1). Five specialized GPT-4o agents collaborate asynchronously so that each reflection emerges with both per-dimension rubric scores and a concise, dialogic feedback comment. The execution logic is summarized in Algorithm 1.

---

**Algorithm 1: Four-Agent Feedback Pipeline**

---

**Require:** reflection  $x$ , rubric  $R$   
**Ensure:** score  $s$ , comment  $t$   
1:  $\mathcal{A} \leftarrow \{\text{Evaluator, Equity, Meta}\}$   
2: **for all**  $a \in \mathcal{A}$  **in parallel do**  
3:      $o_a \leftarrow a.\text{GENERATE}(x, R)$   
4:  $s \leftarrow \text{PARSESCORE}(o_{\text{Evaluator}})$   
5:  $t \leftarrow \text{AGGREGATOR.COMBINE}(o_{\text{Evaluator}}, o_{\text{Equity}}, o_{\text{Meta}})$   
6:  $(c, \delta) \leftarrow \text{REFLEXION}(t)$   
7: **if**  $c = \text{REVISE}$  **then**  
8:      $t \leftarrow \text{REVISE}(t, \delta)$   
9: **return**  $(s, t)$

---

## Role Configuration

Each agent role is motivated by principles from formative assessment, equity-centered pedagogy, and multi-agent LLM design. The workflow separates responsibilities to promote transparency, reduce compounding errors, and allow instructors to audit or adjust individual components. The five roles are:

1. **Evaluator Agent** applies the four dimension rubric to the raw reflection text and returns a structured JSON object. The output includes (i) an integer score from 0 to 3 for each rubric dimension, (ii) a brief natural language explanation of the score in a `reasoning` field, and (iii) an `areas_for_improvement` list that highlights specific suggestions for the learner. This structured format provides an explicit, auditable record of how each score is assigned and why.
2. **Equity Monitor Agent** reviews the Evaluator’s narrative for biased, exclusionary, or culturally insensitive phrasing and proposes revisions. This step supports equitable treatment across learners by reducing unintended linguistic or evaluative bias.
3. **Metacognitive Agent** generates one or two reflective prompts that encourage the learner to examine their reasoning and plan next steps. This aligns the workflow with research in self-regulated learning and dialogic feedback.
4. **Aggregator Agent** synthesizes the outputs of the previous agents into a concise, learner-facing comment capped at 120 words. By highlighting only a small number of actionable next steps, the Aggregator helps prevent feedback overload (Underwood and Tregidgo 2006; Black and Wiliam 2009; Grimes and Warschauer 2010).
5. **Reflexion Agent** performs a lightweight post-hoc check, returning `CONFIDENT` or `REVISE` along with targeted suggestions. This final layer improves reliability by catching inconsistencies or omissions before feedback is released.

Generation temperature is fixed at 0.3 to balance determinism and nuance.

## Experiments

This section details the empirical set-up used to validate the proposed reflection-assessment pipeline.

**Dataset.** Our analysis draws on 336 written reflections produced by 28 adult learners enrolled in an online, synchronous AI literacy program (approximately 12 reflections per learner). Participants ranged in age from 18 to 22 years. All reflections were written in American English, collected under an approved IRB protocol, and released with informed participant consent.

**Evaluation Metrics.** We evaluate scoring accuracy, fairness, and feedback quality using four complementary families of metrics:

- (1) **Score alignment.** Agreement between the *system* score vector and the *human-reference* vector (the majority vote across three rubric-trained annotators) is assessed using two statistics:
  - **Mean Absolute Error (MAE)**, defined in Eqs. (1) and (2), which measures the average  $L_1$  deviation across the four rubric dimensions.
  - **Quadratic Weighted Kappa (QWK)**, defined in Eq. (3), which quantifies chance-corrected ordinal concordance.
- (2) **Inter-rater reliability.** We compute the two-way mixed intraclass correlation coefficient  $\text{ICC}^{(2,1)}$  for both Human–Human and Human–AI pairs to assess consistency at the item level.
- (3) **Fairness disparity.** Following Eq. (4), we measure the error gap  $\Delta_{\text{MAE}}$  across learner proficiency bands to evaluate whether scoring accuracy is comparable for lower- and higher-scoring learners.
- (4) **Feedback quality.** We summarize grader judgments using the aggregate usefulness score  $Q(g)$  defined in Eq. (5), which averages ratings across the five feedback-quality dimensions.

**Participants and Compensation.** Three human annotators were recruited and compensated for their work. Each annotator received a unique Google Sheet along with an instructional video that explained the annotation procedure, rubric dimensions, and task expectations. Annotators were required to view the full video before beginning the task.

**Sampling and Annotation Design.** To balance temporal coverage with cognitive load, we sampled reflections from three key points in the course: Class 1 (early), Class 6 (mid-point), and Class 12 (final). For each of the 28 learners, one reflection from each selected class was included, yielding a total of  $3 \times 28 = 84$  reflections. Each annotator completed six tasks: grading human-written reflections from Classes 1, 6, and 12, and evaluating AI-generated feedback for the same three classes. All annotations were submitted directly in the assigned Google Sheet using the shared four-dimension rubric.

**Rubrics.** The four-dimension grading rubric (Table 1) was identical for humans and AI. After viewing each AI-generated comment, annotators rated its quality on five Likert-scaled dimensions (Table 2); anchors followed best practice in educational-feedback research.

Dimension	Descriptor
Concept Understanding	Accurate, nuanced explanation (3)
	Mostly clear (2)
	Partial/confused (1)
	Missing or off-topic (0)
Real-World Application	Specific, thoughtful (3)
	Reasonable or generic (2)
	Vague (1)
	None (0)
Reflection & Questions	Insightful question or challenge (3)
	Identifies a question (2)
	Surface-level (1)
	None (0)
Clarity & Communication	Clear, polished (3)
	Minor issues (2)
	Hard to follow (1)
	Incoherent (0)

Table 1: Four-dimension grading rubric (0–3 points per dimension).

## Results & Analysis

### RQ1: Scoring Accuracy & Fairness

**Scoring fidelity.** Table 3 reports the mean absolute error (MAE) between model predictions and majority-vote human scores for each rubric dimension ( $n = 84$  reflections per dimension,  $n = 336$  scores in total). As shown, the model aligns most closely with human graders on *Concept Understanding* (MAE = 0.381) and least on *Real-World Application* (MAE = 0.560), yielding an overall average MAE of 0.467.

**Ordinal concordance.** Table 4 lists quadratic weighted kappa (QWK) scores per class and rubric dimension; Table 5 pools the three classes to give the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) for each dimension and overall. Across classes, the model achieves the strongest ordinal agreement with human graders on *Reflection Questions* ( $\mu = 0.483$ ), while *Concept Understanding* lags behind ( $\mu = 0.298$ ). The overall QWK of 0.459 indicates moderate concordance with expert judgment.

**Inter-rater reliability.** Table 6 reports two-way mixed intraclass correlation  $ICC^{(2,1)}$  for each rubric dimension and overall, comparing (i) consistency among human graders and (ii) consistency between the model and the human majority vote ( $n = 12$  pairwise comparisons in each case).

Human graders exhibit *moderate* internal consistency overall ( $ICC^{(2,1)} = 0.573 \pm 0.345$ ), with their strongest agreement on *Reflection Questions* ( $ICC = 0.677$ ). Model–human reliability is lower ( $ICC^{(2,1)} = 0.410 \pm 0.132$ ), yet the model aligns most closely with humans on *Real-World Application* ( $ICC = 0.487$ ). The largest divergence occurs on *Concept Understanding* ( $ICC = 0.305$ ), suggesting that additional calibration on that dimension could further narrow the human–AI reliability gap.

Dimension	Criterion
Correctness	Factually accurate; avoids hallucinations.
Alignment with Rubric	Closely reflects the official grading criteria.
Actionability	Offers specific, constructive suggestions for improvement.
Depth of Insight	Demonstrates nuanced, critical understanding.
Empathy & Tone	Supportive, respectful, and learner-appropriate.

Table 2: Rubric for evaluating AI-generated feedback (1 = poor, 5 = excellent).

Rubric Dimension	MAE
Concept Understanding	0.381
Real-World Application	0.560
Reflection Questions	0.500
Clarity of Communication	0.429
<b>Overall (mean)</b>	<b>0.467</b>

Table 3: Mean absolute error (lower is better) between the model’s rubric scores and the human reference.

**Fairness across proficiency bands.** To test whether overall model error differs between students of different ability levels, we first compute each learner’s *average human rubric score*  $\bar{y}_i = \frac{1}{D} \sum_{d=1}^D y_{i,d}$  across the four dimensions. Reflections are then split into a *low-ability* band ( $\bar{y}_i \in \{0, 1\}$ ) and a *high-ability* band ( $\bar{y}_i \in \{2, 3\}$ ). For every dimension we report the model’s MAE within each band and the worst-band error gap  $\Delta_{MAE}$  from Eq. (4).

Because ability is measured at the *reflection level*, not at the level of individual dimensions, the values in Table 7 indicate whether AI grading errors differ when comparing lower-performing and higher-performing students overall. The model shows lower accuracy for low-ability reflections in three dimensions, most sharply on *Real-World Application* ( $\Delta_{MAE} = 1.10$ ). Only *Clarity of Communication* reverses the trend, showing slightly higher error for high-ability learners. The aggregate gap ( $\Delta_{MAE} = 0.50$ ) signals that students who most need formative support still face larger grading errors.

This pattern highlights the importance of fairness-sensitive evaluation in automated feedback systems. By making the error disparity explicit and measurable, the framework enables instructors and system designers to identify where inaccuracies concentrate, monitor their magnitude over time, and apply targeted calibration strategies that improve accuracy for lower-scoring learners without degrading performance elsewhere.

Class	CU	RWA	RQ	CC	Overall
1	0.455	0.546	0.382	0.208	0.464
6	0.302	0.395	0.525	0.389	0.463
12	0.138	0.494	0.541	0.450	0.449

Table 4: Quadratic weighted kappa (higher is better) between model and human ordinal scores for each class: CU = Concept Understanding, RWA = Real-World Application, RQ = Reflection Questions, CC = Clarity of Communication.

Rubric Dimension	$\mu$ (QWK)	$\sigma$
Concept Understanding	0.298	0.158
Real-World Application	0.479	0.077
Reflection Questions	0.483	0.088
Clarity of Communication	0.349	0.126
<b>Overall</b>	<b>0.459</b>	<b>0.008</b>

Table 5: Pooled QWK across all classes ( $n = 84$  reflections per dimension).

## RQ2: Feedback Usefulness

**Aggregate quality.** Table 8 presents the five-dimension usefulness ratings (1–5 Likert) assigned by three trained graders to each generated comment, summarized as mean  $\pm$  SD. The overall  $Q(g) = 3.967$  (Eq. 5) narrowly misses the target threshold of 4.0.

The role-based ensemble generates feedback perceived as highly *empathetic* ( $\mu = 4.223$ ) and *factually correct* ( $\mu = 4.080$ ), but falls short on *actionability* ( $\mu = 3.760$ ) and *depth of insight* ( $\mu = 3.845$ ). The near-miss on the 4.0 aggregate target suggests that tightening prompts to elicit more concrete next-step guidance could lift usefulness above the desired threshold without sacrificing empathy or correctness.

**Qualitative Example.** To illustrate the nature of learner reflections and the feedback produced by the system, we present an anonymized example below.

### Student Reflection

“i have learn what are large language models,large language models(LLMs)are very large deep learning models that are pre-trained on vast amounts of data,this can ingest massive amounts of data,often from the Internet,but also from sources such as the common crawl.i also learn such as temperature(controls randomness) , top p(controls diversity based on cumulative probability),max length,stop sequences,frequency penalty,presence penalty ,model behavior may vary,basics of prompting(key concepts),prompt structure in chat models,basic prompt example... if i want to applying this in a real-world scenario,in content creation,i can use designed prompts with controlled temperature,for a factual

Rubric Dimension	Human–Human ICC	Human–AI ICC
Concept Understanding	$0.56 \pm 0.41$	$0.31 \pm 0.16$
Real-World Application	$0.62 \pm 0.33$	$0.49 \pm 0.08$
Reflection Questions	$0.68 \pm 0.29$	$0.49 \pm 0.09$
Clarity of Communication	$0.44 \pm 0.50$	$0.36 \pm 0.13$
<b>Overall</b>	<b><math>0.57 \pm 0.35</math></b>	<b><math>0.41 \pm 0.13</math></b>

Table 6: Two-way mixed intraclass correlation  $ICC^{(2,1)}$ .

Rubric Dimension	Low MAE	High MAE	$\Delta_{MAE}$
Concept Understanding	1.000	0.278	0.722
Real-World Application	1.500	0.403	1.097
Reflection Questions	0.917	0.431	0.486
Clarity of Communication	0.167	0.472	0.306
<b>Overall</b>	<b>0.896</b>	<b>0.396</b>	<b>0.500</b>

Table 7: Model MAE stratified by learner ability, where ability is defined *once per reflection* as the average human rubric score across all four dimensions ( $n_{low} = 48$ ,  $n_{high} = 288$ ). Lower  $\Delta_{MAE}$  indicates greater fairness.

blog ,set a loww temperature for precise information my open question is how to precisely calibrate prompt parameters for complex tasks.as small adjust might drastically change outputs.”

### AI-Generated Feedback

“You’ve made a solid start in understanding large language models (LLMs) and their applications in content creation, which is commendable! To enhance your reflection, consider diving deeper into the technical terms you’ve mentioned, explaining their significance and how they interconnect. This will provide a clearer picture of LLMs. When discussing real-world applications, try to include more specific examples that showcase diverse uses across various cultures or industries, as this will enrich your perspective and representation. Your insightful question about calibrating prompt parameters indicates critical thinking; expanding on this could lead to a deeper exploration of the challenges faced in practical applications. Additionally, focusing on clarity and structure in your writing will significantly improve communication—proofreading for grammar and punctuation can make your ideas shine. Reflect on how your understanding of LLMs might evolve by considering their diverse applications and think about how to make your explanations more accessible to a broader audience. Keep up the great work, and I look forward to your continued progress!”

## RQ3: Practicality (Time & Cost)

**Latency and throughput.** A scoring-only pass, which involves a single invocation of the EVALUATOR agent, requires an average of  $7.71 \pm 0.41$  s per reflection. Gener-

Dimension	Mean $\pm$ SD
Correctness	4.080 $\pm$ 0.756
Alignment	3.924 $\pm$ 0.763
Actionability	3.760 $\pm$ 0.845
Depth of Insight	3.845 $\pm$ 0.860
Empathy & Tone	<b>4.223 <math>\pm</math> 0.612</b>
<b>Overall <math>Q(g)</math></b>	<b>3.967</b>

Table 8: Feedback usefulness ratings (higher is better).

Annotator	Class 1	Class 6	Class 12	Total
#1	30	30	20	80
#2	80	60	75	215
#3	27	13	11	51

Table 9: Time (minutes) each human annotator spent assigning rubric-based scores to reflections

ating narrative feedback through the EQUITY MONITOR, METACOGNITIVE, and AGGREGATOR agents, followed by a final REFLEXION check, adds approximately 25.64 s. This results in an end-to-end latency of **33.35 s** per reflection.

For the evaluation corpus of  $N = 84$  reflections (28 students  $\times$  3 reflections each), this translates to **10.8 min** wall-clock time for rubric scoring alone and **46.7 min** for complete feedback generation on a single processing thread. Because both stages are dominated by network I/O rather than compute, throughput scales nearly inversely with the number of parallel threads  $T$ , approaching time/ $T$  under ideal conditions.

For context, the three human evaluators self-reported that they required an average of **1.4 min** to score a single reflection, based solely on rubric-based scoring. The slowest evaluator required **2.6 min**. In contrast, our automated pipeline completes the same scoring task in 7.71s. This corresponds to an approximately  $11\times$  average speed-up and a  $20\times$  improvement relative to the slowest evaluator, while maintaining the level of accuracy and fairness reported in Table 9.

Table 9 reports the wall-clock time each evaluator spent assigning rubric scores. Narrative feedback was not included in this task. Across all 36 evaluations this corresponds to an overall average of  $\approx 9.6$  min per individual evaluation. Parallel self-consistency (three runs, four workers) keeps total latency below 12 min, confirming suitability for real-time classroom feedback.

**Monetary cost.** A single pass of the grading and feedback pipeline using gpt-4o-mini-2024-07-18 consumes 1216 input tokens and 2283 output tokens, for a total of 3499 tokens. At the model’s 18 July 2024 pricing (\$0.15 per million input tokens and \$0.60 per million output tokens), the cost per reflection is  $\$1.55 \times 10^{-3}$ .<sup>1</sup> Grading and providing feedback for the 28-student cohort

<sup>1</sup>  $1216 \times 0.15/10^6 + 2283 \times 0.60/10^6 = 0.0015522$ .

across three sessions ( $28 \times 3 = 84$  reflections) therefore costs \$0.13. Applying the same process to all twelve weekly reflections ( $28 \times 12$ ) yields a total cost of \$0.52.

These figures highlight the practical viability of the approach. The pipeline is time-efficient, low-cost, and suitable for use at scale, satisfying the practicality requirement in RQ3 and making automated feedback accessible even in resource-constrained educational settings.

## Conclusion

This work addresses a central equity challenge in education: providing timely, high-quality formative feedback at a scale that is unattainable through human effort alone. We introduce (i) a scalable and self-consistent scoring pipeline that produces auditable rubric scores with minimal human supervision and (ii) a role-based agent ensemble that generates bias-aware, conversational feedback aligned with formative-assessment and metacognitive principles.

Our evaluation in an authentic AI literacy course shows that the system approaches expert-level scoring fidelity, produces feedback that trained graders judge as helpful and empathetic, and maintains reasonable fairness across learner proficiency levels. The pipeline operates with sub-minute latency and negligible cost, making it feasible for deployment in large or resource-constrained courses.

Although the system does not eliminate all disparities in automated grading, it provides a transparent and measurable framework for monitoring accuracy and fairness, and it offers a practical foundation for future calibration and refinement. By integrating structured scoring, equity-sensitive evaluation, and pedagogically grounded feedback generation, this work contributes a meaningful step toward scalable, accessible, and socially responsible AI support for learning.

The study has limitations. Evaluations are restricted to adult, English-language reflections; low-proficiency writers continue to experience higher error; and model drift may affect performance over time. Future research should extend the framework to multilingual and multicultural contexts, incorporate adaptive fairness objectives that reduce subgroup gaps, and measure causal learning gains in large-scale classroom deployments.

Education, often regarded as one of society’s most powerful engines of social progress, is entering a period of rapid transformation. We hope this work will encourage further research and collaboration toward building learning ecosystems that are inclusive, equitable, and enriched by responsible AI.

## Ethical Statement

**Intended purpose.** The proposed pipeline is conceived as an *augmentation* for instructors rather than a replacement. By automating low-level rubric application and producing concise, dialogic comments, it seeks to make timely formative feedback feasible even in large or resource-constrained courses. The overarching aim is to foster more equitable learning outcomes through broader access to high-quality feedback.

**Potential misuse and risks.** Three principal hazards are anticipated. First, instructors and students might place excessive trust in automated scores that can still exhibit residual bias. Second, if deployed without human oversight, the system could erode teachers' formative role. Third, storing raw reflections indefinitely raises privacy concerns. In addition, applying the model beyond its validated context, for example with younger learners or non-English submissions, could yield pedagogically inappropriate guidance.

**Mitigation strategies.** A human-in-the-loop workflow ensures that instructors review all model outputs before release and retain full authority over grades. Fairness is monitored through an objective that jointly minimizes mean absolute error and worst-subgroup error, accompanied by demographic disparity reports. Transparency is promoted by attaching confidence tags (CONFIDENT or REVISE) and a concise rationale to every comment. Privacy is protected by anonymizing reflections at ingestion and processing them on secure servers under an approved IRB protocol.

**Broader societal impact.** By lowering the cost of equitable, dialogic feedback, the system has the potential to narrow attainment gaps in community colleges, vocational programs, and massive open online courses. Openly releasing resources is intended to catalyze transparent, bias-aware progress in AI-mediated formative assessment while reinforcing the indispensable pedagogical role of human educators.

## References

- Black, P.; and Wiliam, D. 1998. *Inside the black box: Raising standards through classroom assessment*. Granada Learning.
- Black, P.; and Wiliam, D. 2009. Developing the theory of formative assessment. *Educational Assessment, Evaluation and Accountability (formerly: Journal of personnel evaluation in education)*, 21(1): 5–31.
- Brown, G. T. 2022. The past, present and future of educational assessment: A transdisciplinary perspective. In *Frontiers in Education*, volume 7, 1060633. Frontiers Media SA.
- Carlson, M.; Pack, A.; and Escalante, J. 2023. Utilizing OpenAI's GPT-4 for written feedback. *Tesol Journal*, 759: e759.
- Chen, D.; Hebert, M.; and Wilson, J. 2022. Examining human and automated ratings of elementary students' writing quality: A multivariate generalizability theory application. *American Educational Research Journal*, 59(6): 1122–1156.
- Chiang, C.-H.; Chen, W.-C.; Kuan, C.-Y.; Yang, C.; and Lee, H.-y. 2024. Large language model as an assignment evaluator: Insights, feedback, and challenges in a 1000+ student course. *arXiv preprint arXiv:2407.05216*.
- García-Varela, F.; Nussbaum, M.; Mendoza, M.; Martínez-Troncoso, C.; and Bekerman, Z. 2025. ChatGPT as a Stable and Fair Tool for Automated Essay Scoring. *Education Sciences*, 15(8): 946.
- Graham, S.; Harris, K.; and Hebert, M. 2011. Informing Writing: The Benefits of Formative Assessment. A Report from Carnegie Corporation of New York. *Carnegie Corporation of New York*.
- Graham, S.; Hebert, M.; and Harris, K. R. 2015. Formative assessment and writing: A meta-analysis. *The elementary school journal*, 115(4): 523–547.
- Graham, S.; McKeown, D.; Kiuahara, S.; and Harris, K. R. 2012. A meta-analysis of writing instruction for students in the elementary grades. *Journal of educational psychology*, 104(4): 879.
- Grimes, D.; and Warschauer, M. 2010. Utility in a fallible tool: A multi-site case study of automated writing evaluation. *The Journal of Technology, Learning and Assessment*, 8(6).
- Hattie, J.; and Timperley, H. 2007. The power of feedback. *Review of educational research*, 77(1): 81–112.
- Heilman, M.; and Madnani, N. 2013. ETS: Domain adaptation and stacking for short answer scoring. In *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, 275–279.
- Hollingsworth, J. 1960. Automatic graders for programming classes. *Communications of the ACM*, 3(10): 528–529.
- Impey, C.; Wenger, M.; Garuda, N.; Golchin, S.; and Stamer, S. 2025. Using large language models for automated grading of student writing about science. *International Journal of Artificial Intelligence in Education*, 1–35.
- Juwah, C.; Macfarlane-Dick, D.; Matthew, B.; Nicol, D.; Ross, D.; and Smith, B. 2004. Enhancing student learning through effective formative feedback. *The Higher Education Academy*, 140: 1–40.
- Kortemeyer, G. 2023. Can an AI-tool grade assignments in an introductory physics course? *arXiv preprint arXiv:2304.11221*.
- Liu, J.; and Dorans, N. J. 2016. Fairness in score interpretation. *Fairness in educational assessment and measurement*, 77–96.
- Liu, X.; Wang, S.; Wang, P.; and Wu, D. 2019. Automatic grading of programming assignments: an approach based on formal semantics. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)*, 126–137. IEEE.
- MacArthur, C. A. 2016. Instruction in evaluation and revision.
- Messer, M.; Brown, N. C.; Kölling, M.; and Shi, M. 2024. Automated grading and feedback tools for programming education: A systematic review. *ACM Transactions on Computing Education*, 24(1): 1–43.
- Moore, N. S.; and MacArthur, C. A. 2016. Student use of automated essay evaluation technology during revision. *Journal of Writing Research*, 8(1): 149–175.
- Morris, R.; Perry, T.; and Wardle, L. 2021. Formative assessment and feedback for learning in higher education: A systematic review. *Review of Education*, 9(3): e3292.
- Nicolai, K. D.; Koenka, A. C.; and Braxton, D. 2023. A literature review of Black and Latiné youth's experience of inequitable feedback: How can middle school educators provide motivationally-supportive feedback? *Middle School Journal*, 54(4): 14–24.

Ouyang, F.; Dinh, T. A.; and Xu, W. 2023. A systematic review of AI-driven educational assessment in STEM education. *Journal for STEM Education Research*, 6(3): 408–426.

Page, E. B. 1966. The imminence of... grading essays by computer. *The Phi Delta Kappan*, 47(5): 238–243.

Panadero, E.; Jonsson, A.; Pinedo, L.; and Fernández-Castilla, B. 2023. Effects of rubrics on academic performance, self-regulated learning, and self-efficacy: A meta-analytic review. *Educational Psychology Review*, 35(4): 113.

Pishchukhina, O.; and Allen, A. 2021. Supporting learning in large classes: online formative assessment and automated feedback. In *2021 30th Annual Conference of the European Association for Education in Electrical and Information Engineering (EAEEIE)*, 1–4. IEEE.

Riordan, B.; Horbach, A.; Cahill, A.; Zesch, T.; and Lee, C. 2017. Investigating neural architectures for short answer scoring. In *Proceedings of the 12th workshop on innovative use of NLP for building educational applications*, 159–168.

Roscoe, R. D.; Wilson, J.; Johnson, A. C.; and Mayra, C. R. 2017. Presentation, expectations, and experience: Sources of student perceptions of automated writing evaluation. *Computers in Human Behavior*, 70: 207–221.

Sung, C.; Dhamecha, T.; Saha, S.; Ma, T.; Reddy, V.; and Arora, R. 2019. Pre-training BERT on domain resources for short answer grading. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 6071–6075.

Underwood, J. S.; and Tregidgo, A. P. 2006. Improving student writing through effective feedback: Best practices and recommendations. *Journal of Teaching Writing*, 22(2): 73–98.

Ureel II, L. C.; and Wallace, C. 2019. Automated critique of early programming antipatterns. In *Proceedings of the 50th ACM Technical Symposium on Computer Science Education*, 738–744.

Wu, Q.; Bansal, G.; Zhang, J.; Wu, Y.; Li, B.; Zhu, E.; Jiang, L.; Zhang, X.; Zhang, S.; Liu, J.; et al. 2024. Autogen: Enabling next-gen LLM applications via multi-agent conversations. In *First Conference on Language Modeling*.

Xiao, C.; Ma, W.; Song, Q.; Xu, S. X.; Zhang, K.; Wang, Y.; and Fu, Q. 2025. Human-ai collaborative essay scoring: A dual-process framework with llms. In *Proceedings of the 15th International Learning Analytics and Knowledge Conference*, 293–305.

Yeung, C.; Yu, J.; Cheung, K. C.; Wong, T. W.; Chan, C. M.; Wong, K. C.; and Fujii, K. 2025. A Zero-Shot LLM Framework for Automatic Assignment Grading in Higher Education. *arXiv preprint arXiv:2501.14305*.