

Decision-Level Ordinal Modeling for Multimodal Essay Scoring with Large Language Models

Anonymous ACL submission

Abstract

Automated essay scoring (AES) predicts multiple rubric-defined trait scores for each essay, where each trait follows an ordered discrete rating scale. Most LLM-based AES methods cast scoring as autoregressive token generation and obtain the final score via decoding and parsing, making the decision implicit. This formulation is particularly sensitive in multimodal AES, where the usefulness of visual inputs varies across essays and traits. To address these limitations, we propose **Decision-Level Ordinal Modeling (DLOM)**, which makes scoring an explicit ordinal decision by reusing the LM head to extract score-wise logits on pre-defined score tokens, enabling direct optimization and analysis in the score space. For multimodal AES, **DLOM-GF** introduces a gated fusion module that adaptively combines textual and multimodal score logits. For text-only AES, **DLOM-DA** adds a distance-aware regularization term to better reflect ordinal distances. Experiments on the multimodal Essay-Judge dataset show that DLOM improves over a generation-based SFT baseline across scoring traits, and DLOM-GF yields further gains when modality relevance is heterogeneous. On the text-only ASAP/ASAP++ benchmarks, DLOM remains effective without visual inputs, and DLOM-DA further improves performance and outperforms strong representative baselines.

1 Introduction

Automated Essay Scoring (AES) is a longstanding problem in natural language processing (Williamson et al., 1999; Smolentzov, 2013), aiming to automatically assess the quality of student essays according to predefined scoring rubrics (Misgna et al., 2024), often across multiple scoring traits rather than a single holistic score (Ridley et al., 2021; Kumar et al., 2022; Do et al., 2024a). Recent advances in large language models (LLMs) have significantly improved semantic

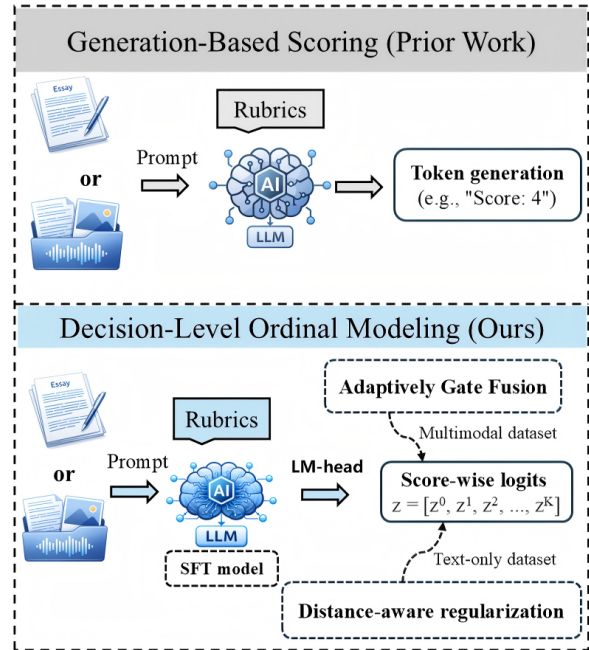


Figure 1: Comparison between generation-based scoring and decision-level ordinal modeling (DLOM).

understanding and instruction following, making them a natural choice for AES. As a result, a growing body of work has explored LLM-based essay scoring (Ridley et al., 2021; Ramesh and Sanampudi, 2022; Xia et al., 2024; Chu et al., 2025), including extensions to multimodal settings where essays are accompanied by visual inputs such as charts or diagrams, which further complicates the scoring decision (Su et al., 2025a,b).

Most existing LLM-based AES approaches formulate scoring as a generation problem (Do et al., 2024a,b), where the model is prompted to produce a score token or short textual response corresponding to a rubric-defined grade. This formulation is appealing due to its simplicity and compatibility with instruction-tuned LLMs (Pack et al., 2024; Song et al., 2024; Bui and Barrot, 2025). However, generation-based scoring makes the scoring

061	decision implicit, as the final score emerges as a	and text-only benchmarks. We demonstrate	111
062	byproduct of language generation (Liu et al., 2023).	the effectiveness and generality of the pro-	112
063	Consequently, the decision process is difficult to	posed formulation on both multimodal and	113
064	control, analyze, or adapt, and is often sensitive to	text-only essay scoring benchmarks, achiev-	114
065	prompt design, decoding strategies, and tokeniza-	ing consistent improvements over generation-	115
066	tion artifacts—factors that are largely orthogonal	based SFT baselines.	116
067	to the scoring task itself (Wang et al., 2018).		
068	In particular, AES is inherently ordinal: score	2 Related Work	117
069	categories follow a fixed and ordered structure, and		
070	errors between adjacent score levels are typically	2.1 Automated Essay Scoring	118
071	less severe than those between distant levels (John-	Automated essay scoring (AES) has progressed	119
072	son, 1996; Taghipour and Ng, 2016; Alikaniotis	from feature-engineered systems based on surface,	120
073	et al., 2016). Motivated by this, we adopt a	syntactic, and discourse cues (Page, 1966; Dikli,	121
074	decision-based formulation of essay scoring to bet-	2006; Dong and Zhang, 2016; Hou et al., 2025) to	122
075	ter match the ordinal nature of rubric-based assess-	neural representation learning with CNN/RNN ar-	123
076	ment. We further argue that, especially with LLMs	chitectures and pre-trained language models (Dong	124
077	that provide strong semantic representations, separ-	et al., 2017; Shibata and Uto, 2025). Recent	125
078	ating semantic understanding from decision mak-	transformer-based models further improve scoring	126
079	ing offers a more appropriate and general modeling	by leveraging contextual representations and trans-	127
080	perspective for essay scoring.	fer learning (Devlin et al., 2019; Yang et al., 2020;	128
081	In this work, we propose Decision-Level Or-	Azhari et al., 2024), enabling end-to-end learning	129
082	dinal Modeling (DLOM) for LLM-based essay	of essay representations and scoring functions.	130
083	scoring. Figure 1 illustrates the difference be-		
084	tween generation-based scoring and our proposed	2.2 LLM-Based Essay Scoring	131
085	decision-level ordinal modeling (DLOM). Build-	With the rise of large language models (LLMs),	132
086	ing on a supervised fine-tuned (SFT) LLM, we	many recent AES methods use instruction-	133
087	cast score prediction as an explicit decision over a	following LLMs as backbones (Dong et al., 2021;	134
088	predefined ordered label set. The model extracts	Li and Liu, 2024). Most work adopts a generation-	135
089	score-wise logits over the ordinal categories, en-	based interface, prompting the model to output a	136
090	abling direct optimization and analysis at the de-	score token or short response aligned with rubric-	137
091	cision stage while preserving the semantic repre-	defined grades (Zhou et al., 2023; Li and Pan,	138
092	sentations learned during SFT. We further extend	2025). This approach is simple and often requires	139
093	DLOM to multimodal scoring with a decision-level	no architectural changes, but the final label is typi-	140
094	gated fusion module (DLOM-GF) that adaptively	cally obtained by decoding and parsing generat-	141
095	combines text-only and multimodal score logits.	ed text, making it sensitive to prompting, decoding,	142
096	For text-only benchmarks, we introduce a distance-	and tokenization details. Several studies improve	143
097	aware regularization variant (DLOM-DA) to better	robustness via prompt design, calibration, or decod-	144
098	respect ordinal distances between score levels.	ing constraints (Wiher et al., 2022; Liu et al., 2024),	145
099	Our main contributions are threefold:	while largely retaining the generation paradigm.	146
100		2.3 Ordinal Modeling for Scoring Tasks	147
101	• Decision-level ordinal modeling for LLM-	Ordinal modeling has been recognized as impor-	148
102	based AES. We cast scoring as an explicit or-	tant for scoring and assessment tasks, where score	149
103	dinal decision by predicting score-wise logits	categories follow a fixed order and misclassifica-	150
104	over predefined score levels, enabling direct	tion costs are asymmetric (Frank and Hall, 2001;	151
105	optimization and analysis in the score space.	Gutiérrez et al., 2016). Prior work outside the LLM	152
106		context has explored ordinal regression, ranking-	153
107	• DLOM extensions for multimodal and text-	based losses, and distance-aware objectives to bet-	154
108	only scoring. We introduce DLOM-GF for	ter reflect the structure of scoring scales (McCul-	155
109	decision-level multimodal gated fusion and	lagh, 1980; Herbrich et al., 2000; Suárez et al.,	156
110	DLOM-DA for distance-aware regularization	2021). In LLM-based AES, however, ordinal struc-	157
	in text-only settings.	ture is often only implicitly encoded through tex-	158
	• Comprehensive evaluation on multimodal		

Paradigm	Decision interface	Multimodal fusion
Generation-based LLM AES	Decode and parse score text	Prompt/decoding dependent
Pooling-head classifier	Pool(H) + ($K+1$)-way classification head	Usually representation-level
DLOM (ours)	LM-head score-logit extraction on fixed score tokens	Decision-level gated fusion

Table 1: Positioning of our approach against common decision interfaces used in AES.

159 tual descriptions of score levels or prompt instruc- 200
160 tions (Stahl et al., 2024). Such an explicit decision 201
161 interface over ordered labels has received relatively 202
162 limited attention in LLM-based AES, motivating 203
163 our decision-level ordinal formulation.

164 **Positioning.** Table 1 summarizes common deci- 204
165 sion interfaces in LLM-based AES. Compared 205
166 to generation-and-parse scoring and pooling-head 206
167 classifiers, DLOM reuses the LM head to extract 207
168 score-wise logits on fixed score tokens and makes 208
169 predictions directly in the ordered score space, 209
170 avoiding both decoding-time artifacts and an addi- 210
171 tional classification head. 211

172 2.4 Multimodal Essay Scoring

173 Beyond text-only AES, recent work has explored 214
174 multimodal essay scoring scenarios in which essays 215
175 are accompanied by visual inputs such as charts, 216
176 tables, or diagrams (Su et al., 2025a). Most ap- 217
177 proaches perform multimodal fusion at the repre- 218
178 sentation level, combining textual and visual infor- 219
179 mation via concatenation, attention mechanisms, or 220
180 cross-modal transformers (Su et al., 2025b). Since 221
181 modality relevance can vary across essays and 222
182 traits, representation-level fusion may be sensitive 223
183 to heterogeneous modality contributions. Decision- 224
184 level fusion that operates on score predictions has 225
185 been less explored in AES, and offers an alterna- 226
186 tive interface for controlling modality influence and 227
187 improving interpretability. 228

188 3 Methodology

189 3.1 Problem Formulation

190 Given an essay e and an optional visual input v , 234
191 the goal of automated essay scoring is to predict a 235
192 discrete score y from a fixed, ordered set of score 236
193 levels $\mathcal{Y} = \{0, 1, \dots, K\}$. The score levels follow 237
194 an inherent ordinal structure, where misclassifica- 238
195 tions between adjacent levels are less severe than 239
196 those between distant levels. We assume access to a 240
197 supervised fine-tuned (SFT) large language model 241
198 that encodes the semantic content of the essay. Our 242
199 objective is to design a decision mechanism that

operates on top of this model and explicitly ac- 200
counts for the ordinal nature of the scoring task at 201
the decision stage. 202

203 3.2 Decision-Level Ordinal Modeling

204 Instead of formulating essay scoring as a language 205
206 generation task, we cast score prediction as an ex- 206
207 plicit decision over a fixed set of ordered score 207
208 levels. A standard AES classifier typically adds a 208
209 ($K+1$)-way linear head on top of pooled hidden 209
210 states (e.g., [CLS]) and predicts scores in that head 210
211 space. In contrast, our decision-level formulation 211
212 reuses the LLM’s LM head (vocabulary projec- 212
213 tion) and derives score-wise logits directly from 213
214 the model outputs without decoding. Concretely, 214
215 we restrict the final-step vocabulary logits to a pre- 215
216 defined set of score tokens, obtaining $\mathbf{z} \in \mathbb{R}^{K+1}$ 216
217 as the decision space. This keeps the backbone 217
218 unchanged while making the scoring decision ex- 218
219 plicit in an ordered score space, and also supports 219
220 decision-level fusion for multimodal scoring. 219

220 Figure 2 provides an overview of the pro- 220
221 posed framework. Given an input essay (and op- 221
222 tional visual input), the SFT LLM produces final- 222
223 step vocabulary logits via its LM head. We se- 223
224 lect the entries corresponding to the predefined 224
225 score tokens to form the score-wise logit vector 225
226 $\mathbf{z} = [z_0, z_1, \dots, z_K] \in \mathbb{R}^{K+1}$. Formally, $\mathbf{z} =$ 226
227 $\text{Logits}(e, v)_L[\mathcal{S}]$, where \mathcal{S} denotes the score-token 227
228 set and L is the final position. This formulation 228
229 decouples semantic understanding from decision 229
230 making: the LLM encodes essay semantics, while 230
231 the scoring decision is made explicitly in the ordi- 231
232 nal score space. At inference time, the predicted 232
233 score is obtained by $\hat{y} = \arg \max_{k \in \mathcal{Y}} z_k$. A de- 233
234 tailed comparison between generation-based scor- 234
235 ing, CLS-head classification, and our decision-level 235
236 interface is provided in Appendix A. 236

237 3.3 DLOM-GF for Multimodal AES

238 For multimodal essay scoring settings where both 238
239 textual and visual inputs are available, the reliabil- 239
240 ity of visual information varies across essays and 240
241 scoring traits. We introduce a decision-level gat- 241
242 ing mechanism that adaptively combines textual 242

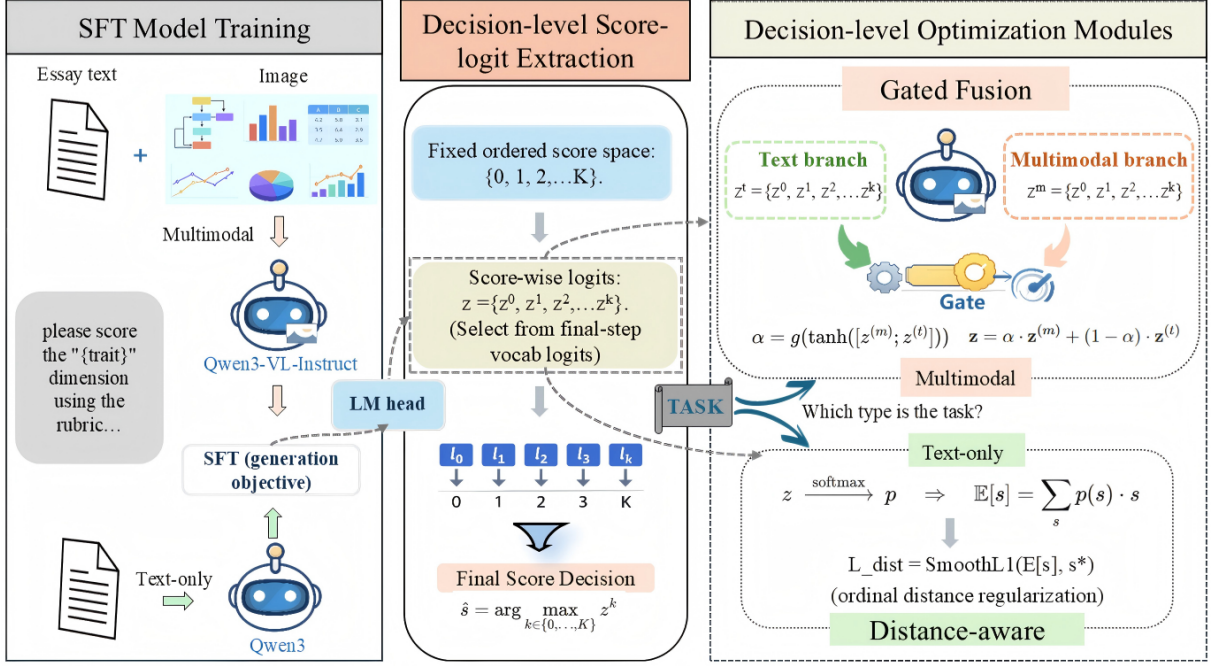


Figure 2: Overview of the proposed decision-level ordinal modeling framework. The framework consists of three stages: (i) supervised fine-tuning (SFT) for semantic encoding, (ii) decision-level score-logit extraction over an ordered score-token set, and (iii) task-specific decision-level objectives: decision-level gated fusion for multimodal scoring and distance-aware regularization for text-only scoring.

and multimodal predictions based on their score-level evidence. Specifically, we obtain two score-wise logit vectors: $\mathbf{z}^{(t)} \in \mathbb{R}^{K+1}$ from the text-only branch and $\mathbf{z}^{(m)} \in \mathbb{R}^{K+1}$ from the multimodal branch, where each logit corresponds to a candidate score level. We concatenate the two logit vectors and apply a non-linear transformation to produce a gating input:

$$\mathbf{h} = \tanh([\mathbf{z}^{(m)}; \mathbf{z}^{(t)}]).$$

A lightweight gating network $g(\cdot)$ then maps \mathbf{h} to a scalar fusion weight:

$$\alpha = g(\mathbf{h}) \in [0, 1],$$

where $g(\cdot)$ is implemented as a linear layer followed by a sigmoid activation. The final score logits are computed as a convex combination in the ordinal score space:

$$\mathbf{z} = \alpha \cdot \mathbf{z}^{(m)} + (1 - \alpha) \cdot \mathbf{z}^{(t)}.$$

By conditioning the gate on score-wise logits rather than intermediate representations, the proposed mechanism operates entirely at the decision level. The gating function is learned implicitly through the scoring objective, allowing the model to estimate modality reliability in an instance-adaptive and task-aligned manner.

3.4 Training Objective

The model is trained to predict an ordinal score label from a fixed and ordered score set. Given the score-wise logits $\mathbf{z} \in \mathbb{R}^{K+1}$ produced by the decision-level score-logit extraction (or the gated fusion module in multimodal settings), we compute a categorical distribution $p = \text{softmax}(\mathbf{z})$ and optimize the cross-entropy loss with respect to the gold score y :

$$\mathcal{L}_{\text{CE}} = -\log p_y.$$

For multimodal essay scoring, the score logits \mathbf{z} are obtained via the proposed decision-level gated fusion module, and the entire model is trained end-to-end by backpropagating \mathcal{L}_{CE} through both modality branches and the gating function.

In text-only settings, where decision-level fusion is not applicable, we additionally introduce a distance-aware regularization term to better capture the ordinal structure of the scoring scale. Specifically, we compute the expected score under the predicted distribution, $\mathbb{E}[s] = \sum_{k \in \mathcal{Y}} p_k \cdot k$, and penalize its deviation from the gold score using a SmoothL1 loss:

$$\mathcal{L}_{\text{dist}} = \text{SmoothL1}(\mathbb{E}[s], y).$$

291 The final objective is

$$292 \quad \mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda \mathcal{L}_{\text{dist}} + \beta(\lambda - 0.5)^2,$$

293 where $\lambda = \sigma(\lambda_{\text{logit}}) \in (0, 1)$ is a learnable weight
294 and β is a small regularization coefficient. This
295 auxiliary objective encourages ordinal consistency
296 by assigning larger penalties to predictions that are
297 farther away from the ground-truth score.

298 4 Experiments

299 4.1 Datasets and Evaluation Metrics

300 We evaluate our approach on both multimodal and
301 text-only multi-trait essay scoring benchmarks and
302 treat each trait as an independent ordinal prediction
303 task and report trait-wise results as well as their
304 macro-averaged performance.

305 **Multimodal Dataset.** For multimodal essay scor-
306 ing, we use the EssayJudge dataset (Su et al.,
307 2025a), which consists of student essays accompa-
308 nied by visual inputs such as charts or diagrams.
309 Each essay is annotated with multiple trait-level
310 scores according to a predefined rubric, covering
311 diverse aspects of writing quality. This dataset is
312 specifically designed to assess multimodal reason-
313 ing in essay scoring, where visual information may
314 contribute unevenly across essays and traits (more
315 information is in the Appendix B.1).

316 **Text-only Datasets.** For text-only evaluation, we
317 use the combined ASAP/ASAP++ benchmark,
318 which is a standard dataset for multi-trait auto-
319 mated essay scoring. The original ASAP dataset
320 provides holistic scores for all prompts, while only
321 a subset of prompts includes trait-level annota-
322 tions. ASAP++ extends ASAP by providing trait
323 scores for the remaining prompts (Mathias and
324 Bhattacharyya, 2018, 2020), enabling compre-
325 hensive multi-trait evaluation across all prompts. The
326 final dataset consists of eight writing prompts, each
327 associated with its own scoring rubric and traits
328 to be evaluated (more information is in the Ap-
329 pendix B.2). Following prior work, we use the
330 combined ASAP/ASAP++ dataset for trait-level
331 and cross-prompt evaluation.

332 **Evaluation Metrics.** We adopt Quadratic
333 Weighted Kappa (QWK) as the primary evaluation
334 metric (Cohen, 1968), which is widely used in AES
335 to measure agreement between model predictions
336 and human raters while accounting for the ordinal
337 nature of score scales. More information is shown
338 in Appendix B.3.

4.2 Experimental Setup 339

340 **Backbone Models.** We use Qwen3-VL-8B-
341 Instruct (Qwen Team, 2025b) as the backbone
342 language model for multimodal experiments and
343 Qwen3-8B (Qwen Team, 2025a) for text-only ex-
344 periments. Both models are initialized from pub-
345 licly released checkpoints and further adapted to
346 the essay scoring task through supervised fine-
347 tuning. For further details on the model check-
348 points, please refer to Appendix B.4.

349 **Training Protocol.** For both multimodal and text-
350 only settings, we first perform generation-based
351 supervised fine-tuning using a unified prompt tem-
352 plate based on the official rubric descriptions. We
353 then apply the proposed DL0M training on top of
354 the SFT models, details of the prompt template and
355 rubrics are provided in Appendix B.5.

356 In multimodal experiments, the text-only and
357 multimodal branches are trained jointly, with score
358 predictions from both branches combined through
359 the decision-level gated fusion mechanism. The
360 entire model, including the gating function, is opti-
361 mized end-to-end under the same ordinal classifica-
362 tion objective. All experiments are conducted using
363 five-fold cross-validation, and results are reported
364 by averaging performance across folds. For the text-
365 only datasets (ASAP/ASAP++), we follow the pre-
366 viously established partitioning scheme (Taghipour
367 and Ng, 2016). For the multimodal dataset, which
368 lacks an existing standard split, we construct five
369 folds at the essay level, ensuring that no essay ap-
370 pears in more than one fold. We do not apply early
371 stopping during training. Instead, all models are
372 trained for a fixed number of epochs under the same
373 training configuration.

374 **Inference.** At inference time, for both text-only
375 and multimodal settings, scores are predicted di-
376 rectly from score-wise logits, and each trait is eval-
377 uated independently.

4.3 Main Results on EssayJudge 378

379 To enable supervised training and reproducible
380 comparison on EssayJudge, we establish an ex-
381 plicit data partition and report training-based re-
382 sults. Prior work (CAFES) evaluates EssayJudge
383 mainly through inference-time agentic pipelines
384 over off-the-shelf MLLMs (Su et al., 2025b). We
385 therefore report their best-reported results in our
386 table as a reference. Table 2 shows the main results
387 on the multimodal EssayJudge dataset under the

Model	Lexical Level		Sentence Level				Discourse Level				Avg
	LA	LD	CH	GA	GD	PA	AC	JP	OS	EL	
CAFES	0.510	0.500	0.520	0.570	0.540	0.490	0.370	0.440	0.480	0.280	0.470
SFT-Gen (baseline)	0.595	0.547	0.570	0.569	0.591	0.503	0.205	0.404	0.486	0.455	0.492
DLOM	0.594	0.544	0.562	0.580	0.587	0.498	0.251	0.425	0.517	0.477	0.504
DLOM-GF	0.624	0.551	0.569	0.589	0.613	0.506	0.251	0.447	0.527	0.479	0.516

Table 2: Main results on the EssayJudge dataset under the multi-trait scoring setting. LA: lexical accuracy, LD: lexical diversity, CH: coherence, GA: grammatical accuracy, GD: grammatical diversity, PA: punctuation accuracy, AC: argument clarity, JP: justifying persuasiveness, OS: organizational structure, EL: essay length.

Model	Traits										Avg
	Content	PA	Lang	Nar	Org	Conv	WC	SP	Style	Voice	
HISK	0.679	0.697	0.605	0.659	0.610	0.527	0.579	0.553	0.609	0.489	0.601
MTL-BiLSTM	0.685	0.701	0.604	0.668	0.615	0.560	0.615	0.598	0.632	0.582	0.626
ArTS	0.730	0.751	0.698	0.725	0.672	0.668	0.679	0.678	0.721	0.570	0.689
SFT-Gen (baseline)	0.705	0.752	0.698	0.732	0.638	0.582	0.632	0.614	0.521	0.585	0.646
DLOM	0.738	0.776	0.716	0.753	0.692	0.655	0.659	0.662	0.615	0.581	0.685
DLOM-DA	0.745	0.778	0.720	0.756	0.699	0.659	0.673	0.669	0.644	0.634	0.697

Table 3: Trait-level results on the ASAP/ASAP++ benchmark. PA: prompt adherence, Lang: language, Nar: narrativity, Org: organization, Conv: conventions, WC: word choice, SP: sentence fluency.

multi-trait scoring setting. We compare three scoring paradigms: generation-based supervised fine-tuning (SFT-Gen, baseline), decision-level ordinal modeling (DLOM), and decision-level modeling with gated multimodal fusion (DLOM-GF). Compared to SFT-Gen, decision-level ordinal modeling results in a higher average QWK, while exhibiting differentiated effects across scoring traits. In particular, the decision-level formulation yields clear gains on discourse-level dimensions such as argument clarity and organizational structure, which require holistic reasoning over essay content. In contrast, improvements on lexical-level and sentence-level traits are comparatively limited, suggesting that the benefits of explicit ordinal decision modeling are more pronounced for higher-level scoring criteria.

Building upon the decision-level formulation, incorporating a gated multimodal fusion mechanism further improves overall performance and achieves the best average QWK among all compared methods. Beyond improving individual lower-level traits, the gating mechanism consistently strengthens performance on discourse-level dimensions, further stabilizing gains on traits that rely on holistic essay evaluation. These results suggest that adaptively integrating textual and visual evidence at the decision stage not only alleviates weaknesses in specific traits, but also reinforces the advantages of decision-level modeling for high-level essay scor-

ing. You can find more information about main results in Appendix B.6

4.4 Text-only Results on ASAP/ASAP++

To examine whether the proposed decision-level formulation generalizes beyond the multimodal setting, we further evaluate our approach on the text-only ASAP/ASAP++ benchmark. For completeness, we also include several representative baselines (HISK (Cozma et al., 2018), MTL-BiLSTM (Kumar et al., 2022), and ArTS (Do et al., 2024a)). Overall, our variant DLOM-DA yields the best average performance among all compared methods.

Table 3 reports trait-level results on the dataset. Compared to the baseline, our DLOM leads to improvements on most traits and increases the average QWK from 0.646 to 0.685, indicating the effectiveness of explicitly modeling ordinal score decisions. Incorporating distance-aware regularization further improves performance and achieves the best average results among all compared methods, raising the average QWK to 0.697. Notably, the distance-aware variant yields consistent gains across all traits, suggesting that the additional ordinal regularization helps stabilize score predictions in the text-only setting.

The prompt-level results on the benchmark are in Table 4. The DLOM increases the average QWK from 0.676 to 0.713, and incorporating distance-

Model	Prompts								Avg
	1	2	3	4	5	6	7	8	
HISK	0.674	0.586	0.651	0.681	0.693	0.709	0.641	0.516	0.644
MTL-BiLSTM	0.670	0.611	0.647	0.708	0.704	0.712	0.684	0.581	0.665
ArTS	0.708	0.706	0.704	0.767	0.723	0.776	0.749	0.603	0.717
SFT-Gen (baseline)	0.675	0.651	0.721	0.770	0.706	0.753	0.556	0.572	0.676
DLOM	0.682	0.702	0.735	0.789	0.730	0.778	0.664	0.621	0.713
DLOM-DA	0.681	0.696	0.743	0.795	0.728	0.780	0.688	0.646	0.720

Table 4: Prompt-level results on the ASAP/ASAP++ benchmark.

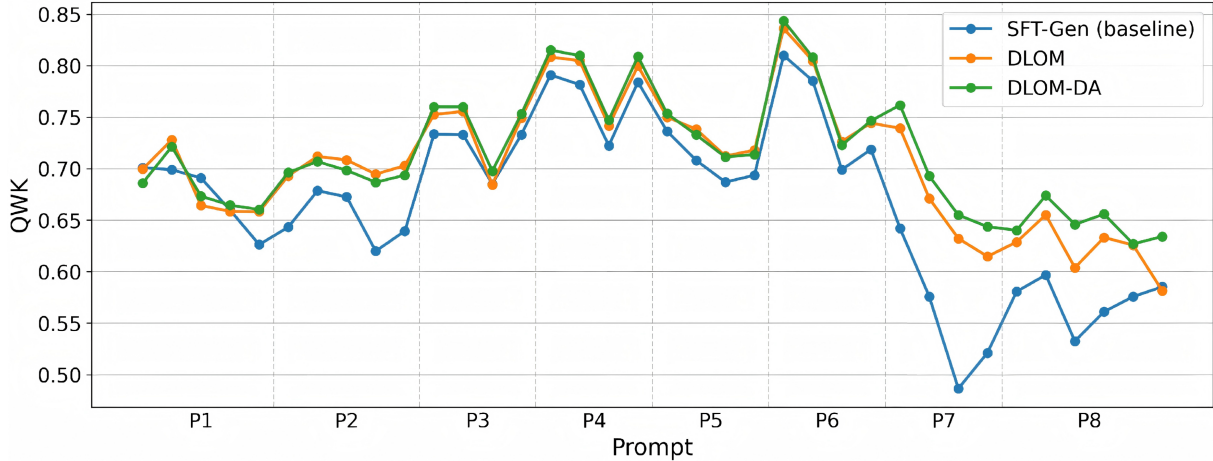


Figure 3: Prompt-wise QWK trends across different models on ASAP/ASAP++. Each point corresponds to a specific trait under a given prompt. Vertical dashed lines indicate boundaries between prompts.

aware regularization further improves performance to 0.720. Specifically, the decision-level ordinal modeling achieves gains on all prompts. Overall, the text-only results corroborate our findings in the multimodal setting. Even without visual inputs, explicitly modeling essay scoring at the decision level leads to consistent performance improvements, demonstrating that the proposed formulation’s generalization ability. More details are provided in Appendix B.7.

4.5 Analysis and Discussion

Prompt-wise Analysis. To better understand the robustness of different modeling paradigms with respect to prompt variations, we conduct a prompt-wise analysis on the text-only dataset. Figure 3 presents QWK trends across prompts, where we evaluate the corresponding traits under each prompt. While absolute performance varies across prompts and traits, SFT-Gen exhibits noticeable fluctuations, with a sustained performance degradation on P7 and P8 (i.e., consistently lower QWK across the associated traits). In contrast, DLOM yields more consistent trends across prompts, substantially mitigating extreme degra-

Method	Avg. QWK
SFT-Gen (baseline)	0.492
SFT + Ordinal Loss ($w=0.01$)	0.472
SFT + Ordinal Loss ($w=0.1$)	0.481
DLOM (ours)	0.504

Table 5: Effect of decision-level formulation on the multimodal EssayJudge dataset.

dations observed in the SFT baseline. This effect is further strengthened by incorporating distance-aware ordinal modeling, while DLOM-DA consistently maintains higher QWK scores across most prompts. These observations suggest that explicitly modeling scoring decisions in the ordinal space reduces sensitivity to prompt-specific variations and leads to more robust essay scoring behavior.

Necessity of Decision-level Formulation. A natural question is whether ordinal structure can be incorporated into generation-based essay scoring by simply adding an ordinal loss term, without modifying the underlying modeling formulation. In this setting, the ordinal loss is introduced as an auxiliary regularization term, with w controlling its relative weight to the token-level generation objective. To examine this alternative, we compare decision-level

Inference Strategy	LA	LD	CH	GA	GD	PA	AC	JP	OS	EL	Avg.
Text-only Decision	0.524	0.551	0.474	0.541	0.559	0.414	0.249	0.342	0.525	0.363	0.454
Multimodal-only Decision	0.616	0.520	0.544	0.580	0.602	0.499	0.242	0.416	0.515	0.472	0.501
DLOM-GF (ours)	0.624	0.551	0.569	0.589	0.613	0.506	0.251	0.447	0.527	0.479	0.516

Table 6: Per-trait comparison of different decision-level inference strategies on the multimodal EssayJudge dataset.

ordinal modeling with generation-based SFT augmented with ordinal loss under different values of w on the multimodal dataset. As shown in Table 5, across the tested loss weights, generation-based models with ordinal regularization consistently underperform decision-level ordinal modeling, and even fall below the original SFT baseline.

These results demonstrate that explicitly reformulating essay scoring as a decision over a fixed, ordered score space leads to a higher average QWK. Unlike loss-level ordinal constraints, which act as auxiliary signals on token-level generation, decision-level ordinal modeling embeds ordinality directly into the prediction space and the decision process itself. This formulation encourages the model to perform score-wise comparisons and reason explicitly over relative score preferences, rather than relying on indirect token-level supervision. These results suggest that, under our experimental setting, ordinal structure in essay scoring is more effectively leveraged when modeled at the decision level rather than imposed solely through loss-level regularization.

Effect of Decision-Level Fusion. Building on the decision-level formulation, we further investigate whether fusing modality-specific score logits can yield additional gains in multimodal scoring. Table 6 presents a per-trait comparison of different decision-level inference strategies on the EssayJudge dataset. Overall, decision-level fusion achieves the highest average QWK (0.516), outperforming both text-only (0.454) and multimodal-only (0.501) inference. This indicates that, even when the backbone and the decision objective are fixed, fusing modality-specific score logits at the decision stage provides a clear additional benefit.

At the trait level, decision-level fusion attains the best performance across all scoring dimensions. While multimodal-only inference generally outperforms text-only inference, it still consistently lags behind decision-level fusion. These results suggest that textual and visual cues provide complementary information, and that the adaptive decision-level gating mechanism can selectively leverage

the more informative modality, leading to more reliable multimodal essay scoring decisions.

Design Implications. Our analysis suggests several practical insights for LLM-based essay scoring. First, in the text-only setting, decision-level ordinal modeling exhibits improved robustness to prompt variations, and additional distance-aware regularization can further enhance performance stability. Second, the key to leveraging ordinality lies in the decision-level formulation: by casting scoring as an explicit decision over a fixed ordered score set, ordinality is encoded directly in the prediction space and the decision mechanism. Third, for multimodal essay scoring, decision-level fusion with an adaptive gate provides an effective way to integrate textual and visual cues in the ordinal score space without requiring representation-level fusion. Overall, these observations highlight the importance of aligning modeling choices with the structure of the scoring decision itself.

5 Conclusion

In this work, we revisited LLM-based automated essay scoring from a decision-centric perspective. We reformulated scoring as an explicit ordinal decision over a fixed, ordered score set, enabling direct optimization in the score space without relying on decode-and-parse generation. On top of this formulation, we introduced a decision-level gated fusion module for multimodal AES and a distance-aware regularization variant for text-only scoring.

Empirically, decision-level ordinal modeling improves over generation-based SFT on the multimodal EssayJudge dataset, and gated fusion yields further gains when visual relevance is heterogeneous. On the text-only ASAP/ASAP++ benchmark, the same formulation remains effective, with distance-aware regularization providing additional improvements. Overall, our results highlight that explicitly modeling the ordinal nature of scoring can yield more reliable LLM-based assessment, and the same view may extend to other ordinal or decision-centric evaluation tasks.

574 Limitations

575 This work focuses on decision-level modeling for
576 essay scoring and evaluates its effectiveness on a
577 limited set of backbone models and datasets. Al-
578 though the proposed formulation and fusion strat-
579 egy demonstrate consistent improvements, our ex-
580 periments are conducted primarily with Qwen-
581 based large language models. Further validation on
582 a broader range of architectures may be necessary
583 to assess the generality of the approach.

584 In addition, while our decision-level fusion
585 framework is applicable to multimodal settings,
586 distance-aware regularization is only explored in
587 the text-only scenario. We leave a systematic inves-
588 tigation of integrating additional ordinal constraints
589 with multimodal fusion to future work.

590 References

591 Dimitrios Alikaniotis, Helen Yannakoudakis, and Marek
592 Rei. 2016. Automatic text scoring using neural net-
593 works. *arXiv preprint arXiv:1606.04289*.

594 Azhari Azhari, Agus Santoso, Anak Agung Putri Ratna,
595 and Jasson Prestiliano. 2024. Optimization of aes
596 using bert and bilstm for grading the online exams.
597 *International Journal of Intelligent Engineering &*
598 *Systems*, 17(5).

599 Ngoc My Bui and Jessie S Barrot. 2025. Chatgpt as
600 an automated essay scoring tool in the writing class-
601 rooms: how it compares with human scoring. *Ed-*
602 *ucation and Information Technologies*, 30(2):2041–
603 2058.

604 SeongYeub Chu, Jong Woo Kim, Bryan Wong, and
605 Mun Yong Yi. 2025. Rationale behind essay scores:
606 Enhancing s-llm’s multi-trait essay scoring with ra-
607 tionale generated by llms. In *Findings of the Associ-*
608 *ation for Computational Linguistics: NAACL 2025*,
609 pages 5796–5814.

610 J. Cohen. 1968. Weighted kappa: nominal scale agree-
611 ment with provision for scaled disagreement or par-
612 tial credit. *Psychological Bulletin*, 70(4):213–220.

613 Mădălina Cozma, Andrei Butnaru, and Radu Tudor
614 Ionescu. 2018. Automated essay scoring with string
615 kernels and word embeddings. In *Proceedings of the*
616 *56th Annual Meeting of the Association for Com-*
617 *putational Linguistics (Volume 2: Short Papers)*,
618 pages 503–509, Melbourne, Australia. Association
619 for Computational Linguistics.

620 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
621 Kristina Toutanova. 2019. Bert: Pre-training of deep
622 bidirectional transformers for language understand-
623 ing. In *Proceedings of the 2019 conference of the*
624 *North American chapter of the association for com-*
625 *putational linguistics: human language technologies,*
626 *volume 1 (long and short papers)*, pages 4171–4186.

Semire Dikli. 2006. An overview of automated scoring
of essays. *The Journal of Technology, Learning and*
Assessment, 5(1).

Heejin Do, Yunsu Kim, and Gary Lee. 2024a. Autore-
gressive score generation for multi-trait essay scoring.
In *Findings of the Association for Computational Lin-*
guistics: EACL 2024, pages 1659–1666.

Heejin Do, Sangwon Ryu, and Gary Lee. 2024b. Au-
toregressive multi-trait essay scoring via reinforc-
ement learning with scoring-aware multiple rewards.
In *Proceedings of the 2024 Conference on Empiri-*
cal Methods in Natural Language Processing, pages
16427–16438.

Fei Dong and Yue Zhang. 2016. Automatic features for
essay scoring—an empirical study. In *Proceedings of*
the 2016 conference on empirical methods in natural
language processing, pages 1072–1077.

Fei Dong, Yue Zhang, and Jie Yang. 2017. Attention-
based recurrent convolutional neural network for au-
tomatic essay scoring. In *Proceedings of the 21st*
conference on computational natural language learn-
ing (CoNLL 2017), pages 153–162.

Lulu Dong, Lin Li, HongChao Ma, YeLing Liang, and 1
others. 2021. Automated chinese essay scoring using
pre-trained language models. In *CS & IT Confer-*
ence Proceedings, volume 11. CS & IT Conference
Proceedings.

Eibe Frank and Mark Hall. 2001. A simple approach to
ordinal classification. In *Machine Learning: ECML*
2001, pages 145–156, Berlin, Heidelberg. Springer
Berlin Heidelberg.

Pedro Antonio Gutiérrez, María Pérez-Ortiz, Javier
Sánchez-Monedero, Francisco Fernández-Navarro,
and César Hervás-Martínez. 2016. Ordinal regres-
sion methods: Survey and experimental study. *IEEE*
Transactions on Knowledge and Data Engineering,
28(1):127–146.

Ralf Herbrich, Thore Graepel, and Klaus Obermayer.
2000. Large margin rank boundaries for ordinal re-
gression. In *Advances in Large-Margin Classifiers*.
The MIT Press.

Zhaoyi Joey Hou, Alejandro Ciuba, and Xiang Lor-
raine Li. 2025. Improve llm-based automatic es-
say scoring with linguistic features. *arXiv preprint*
arXiv:2502.09497.

Valen E Johnson. 1996. On bayesian analysis of multi-
rater ordinal data: An application to automated essay
grading. *Journal of the American Statistical Associa-*
tion, 91(433):42–51.

Rahul Kumar, Sandeep Mathias, Sriparna Saha, and
Pushpak Bhattacharyya. 2022. Many hands make
light work: Using essay traits to automatically score
essays. In *Proceedings of the 2022 Conference of*
the North American Chapter of the Association for
Computational Linguistics: Human Language Tech-
nologies, pages 1485–1495.

791 David M Williamson, Isaac I Bejar, and Anne S Hone.
792 1999. ‘mental model’ comparison of automated and
793 human scoring. *Journal of Educational Measure-*
794 *ment*, 36(2):158–184.

795 Wei Xia, Shaoguang Mao, and Chanjing Zheng. 2024.
796 Empirical study of large language models as au-
797 tomated essay scoring tools in english composi-
798 tion_taking toefl independent writing task for exam-
799 ple. *arXiv preprint arXiv:2401.03401*.

800 Ruosong Yang, Jiannong Cao, Zhiyuan Wen, Youzheng
801 Wu, Xiaodong He, and 1 others. 2020. Enhancing
802 automated essay scoring performance via fine-tuning
803 pre-trained language models with combination of
804 regression and ranking. In *Proceedings of ACL 2020*.

805 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-
806 dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
807 and Le Hou. 2023. Instruction-following evalu-
808 ation for large language models. *arXiv preprint*
809 *arXiv:2311.07911*.

A Additional Clarification on Decision Interfaces

This appendix clarifies how our decision-level formulation differs from two common alternatives for LLM-based AES: generation-based scoring and standard classification heads.

Generation-based scoring. A generation-based scorer predicts a score by decoding a textual output (e.g., a score token or a short phrase) and then mapping the decoded text to an ordinal label. Formally, it relies on $\hat{y} = \text{Parse}(\text{Decode}(p_\theta(\cdot | e, v)))$, which introduces sensitivity to decoding strategies (e.g., sampling/greedy) and prompt/output formatting.

CLS-head classification (with/without ordinal losses). A standard classifier treats the LLM as an encoder, pools hidden states, and applies an additional $(K+1)$ -way head:

$$\mathbf{o} = W_{\text{cls}} \cdot \text{pool}(H) + b, \quad \hat{y} = \arg \max_k o_k,$$

optionally augmented with ordinal regression objectives (e.g., CORAL-style losses or distance-based regularizers) on top of \mathbf{o} . This design introduces a separate decision space parameterized by W_{cls} , which is not tied to the LM output distribution.

Ours: decision-level score-logit extraction. In contrast, we reuse the LM head and make the decision directly in a predefined ordered score space by restricting the final-step vocabulary logits to a fixed score-token set \mathcal{S} :

$$\mathbf{z} = \text{Logits}(e, v)_L[\mathcal{S}] \in \mathbb{R}^{K+1}, \quad \hat{y} = \arg \max_{k \in \mathcal{Y}} z_k.$$

This removes the need for decoding and avoids introducing an additional classification head, while yielding an explicit score-wise decision representation. Operating in this score-logit space also enables our decision-level gated fusion, which adaptively combines modality-specific evidence via a convex combination of $\mathbf{z}^{(m)}$ and $\mathbf{z}^{(t)}$ under the same decision objective. Moreover, this interface is LLM-native: it leverages the model’s learned token-preference distribution and avoids relying on prompt- or decoding-specific behaviors when producing a final score.

B Experiments

B.1 EssayJudge

Table 7 shows the detailed information about EssayJudge dataset (Su et al., 2025a). It consists of 1,054 multimodal essays written as part of IELTS Writing Task 1. In this task, candidates are required to write a report based on visual data such as charts, graphs, or diagrams. The dataset is specifically designed to assess multimodal reasoning in automated essay scoring, where visual inputs play a significant role in the essay’s content and quality.

Statistic	Number
Total Multimodal Essays	1,054
Image Type	
- Single-Image	703 (66.7%)
- Multi-Image	351 (33.3%)
Multimodal Essay Type	
- Flow Chart	305 (28.9%)
- Bar Chart	211 (20.0%)
- Table	153 (14.5%)
- Line Chart	145 (13.8%)
- Pie Chart	71 (6.7%)
- Map	62 (5.9%)
- Composite Chart	107 (10.2%)

Table 7: Key statistics of EssayJudge dataset.

B.2 ASAP/ASAP++

We use the open-sourced ASAP/ASAP++, as summarized in Table 8, different prompts are assessed with distinct traits, each having varied score ranges (Mathias and Bhattacharyya, 2018, 2020). It consists of English essays written by American 7th to 10th-grade high school students across eight prompts.

Prompt	# Essays	Traits	Score Range
1	1785	Content, WC, Org, SF, Conv	1 - 6
2	1800	Content, WC, Org, SF, Conv	1 - 6
3	1726	Content, PA, Nar, Lang	0 - 3
4	1772	Content, PA, Nar, Lang	0 - 3
5	1805	Content, PA, Nar, Lang	0 - 4
6	1800	Content, PA, Nar, Lang	0 - 4
7	1569	Content, Org, Conv, Style	0 - 6
8	723	Content, WC, Org, SF, Conv, Voice	2 - 12

Table 8: Composition of the ASAP/ASAP++ combined dataset. The prompt defines the writing theme. Abbreviations: WC (Word Choice), Org (Organization), SF (Sentence Fluency), Conv (Conventions), PA (Prompt Adherence), Nar (Narrativity), Lang (Language).

Model	Source	URL
Qwen3-VL-8B-Instruct	local checkpoint	https://huggingface.co/Qwen/Qwen3-VL-8B-Instruct
Qwen3-8B	local checkpoint	https://huggingface.co/Qwen/Qwen3-8B

Table 9: Model sources for the multimodal and text-only essay scoring tasks.

B.3 Quadratic Weighted Kappa (QWK)

Quadratic Weighted Kappa (QWK) is a commonly used metric in automated essay scoring (AES) to measure the agreement between model predictions and human raters. It accounts for the ordinal nature of the score scales, where errors between adjacent score levels are less severe than those between distant levels.

The QWK is defined as:

$$QWK = 1 - \frac{\sum_{i,j} w_{ij} \cdot O_{ij}}{\sum_{i,j} w_{ij} \cdot E_{ij}}$$

where: - O_{ij} is the observed agreement between categories i and j , - E_{ij} is the expected agreement between categories i and j , - w_{ij} is the weight assigned to the difference between score levels i and j .

The weight w_{ij} is typically defined as:

$$w_{ij} = \frac{(i - j)^2}{(K - 1)^2}$$

where K is the number of possible score categories.

B.4 Model Selecting and Sources

For our experiments, we selected two large language models: Qwen3-VL-8B-Instruct and Qwen3-8B, chosen for their distinct strengths in multimodal and text-only tasks, respectively.

Qwen3-VL-8B-Instruct : This is a multimodal model designed to handle both textual and visual inputs, making it ideal for essay scoring tasks where essays are accompanied by images, charts, or other forms of visual data. And its open-source availability makes it accessible for replication and future work in this area.

Qwen3-8B : This is a text-only model, suitable for tasks that involve only textual input. We use it as a baseline to compare the performance of our multimodal model, allowing us to isolate the effect of incorporating visual information into the scoring process.

These two models represent a balanced comparison between multimodal and text-only settings,

allowing us to evaluate the robustness and generality of our proposed method across different types of inputs. Table 9 shows their checkpoints' source.

B.5 Prompt Template and Rubrics

Training on multimodal dataset. The training was conducted using a unified prompt template as shown in Figure 4 and the rubric for each trait is publicly available and can be accessed in the original dataset repository (Su et al., 2025a).

```

prompt_template = """
You are an IELTS examiner. Based on the image, question and essay, please
score the "{trait}" dimension using the rubric.

Rubric: {rubric}
Essay title: "{question}"
Essay: "{essay}"
Image: <image>

Your only output must be a single number (the score) from the rubric. No other
text.

Score:
"""

```

Figure 4: Prompt Template for Multimodal Dataset

Training on text-only datasets. We train all text-only models using a unified prompt template (Figure 5). Trait rubrics and prompt definitions for ASAP/ASAP++ are publicly available.¹

```

prompt = (
    f"You must score the essay's \"{trait}\" dimension using the rubric.\n"
    f"Rubric: {rubric}\n"
    f"Essay: {essay}\n"
    f"Your only output must be a single number (the score) from the rubric.\n"
    f"Score:"
)

```

Figure 5: Prompt Template for Text-only Dataset

B.6 Main Results

We observe relatively large variance across cross-validation folds for several traits, including the generation-based baseline. This variance is mainly attributed to the limited size and heterogeneous nature of the EssayJudge dataset, rather than instability introduced by the proposed methods. Importantly, decision-level modeling and gated fusion do not exhibit higher variance than the baseline,

¹ASAP: <https://www.kaggle.com/c/asap-aes>.
ASAP++: <https://lwsam.github.io/ASAP++/Irec2018.html>.

Model	Content	PA	Lang	Nar	Org	Conv	WC	SP	Style	Voice	Avg
SFT (Generation)	± 0.031	± 0.026	± 0.026	± 0.025	± 0.051	± 0.066	± 0.044	± 0.049	± 0.030	± 0.056	± 0.040
Decision-level (Ordinal)	± 0.026	± 0.025	± 0.028	± 0.025	± 0.036	± 0.036	± 0.044	± 0.032	± 0.028	± 0.083	± 0.036
Decision-level + Distance-aware	± 0.030	± 0.023	± 0.035	± 0.031	± 0.035	± 0.038	± 0.037	± 0.032	± 0.020	± 0.083	± 0.036

Table 10: Standard deviation of QWK across cross-validation folds on the ASAP/ASAP++ benchmark, reported at the trait level.

Model	1	2	3	4	5	6	7	8	Avg
SFT (Generation)	± 0.053	± 0.036	± 0.032	± 0.028	± 0.021	± 0.017	± 0.054	± 0.057	± 0.037
Decision-level (Ordinal)	± 0.020	± 0.033	± 0.033	± 0.018	± 0.035	± 0.016	± 0.034	± 0.054	± 0.031
Decision-level + Distance-aware	± 0.029	± 0.031	± 0.040	± 0.020	± 0.036	± 0.017	± 0.030	± 0.054	± 0.032

Table 11: Standard deviation of QWK across cross-validation folds on the ASAP/ASAP++ benchmark, reported at the prompt level.

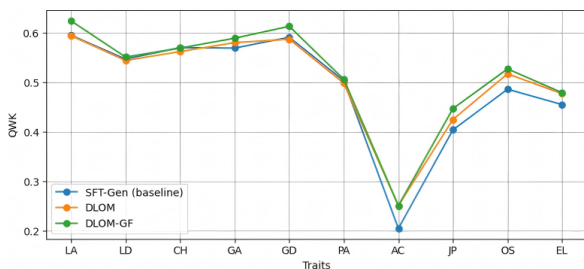


Figure 6: Trait-wise QWK comparison on the Essay-Judge dataset.

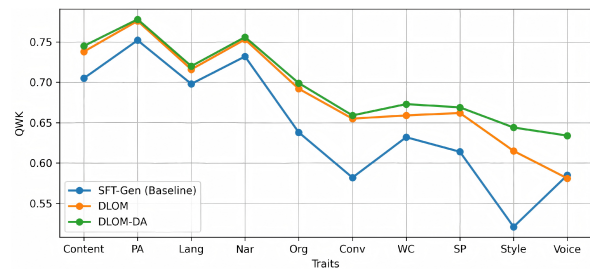


Figure 7: Trait-wise QWK trends on the ASAP/ASAP++ benchmark.

933 indicating comparable stability under the same ex-
 934 perimental setting.

935 Figure 6 provides a visualization of the trait-
 936 level QWK results reported in Table 2. The fig-
 937 ure highlights consistent improvements of decision-
 938 level modeling over the generation-based baseline
 939 across most traits, with more pronounced gains
 940 on discourse-level dimensions. The gated fusion
 941 mechanism further strengthens performance, par-
 942 ticularly on traits with lower baseline scores.

943 B.7 Text-only Results on ASAP/ASAP++

944 **Variance analysis.** Tables 10 and 11 report the
 945 standard deviation of QWK across cross-validation
 946 folds at the trait and prompt levels. Across both
 947 views, decision-level modeling and the distance-
 948 aware variant achieve higher mean performance
 949 without increasing variability: the average standard
 950 deviation remains comparable to the generation
 951 baseline. This suggests that the observed improve-
 952 ments are not obtained at the cost of reduced stabil-
 953 ity.

954 **Visualizations of text-only results.** Figures 7
 955 and 8 visualize the trait-level and prompt-level
 956 QWK results from Tables 3 and 4, respectively.

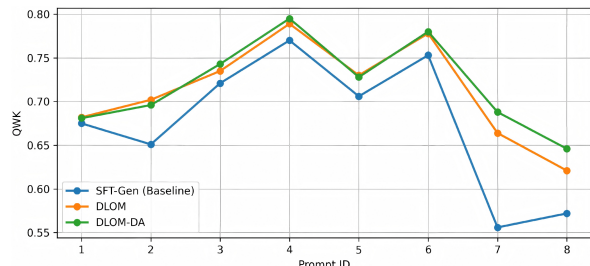


Figure 8: Prompt-wise QWK trends on the ASAP/ASAP++ benchmark.

957 Overall, decision-level modeling improves over
 958 the generation-based baseline on most traits and
 959 prompts, and the distance-aware variant provides
 960 further (often more consistent) gains. Improve-
 961 ments are particularly pronounced on higher-level
 962 traits that require holistic judgments of content and
 963 organization.