

# OPENREVIEWER: PREDICTING CONFERENCE DECISIONS WITH LLMs AND BEYOND

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Anonymous authors  
Paper under double-blind review

## ABSTRACT

The rapid growth of AI conference submissions has strained the peer-review system, motivating interest in AI-assisted review. Yet it remains unclear how reliably such systems approximate human judgment, which relies on domain expertise and nuanced reasoning. To address this challenge, we introduce OpenReviewer, a model designed to directly predict conference acceptance decisions rather than generate full reviews. Using ICLR 2024–2025 data, we evaluate large language models (LLMs), vision–language models (VLMs), and interpretable statistical models. Results show that text-only LLMs with continual pre-training outperform multimodal counterparts, achieving up to 78.5% accuracy on balanced datasets (vs. 50% random baseline). White-box statistical models further provide interpretability through feature analysis, revealing that structural attributes (e.g., paper length, section balance, citation engagement) are consistently predictive. Beyond average accuracy, a confidence-stratified utility analysis shows that the top 10% most confident predictions reach 92.92% overall precision, enabling reliable triage of “obvious” accepts and rejects while exposing areas of uncertainty. Overall, our findings demonstrate both the promise and limitations of AI-involved peer review: current models can reduce workload and aid submission reviewing, but fall short of reliably replacing expert judgment.

## 1 INTRODUCTION

The peer-review process is becoming increasingly unsustainable as submissions to top-tier AI conferences continue to grow at an unprecedented pace, as shown in Figure 1<sup>1</sup>. This explosive growth places pressure on program committees and reviewers, leading to heavier workloads and concerns over the quality and consistency of reviews (Lawrence, 2022; Beygelzimer et al., 2023; Kim et al., 2025; Schaeffer et al., 2025). For authors, uncertainty around submission outcomes and suboptimal venue choices can negatively influence research trajectories and academic career development (e.g., timely PhD graduation) (Kousha & Thelwall, 2024; Yang, 2025).

Recent work has explored AI-assisted review generation as a potential solution, where models take papers as inputs to generate reviews. (Sukpanichnant et al., 2024; Ye et al., 2024; Shin et al., 2025). To our knowledge, no existing work uses LLMs to predict acceptance directly from the paper content itself. Reliable acceptance prediction could guide authors in developing submission strategies, while helping committees triage obviously good/low-quality papers and allocate human review resources more effectively. Therefore, we propose

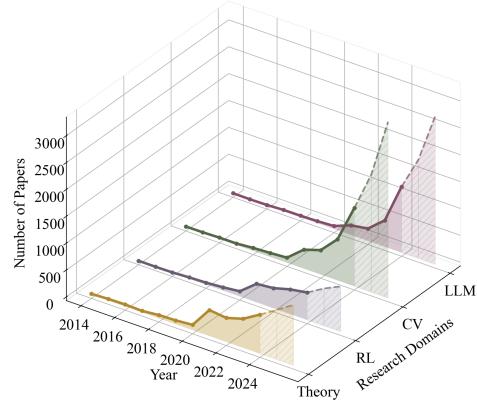


Figure 1: Number of papers accepted by NeurIPS from 2015 to 2024 across four major research domains, with the dashed line indicating the predicted trend.

<sup>1</sup>Our data sources include official conference announcements and the Paper Copilot platform (<https://papercopilot.com/>).

054 OpenReviewer, a LLM-based model that predicts the acceptance of papers submitted to AI con-  
 055 ferences.

056 OpenReviewer is developed using conference submissions and corresponding acceptance infor-  
 057 mation collected from the OpenReview platform<sup>2</sup>. In particular, our dataset consists of submission  
 058 records from the International Conference on Learning Representations (ICLR), chosen for its  
 059 broad coverage of AI topics and the openness of its submission and decision records. The training  
 060 dataset includes three components of papers: (1) *textual content*, consisting primarily of anonymized  
 061 manuscript text; (2) *visual information*, such as system figures and charts; and (3) interpretable, man-  
 062 ually engineered *statistical features*.

063 We adopt prompt-based fine-tuning (Shi & Lipani, 2023) to help LLMs understand the given task,  
 064 combined with a decoupled label loss (Tam et al., 2021) to encourage the use of vocabulary tokens  
 065 (e.g. Yes, No) as labels during training. We explore three approaches for this task: text-only large  
 066 language models, vision-language models, and white-box statistical classifiers. Among text-only  
 067 models, continued pretraining (CPT) on unlabeled corpora prior to fine-tuning yields the best result,  
 068 achieving **78.5%** accuracy on a label-balanced dataset (50% random baseline). For VLMs mod-  
 069 els, unsurprisingly, incorporating image inputs consistently outperforms text-only inputs. We also  
 070 provide qualitative analyses highlighting cases where images help and where they mislead. In addi-  
 071 tion, we conduct a white-box analysis of statistical features by extracting 29 heuristic quantitative  
 072 features across eight categories. Using only these features, a Random Forest classifier (Breiman,  
 073 2001) attains a surprisingly strong 74.2% accuracy, surpassing VLMs. Finally, a model confidence-  
 074 stratified analysis for OpenReviewer shows that within the top 10% confidence slice, covering  
 075 up to 53.06% of predictions, LLMs achieve 93.09% precision on the Accept class, with a compa-  
 076 rable trend observed for the Reject class. This enables reliable triage of clear accepts and rejects  
 077 while routing uncertain cases for human review. Overall, these findings indicate that AI can support  
 078 peer review by reducing the workload on straightforward submissions, while human experts remain  
 079 essential for more nuanced judgments.

## 081 2 RELATED WORK

082 **Peer Review Analysis.** The peer review process, particularly in rapidly evolving fields like AI, is  
 083 facing a sustainability crisis with reviewer overload and declining quality (Chen et al., 2025; Kim  
 084 et al., 2025). While LLMs have been explored to automate or assist reviewing, their readiness lacks  
 085 validation. Large-scale experiments show LLMs can distinguish paper quality but exhibit significant  
 086 biases (Pataranutaporn et al., 2025), with researchers warning against premature deployment due to  
 087 risks in factual accuracy and logical reasoning (Ye et al., 2024). New evaluation methods identify  
 088 “blind spots” in LLM reviews, revealing that they often miss crucial methodological flaws such as  
 089 experimental design issues, statistical significance problems, and logical inconsistencies in argu-  
 090 mentation (Shin et al., 2025). Improvement efforts include structured argumentative review frame-  
 091 works (Sukpanichnant et al., 2024) and graph reasoning systems over reviewer-author debates (Tae-  
 092 choyotin & Acuna, 2025). Additionally, AI-assisted reviews create an “AI review lottery,” inflating  
 093 scores and masking weaknesses (Latona et al., 2024). These challenges prompt calls for systemic  
 094 reform, including transparent processes and reviewer rewards (Yang, 2025; Ye et al., 2024), dedi-  
 095 cated critique tracks (Schaeffer et al., 2025), and lessons from platforms like OpenReview (Wang  
 096 et al., 2023).

097 **Paper Quality Modeling.** Recent advancements in LLMs have spurred significant research into  
 098 computationally modeling the quality of scholarly papers in the form of evaluation, revision, and  
 099 generation. Research focuses on automated assessment using domain-aware retrieval and latent  
 100 reasoning (Zheng et al., 2025), verifiable claim extraction (Song et al., 2024), and retraction pre-  
 101 diction for scientific integrity (Yang & Jia, 2025). Beyond evaluation, quality models support pa-  
 102 per improvement through human-AI collaborative revision frameworks (Fragiadakis et al., 2024;  
 103 Dong et al., 2022) and fully automated generation systems like ARISE, which uses explicit quality  
 104 rubrics (Schneider, 2025).

105  
 106  
 107 <sup>2</sup><https://openreview.net/>

108 **LLM-based Document Classification.** LLMs shift document classification from traditional fine-  
 109 tuning to prompt-based, few-shot learning. This involves reformulating classification tasks as cloze  
 110 questions, enabling strong performance with minimal labeled data (Schick & Schütze, 2021). Re-  
 111 cent studies further demonstrate that continued pretraining can significantly enhance prompt-tuning  
 112 effectiveness, making it an even more powerful learning approach (Chen et al., 2022). Despite these  
 113 advances, LLM-based classification faces several challenges. Raw LLM outputs often suffer from  
 114 miscalibration issues, necessitating the development of context-aware calibration techniques (Zhao  
 115 et al., 2021). Additionally, the direct application of LLMs as zero-shot or few-shot classifiers shows  
 116 promise but remains task-dependent, requiring careful model selection for specialized domains such  
 117 as classifying scientific revision intents (Ruan et al., 2024). To address these limitations, researchers  
 118 have developed hybrid and advanced approaches. Hybrid models like DeepCCP successfully inte-  
 119 grate semantic understanding with citation network structure to achieve more accurate classifica-  
 120 tion (Zhao & Feng, 2022). Furthermore, advanced approaches explore classification through gen-  
 121 eration tasks, including benchmarking LLMs on writing paper sections (Garg et al., 2025) and de-  
 122 veloping multi-agent frameworks for paper reproduction (Miao et al., 2025). These developments  
 123 highlight the evolution towards deeper, context-aware reasoning.

### 124 3 OPENREVIEWER FOR PREDICTING ACCEPTANCE

125 We formulate paper acceptance prediction as a binary classification problem. State-of-the-art LLMs  
 126 and VLMs are inherently generative, making them not directly applicable to traditional classification  
 127 tasks. To use the capabilities of these powerful pre-trained generative models without training a  
 128 new classification head from scratch<sup>3</sup>, we adopt a prompt-based fine-tuning strategy Ruan et al.  
 129 (2024); Schick & Schütze (2021); Shi & Lipani (2023). Specifically, we design an instructive prompt  
 130 template  $\mathcal{T}$  that presents the paper’s features within a natural-language query and guides the model  
 131 to generate a decision token corresponding to one of the two target classes: *accept* or *reject*. The  
 132 template example is given in App. D.

#### 133 3.1 CONTINUAL PRE-TRAINING

134 Continual pre-training (CPT) extends the training of large generative models on additional unlabeled  
 135 corpora to improve their adaptability to new domains and evolving data distributions (Gururangan  
 136 et al., 2020; Chen et al., 2023). It is widely adopted in industry-scale generative systems, where  
 137 models are periodically updated with fresh data to sustain relevance and maintain competitive per-  
 138 formance (Gururangan et al., 2020; Chen et al., 2023; Ke et al., 2023; Elhady et al., 2025). The  
 139 training objective typically follows next-token prediction, formalized as

$$140 \mathcal{L}_{\text{CPT}} = - \sum_{t=1}^T \log P_{\theta}(x_t \mid x_{<t}), \quad (1)$$

141 which maximizes the likelihood of generating each token  $x_t$  given its preceding context  $x_{<t}$  and  
 142 model parameters  $\theta$ . In this paper, we also explore continual pre-training to adapt general-purpose  
 143 base models to the academic peer-review scenario before fine-tuning on the classification task. We  
 144 present the effectiveness of CPT in Section 4.3, with further training details provided in the App. E.  
 145 Unless otherwise specified, CPT is used as the default post-training strategy for our textual models  
 146 before fine-tuning.

#### 147 3.2 INPUT SETTINGS

148 Given a paper input instance  $x$  and prompt template  $\mathcal{T}(x)$ , the model defines a conditional proba-  
 149 bility over the label verbalizer (Tam et al., 2021). We consider two input configurations for  $\mathcal{T}(x)$ :  
 150 *text-only* and *text-image multimodal*. Text-only inputs are anonymized main-body texts from the  
 151 paper manuscripts. The multimodal setting extends the text-only configuration by additionally in-  
 152 corporating visual features extracted from figures in the paper. Details of the PDF preprocessing and  
 153 figure extraction procedure are provided in App. C. Formally,

$$154 \mathcal{T}(x) = \phi(x^{(\text{text})} \oplus x^{(\text{figure})}) \quad (2)$$

155 <sup>3</sup>Our initial experiments with training a classification head on top of a pre-trained LLMs resulted in lower  
 156 accuracy and slower convergence compared to prompt-based generation.

162 where  $\phi$  is modality-specific encoding determined by the multimodal model, and  $x^{(\text{figure})}$  is an  
 163 optional input.  $\mathcal{T}(x)$  packs all modalities into a single token sequence consumable by the model.  
 164

165 We then append a designated decision slot and generate only at this position, defining  $\mathcal{T}^{\text{dec}} =$   
 166  $\mathcal{T}(x) \oplus [\text{label\_mask}]$ . *Verbalizer*  $\mathcal{V}$  maps each candidate label token  $v$  to a class; this allows many-  
 167 to-one mappings (e.g.,  $\{\text{yes}, \text{accept}\} = 1$ ;  $\{\text{reject}, \text{no}\} = 0$ ).

### 168 3.3 MODEL TRAINING OBJECTIVE

170 Given an input  $\mathcal{T}^{\text{dec}}$  and its corresponding ground-truth label  $y^*$ , we apply supervision *only* at the  
 171 decision position, masking all other positions. Following ADAPET (Tam et al., 2021), we define the  
 172 Vocabulary Decoupled Label Loss (VDLL). Let  $z_\theta(t \mid \mathcal{T}^{\text{dec}})$  denote candidate labels logits at the  
 173 decision slot  $t$ , we then define the (3) *restricted softmax* and (4) training objective as:

$$174 \quad \tilde{p}_\theta(t \mid \mathcal{T}^{\text{dec}}) = \frac{\exp(z_\theta(t \mid \mathcal{T}^{\text{dec}}))}{\sum_{a \in \mathcal{V}} \exp(z_\theta(a \mid \mathcal{T}^{\text{dec}}))} \quad (3)$$

$$177 \quad \mathcal{L}_{\text{VDLL}}(\theta) = -\log \sum_{t \in \mathcal{V}_{y^*}} \tilde{p}_\theta(t \mid \mathcal{T}^{\text{dec}}), \quad (4)$$

### 180 3.4 INFERENCE MECHANISM

182 At inference time, we determine the predicted class by comparing the *logit-based* scores of all  
 183 verbalizer candidates at the decision slot. Let  $z_\theta(t \mid \mathcal{T}^{\text{dec}})$  denote the pre-softmax logit assigned  
 184 by the model to token  $t$  at the decision position. We first obtain the token IDs of all verbalizer  
 185 candidates  $\mathcal{V}$ . For each class  $y$ , the score is defined as the maximum logit among its associated  
 186 verbalizer tokens. The final prediction  $\hat{y}$  is then obtained by selecting the class with a higher score,  
 187 for example predicting *Accept* if  $\text{score}(\text{yes}) > \text{score}(\text{no})$  and vice versa:

$$188 \quad \hat{y} = \arg \max_{y \in \mathcal{Y}} \text{score}_\theta(y \mid \mathcal{T}^{\text{dec}}) \quad (5)$$

189 We report a binary decision  $b(\hat{y}) \in \{0, 1\}$ .  
 190

## 191 4 EXPERIMENTS

### 193 4.1 DATA COLLECTION AND PRE-PROCESSING

195 We collect all ICLR 2025 and 2024 submissions and their corresponding final decisions (*accepted*  
 196 or *rejected*) via the OpenReview API-V2. The papers were further partitioned into four main  
 197 subfields based on title keywords: Large Language Models (LLM), Computer Vision (CV), Rein-  
 198 forcement Learning (RL), and Theoretical (Theory). Papers that do not fall into these categories are  
 199 left for future discussion. We build two datasets: the ICLR 2025 dataset, which is naturally imbal-  
 200 anced with a 34/66 accepted-to-rejected split and balanced domain-specific sets from ICLR 2024  
 201 and 2025 with a 50/50 split. Table 6 summarizes the differences. More implementation details are  
 202 explained in App. C.

### 204 4.2 MODELS AND INPUTS

206 We include two categories of models: **text-only** LLMs and **vision-language** models. For the  
 207 text-only LLMs, we select the Qwen-3 family at Qwen3-4B and Qwen3-8B parameter scales  
 208 (Yang et al., 2025). For VLMS, we include Qwen2.5-VL-3B-Instruct (Bai et al., 2025) and  
 209 Gemma-3-4b-bit (Team et al., 2025). Both of these multimodal models are instruction-tuned  
 210 variants. We take the vanilla non-fine-tuned version of each model in a zero-shot setting as the baseline.  
 211 After collecting and preprocessing the papers along with their acceptance outcomes, we fine-tune  
 212 and evaluate the two categories of selected models using the following inputs.

213 **Text-only LLMs:** We first anonymize each paper by removing all information that could reveal  
 214 author identity or acceptance status, including author names, affiliations, email addresses, URLs,  
 215 and header or footer text. Beyond these removals, the input consists of the full manuscript body text  
 and mathematical formulas, but excludes tables and figure captions.

216 Table 1: Performance (%) of LLMs across four domains and the overall aggregation (ALL) on the  
 217 balanced dataset. All models use Qwen3 as the backbone. We compare fine-tuning with CPT against  
 218 the original checkpoints (Orig) at the 4B and 8B parameter scales.

| 220 | SUB-<br>221 DOMAIN | ALL |             |             |             | LLM         |       |       |      | CV   |       |       |             | RL          |             |             |             | THEORY      |             |       |             |             |
|-----|--------------------|-----|-------------|-------------|-------------|-------------|-------|-------|------|------|-------|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------|-------------|-------------|
|     |                    | ACC | MAC-P       | MAC-R       | F1          | ACC         | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1          | ACC         | MAC-P       | MAC-R       | F1          | ACC         | MAC-P       | MAC-R | F1          |             |
| 222 | CPT                | 4B* | 48.7        | 23.8        | 51.0        | 34.7        | 54.2  | 26.9  | 52.4 | 36.5 | 51.0  | 26.6  | 49.5        | 33.7        | 50.6        | 22.2        | 47.4        | 32.6        | 46.1        | 24.3  | 48.6        | 30.2        |
|     |                    | 4B  | 76.4        | 76.3        | 76.4        | 76.4        | 70.2  | 70.1  | 70.2 | 70.1 | 70.2  | 75.3  | 70.4        | 68.7        | 57.4        | 62.5        | 57.3        | 52.4        | 55.9        | 61.5  | 55.5        | 49.1        |
|     |                    | 8B  | <b>78.5</b> | <b>78.5</b> | <b>78.3</b> | <b>78.5</b> | 70.0  | 70.1  | 70.0 | 70.1 | 73.9  | 74.1  | <b>74.0</b> | <b>73.9</b> | <b>67.5</b> | <b>68.0</b> | <b>67.6</b> | <b>67.3</b> | 53.4        | 53.6  | 53.2        | 51.0        |
| 224 | Orig               | 4B* | 51.9        | 24.5        | 52.1        | 33.5        | 54.4  | 28.1  | 51.6 | 38.0 | 51.1  | 27.2  | 49.6        | 37.5        | 51.5        | 21.6        | 49.6        | 31.7        | 47.4        | 22.3  | 47.6        | 31.5        |
|     |                    | 4B  | 67.3        | 71.4        | 68.2        | 66.3        | 68.8  | 69.1  | 68.8 | 68.6 | 71.0  | 71.0  | 71.1        | 71.0        | 59.3        | 59.7        | 59.3        | 58.9        | <b>57.2</b> | 59.5  | <b>56.9</b> | 53.9        |
|     |                    | 8B  | 69.0        | 69.0        | 69.0        | 69.0        | 69.7  | 69.9  | 69.8 | 69.7 | 72.5  | 73.2  | 72.6        | 72.4        | 62.6        | 64.3        | 62.6        | 61.5        | 55.4        | 55.9  | 55.2        | <b>54.0</b> |

226 \* Baseline models. Random guess baseline accuracy is 50%.

228 Table 2: Performance (%) of VLMs across four domains on the balanced dataset. Mac-P and Mac-R  
 229 denote Macro Precision and Macro Recall, respectively.

| 231 | SUB-<br>232 DOMAIN | ALL     |             |             |             | LLM         |             |             |             | CV          |             |             |             | RL          |             |             |             | THEORY      |             |             |             |             |
|-----|--------------------|---------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|     |                    | ACC     | MAC-P       | MAC-R       | F1          | ACC         | MAC-P       | MAC-R       | F1          | ACC         | MAC-P       | MAC-R       | F1          | ACC         | MAC-P       | MAC-R       | F1          | ACC         | MAC-P       | MAC-R       | F1          |             |
| 233 | Qwen-VL            | txt&img | 46.0        | 23.0        | 50.0        | 31.5        | 54.6        | 27.1        | 50.0        | 35.3        | 51.6        | 25.8        | 50.0        | 34.0        | 50.2        | 25.1        | 50.0        | 33.4        | 47.6        | 23.8        | 50.0        | 32.3        |
|     |                    | txt     | 48.0        | 24.0        | 50.0        | 32.4        | 54.6        | 27.3        | 50.0        | 35.3        | 51.6        | 25.8        | 50.0        | 34.0        | 50.2        | 25.1        | 50.0        | 33.4        | 47.6        | 23.8        | 50.0        | 32.3        |
|     |                    | txt&img | 68.2        | 68.6        | <b>67.8</b> | <b>67.7</b> | <b>74.2</b> | <b>75.5</b> | <b>74.5</b> | <b>74.1</b> | <b>70.0</b> | <b>70.1</b> | <b>69.8</b> | <b>69.8</b> | <b>65.7</b> | <b>66.2</b> | <b>65.8</b> | <b>65.5</b> | <b>61.5</b> | <b>62.6</b> | <b>60.5</b> | <b>59.3</b> |
| 235 | Qwen-3             | txt     | 64.4        | <b>69.2</b> | 65.3        | 62.7        | 69.9        | 70.2        | 70.3        | 69.9        | 69.1        | 68.8        | 68.8        | 60.6        | 62.5        | 60.5        | 58.8        | <b>61.5</b> | 61.4        | <b>61.1</b> | <b>61.1</b> |             |
|     |                    | txt&img | 34.4        | 17.2        | 50.0        | 25.6        | 50.6        | 53.1        | 50.0        | 33.9        | 50.4        | 58.5        | 50.0        | 33.6        | 50.0        | 25.0        | 50.0        | 33.3        | 49.6        | 41.5        | 49.9        | 33.3        |
|     |                    | txt     | 34.4        | 17.2        | 50.0        | 25.6        | 50.0        | 35.0        | 50.0        | 33.4        | 50.0        | 50.0        | 50.0        | 33.5        | 50.0        | 25.0        | 49.9        | 33.3        | 50.0        | 41.7        | 50.0        | 33.5        |
| 236 | Gemma-3            | txt&img | 61.9        | 58.0        | 53.2        | 44.3        | 61.9        | 60.5        | 60.3        | 60.3        | 56.5        | 56.4        | 56.3        | 56.2        | 55.7        | 57.5        | 55.6        | 52.7        | 55.9        | 55.9        | 55.9        | 55.8        |
|     |                    | txt     | <b>71.2</b> | 67.7        | <b>66.0</b> | <b>66.5</b> | 57.5        | 59.3        | 57.3        | 54.8        | 58.4        | 59.0        | 58.7        | 58.2        | 58.4        | 61.2        | 58.5        | 55.8        | 59.0        | 58.5        | 57.3        | 56.6        |

238 \* Baseline models. Random guess baseline accuracy is 50%.

241 **Multimodal Models:** For VLMs, the input consists of only only the *Abstract* and *Introduction*  
 242 text, together with the first two figures from each paper. To disentangle the contributions of textual  
 243 and visual information in VLMs, we consider two input configurations: **text+image** and **text-only**.

### 245 4.3 RESULTS

248 We evaluate prediction performance using Accuracy, Macro-Precision (Mac-P), Macro-Recall  
 249 (Mac-R), and F1. Mac-P and Mac-R average class-wise precision and recall, while F1 is the  
 250 harmonic mean of precision and recall. As shown in Table 1 and 2, text-only unimodal models generally  
 251 outperform multimodal text-image models of comparable size. For example, within the Qwen  
 252 family, Qwen3-4B achieves 76.4% accuracy, surpassing multimodal Qwen2.5-VL-3B-Instruct  
 253 at 68.2% and also shows consistently higher Mac-R, Mac-P, and F1.

255 **Text-only models** We evaluate two training strategies: (i) prompt-based fine-tuning on the original  
 256 models, and (ii) CPT followed by prompt-based fine-tuning. As shown in Table 1, CPT yields  
 257 clear improvements for downstream classification. On the aggregated ALL domain, CPT consistently  
 258 outperforms fine-tuning from the original checkpoints (Orig) at the same parameter scale, **improving  
 259 accuracy from 67.3% to 76.4% at 4B and from 69.0% to 78.5% at 8B**. Moreover, CPT is  
 260 more effective at larger scales. For instance, CPT yields a 9.1% improvement at 4B while a 9.5%  
 261 improvement at 8B on the all domain, with consistent increases in LLM, CV, and RL at the 8B scale.

263 **Vision-Language models** We evaluate various VL models, i.e., Qwen2.5-VL-3B-Instruct  
 264 and Gemma-3-4B-it. As shown in Table 2, **Qwen2.5-VL outperforms Gemma-3**, achieving  
 265 68.2% versus 61.9% with text-image input, and consistently higher accuracy across all four sub-  
 266 domains. Second, for Qwen2.5-VL, incorporating text-image input consistently improves perfor-  
 267 mance over text-only input. For example, in the LLM domain it achieves 74.2% compared to 69.9%  
 268 with text-only, and this trend holds across the other three subdomains as well as the aggregated all  
 269 domain. More results on the imbalanced dataset in-domain result in-domain result are provided in  
 App. F.1

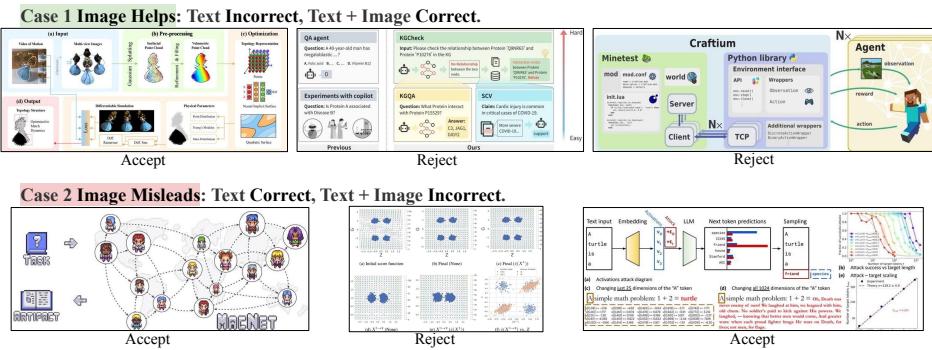
270  
271  
272  
273  
274  
275

Figure 2: Examples of vision–language model predictions on previous submission. Case Group 1: text alone leads to incorrect predictions, while the image provides complementary cues that correct the outcome. Case Group 2: text alone yields the correct answer, but adding the image introduces misleading signals and causes errors.

286

| White Box Features     |                    |                           |                      |                     |                         |                   |  |
|------------------------|--------------------|---------------------------|----------------------|---------------------|-------------------------|-------------------|--|
| Structure              | Visual Content     | Citation Engagement       | Methodological Rigor | Writing Quality     | Novelty & Contribution  | Appendix Material |  |
| <b>total words*</b>    | avg caption length | <b>citations in text*</b> | dataset mentions     | abstract            | novel method claims     | word count        |  |
| <b>total pages*</b>    | image density      | citation density          | metrics mentions     | word count          | comparison studies      | header count      |  |
| <b>header count*</b>   | table density      | baseline mentions         | baseline mentions    | avg sentence length | contribution statements | images count      |  |
| <b>section balance</b> | image count        | statistical tests         | statistical tests    | experiment count    |                         | table count       |  |
| <b>variance*</b>       | equation density   |                           |                      |                     |                         | equation count    |  |
| words/page             | table count        |                           |                      |                     |                         |                   |  |
|                        | equation count     |                           |                      |                     |                         |                   |  |

Figure 3: White box features organized by category and ranked by importance within each group. Features marked with asterisks (\*) represent the top 5 most important features from the Random Forest model

#### 4.4 QUALITATIVE ANALYSIS FOR VL-MODELS

To better understand the role of visual inputs, we qualitatively analyze two outcomes: *Image Helps*, where the model fails with text but succeeds with text–image inputs, and *Image Misleads*, where the addition of images reduces accuracy. Figure 2 illustrates these patterns on prior submissions (Xiong et al., 2025; Lin et al., 2024; Malagón et al.; Qian et al.; Wang et al., 2025; Fort). *Image Helps* (first row) show that schematic figures conveying high-level methodology or motivation, such as pipelines or dataset overviews, help predict the acceptance. In contrast, *Image Misleads* cases often involve detailed result visualizations that are difficult to interpret from figures alone. Additional examples are in App. F.2. We further evaluate models using images as the sole input modality with detail provided in App. F.3.

#### 4.5 STATISTICAL FEATURE ANALYSIS

We train white-box statistical models on manually engineered features to provide an alternative performance baseline and interpretable insights into the structural characteristics that distinguish accepted papers from rejected (Wang et al., 2023).

**Models and Features** We extract 29 quantitative features from each submission PDF across seven categories, as illustrated in Figure 3. A comprehensive list of all features can be found in App. H. These features are then used to train four supervised classifiers, namely Random Forest (Breiman, 2001), Support Vector Machine (Schölkopf et al., 1999), Logistic Regression (Hosmer Jr et al., 2013), and Gradient Boosting (Friedman, 2002).

324

325

326

327

328

329

330

331

332

333

334

335

336

337

| Domain   | Imbalanced Dataset |       |      |      |      | Balanced Dataset |       |      |      |      | Out-of-Distribution Test |       |      |     |      |
|----------|--------------------|-------|------|------|------|------------------|-------|------|------|------|--------------------------|-------|------|-----|------|
|          | Size               | Model | Acc  | F1   | AUC  | Size             | Model | Acc  | F1   | AUC  | Size                     | Model | Acc  | F1  | AUC  |
| LLM      | 3,716              | SVM   | 70.6 | 34.3 | 72.2 | 3,238            | RF    | 66.3 | 67.8 | 71.5 | 2,121                    | GB    | 49.5 | 4.6 | 57.4 |
| CV       | 2,776              | RF    | 71.5 | 49.0 | 72.2 | 3,520            | GB    | 68.3 | 69.6 | 73.3 | 2,230                    | LR    | 51.4 | 2.8 | 55.8 |
| RL       | 1,251              | LR    | 70.1 | 44.4 | 72.1 | 1,526            | GB    | 65.7 | 67.7 | 70.5 | 1,008                    | GB    | 51.5 | 1.1 | 58.3 |
| Theory   | 1,735              | SVM   | 68.5 | 36.9 | 70.2 | 1,974            | SVM   | 63.0 | 64.7 | 71.2 | 1,228                    | GB    | 49.8 | 2.9 | 50.8 |
| Combined | 9,478              | RF    | 77.3 | 60.9 | 83.0 | 10,258           | RF    | 74.2 | 74.9 | 83.1 | 6,587                    | GB    | 53.1 | 2.2 | 61.8 |

GB = Gradient Boosting, RF = Random Forest, LR = Logistic Regression, SVM = Support Vector Machine

Table 3: Performance of statistical models on (i) the imbalanced ICLR 2025 dataset, (ii) balanced domain-specific datasets, and (iii) the Out-of-Distribution Test: models trained on the imbalanced ICLR 2025 data and evaluated on the balanced 50/50 test set.

**Classification Performance** Table 3 reveals distinct performance patterns across all dataset configurations. Random Forest emerges as the best-performing white-box model across both imbalanced and balanced datasets, achieving 77.3% accuracy with an F1-score of 60.9 on imbalanced data, and 74.2% accuracy with a substantially improved F1-score of 74.9 on balanced data. The out-of-domain study demonstrates that *models trained on imbalanced data but evaluated on balanced datasets suffer significant performance degradation*, with Random Forest achieving only 53.1% accuracy and immensely low F1-scores across all models, as the models classified nearly all papers as rejected due to their bias toward the majority class learned from the rejection-heavy imbalanced training data.

The balanced dataset yields on average slightly lower accuracy but significantly higher F1-scores compared to imbalanced, despite having less training data, indicating that *class balance is more critical than dataset size for effective predicting minority research domain*. Across both balanced and imbalanced configurations, combined domain models consistently achieve the best performance compared to individual domains, demonstrating that *cross-domain feature interactions enhance predictive capability*. However, all white-box model results remain significantly below those achieved by fine-tuned LLMs and VLMs, showing the limitations of traditional machine learning approaches in capturing the semantic complexity inherent in peer review decisions.

**Feature Importance Analysis** Random Forest feature importance analysis reveals that structural characteristics dominate acceptance prediction across both dataset configurations, as measured by Gini impurity-based importance scores (Nembrini et al., 2018). As shown in Table 4, the same core structural features consistently appear in the top five most discriminative features across both imbalanced and balanced datasets, suggesting that **paper acceptance favors structure quality rather than domain-specific content**.

Examining the feature rankings reveals several patterns. Content length indicators (*total words*, *total pages*) consistently dominate both configurations, with *total words* ranking first in both cases but showing increased importance (0.079 vs 0.073) in the balanced dataset. Organizational structure features (*header count*, *section balance variance*) maintain high importance across configurations. Most notably, *citations in text* replaces *avg caption length* in the balanced dataset’s top five, suggesting that scholarly engagement becomes more discriminative when class imbalance is addressed.

These patterns indicate that accepted papers consistently tend to be more comprehensive (evidenced by length-based features), better organized (reflected in structural balance metrics), and demonstrate

| (a) Imbalanced Dataset |                          |                |
|------------------------|--------------------------|----------------|
| Rank                   | Feature                  | Imp. Cat.      |
| 1                      | total words              | 0.0739 Struct. |
| 2                      | header count             | 0.0587 Struct. |
| 3                      | total pages              | 0.0575 Struct. |
| 4                      | section balance variance | 0.0492 Struct. |
| 5                      | avg caption length       | 0.0465 Visual  |

| (b) Balanced Dataset |                          |                 |
|----------------------|--------------------------|-----------------|
| Rank                 | Feature                  | Imp. Cat.       |
| 1                    | total words              | 0.0792 Struct.  |
| 2                    | total pages              | 0.0658 Struct.  |
| 3                    | header count             | 0.0569 Struct.  |
| 4                    | section balance variance | 0.0474 Struct.  |
| 5                    | citations in text        | 0.0448 Citation |

Table 4: Top five most discriminative features for paper acceptance prediction from Random Forest analysis across both dataset configurations.

378 stronger scholarly engagement (particularly evident in balanced datasets where citation patterns  
 379 emerge as discriminative). However, the modest importance scores (all  $< 0.08$ ) across both config-  
 380 urations indicate that **no single structural characteristic serves as a strong predictor**, explaining  
 381 why semantic understanding via LLMs significantly outperforms purely structural approaches.  
 382

## 383 5 UTILITY ANALYSIS FOR RECOGNIZING “OBVIOUS” PAPERS

385 In Section 4, we set a default acceptance threshold using  $\text{score}(\text{yes}) > \text{score}(\text{no})$ , though this  
 386 can be adjusted in practical peer-review workflows. In practice, if the model can confidently triage  
 387 “clearly good” and “clearly bad” submissions with minimal errors, it can both reduce reviewer work-  
 388 load and discourage authors from making redundant submission attempts. This section provides a  
 389 *confidence-based utility analysis* to accommodate this need.  
 390

### 391 5.1 CONFIDENCE-BASED STRATIFICATION

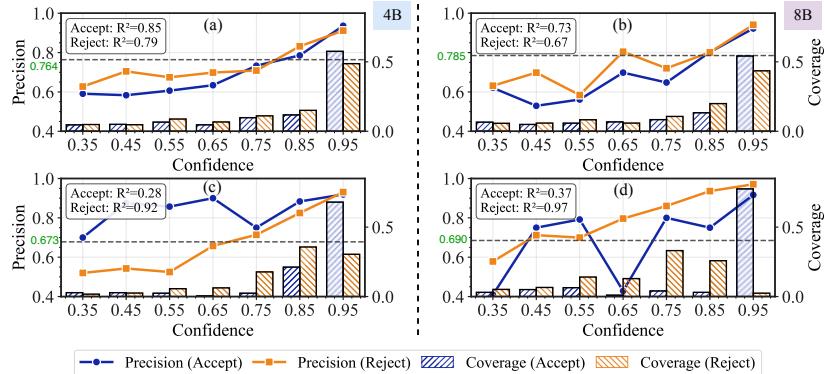
393 **Decision confidence.** At the designated decision slot (cf. §3), let  $l_{\text{yes}}$  and  $l_{\text{no}}$  be the pre-softmax  
 394 logits for the tokens associated with the labels ACCEPT and REJECT, respectively. We define  $p$  as the  
 395 softmax-normalized probability assigned to a class, accept or reject, when considering only these  
 396 two logits. Then we define a scalar *confidence*  $c$  with  $c \approx 0$  indicates indecision ( $\approx 0.5/0.5$ ) and  
 397  $c \approx 1$  indicates near-certainty. Formally,

$$398 c = |p_{\text{yes}} - p_{\text{no}}| = |2p_{\text{yes}} - 1| \in [0, 1] \quad (6)$$

400 **The coverage metric.** Next, we define *coverage* as the fraction of a class’s falling within a given  
 401 confidence bin. Predictions are partitioned into disjoint bins  $B_k$  (e.g.,  $[0.0, 0.1), \dots, [0.9, 1.0]$ ). For  
 402 a set  $\mathcal{S}$  of examples (restricted to a predicted class), the *coverage* of bin  $B_k$  is:

$$403 \text{Cov}(B_k; \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \mathbb{1}\{c_i \in B_k\} = \frac{|\{i \in \mathcal{S} : c_i \in B_k\}|}{|\mathcal{S}|}. \quad (7)$$

405  $c_i \in [0, 1]$  is confidence value of  $i$ . Using these definition together with precision per class, we then  
 406 examine how reliability scales with the model’s self-reported certainty.  
 407



421 Figure 4: Coverage and precision across confidence bins for ACCEPT and REJECT predictions, shown  
 422 for 4B and 8B CPT models. Each panel reports the linear coefficient  $R^2$  of a least-squares fit of  
 423 precision vs. confidence. Panels (a) and (b) correspond to models trained on the balanced dataset,  
 424 while panels (c) and (d) correspond to models trained on the imbalanced.

### 426 5.2 STRATIFIED RESULTS AND OBSERVATIONS

428 We analyze four CPT models: *Qwen3-4B* and *Qwen3-8B* with each trained on the balanced and  
 429 imbalanced datasets, and summarize their behavior in Fig. 4, which plots coverage and precision for  
 430 both predicted classes across confidence bins. Overall, we observe that high-confidence regions ( $c \geq$   
 431 0.9) achieve high precision with substantial coverage, while class imbalance reduces the coverage  
 of confident rejects.

432 **Confidence Concentration and Coverage–Confidence Patterns** Across all models, an average  
 433 of 81.3% of predictions fall within the high-confidence range  $c \in [0.8, 1.0]$ . In the most confident  
 434 interval  $c \in [0.9, 1.0]$  ( $c = 0.95$  in the table), both ACCEPT and REJECT achieve precision above  
 435 91%. This indicates the presence of substantial “obvious tails” that can be triaged with minimal error:  
 436 **when models are highly confident, they are usually correct.** For ACCEPT, coverage increases  
 437 *monotonically* with confidence: it always exceeds 50% and reaches 75.2% for the Qwen3–8B model  
 438 on the imbalanced dataset (Figure 4d), suggesting that most acceptance predictions are made with  
 439 high certainty. In contrast, although REJECT precision improves as  $c$  increases, its coverage is not  
 440 consistently monotonic under imbalanced training, reflecting the relative scarcity of confidently  
 441 identified rejections.

442 We further assess how precision scales with confidence by fitting a least-squares regression sepa-  
 443 rately for ACCEPT and REJECT. The coefficient of determination ( $R^2$ ) (Piepho, 2019), reported in  
 444 the figures, characterizes the degree of linearity in this relationship. The results of  $R^2$  indicate that  
 445 for the minority class, models trained on imbalanced data exhibit markedly poorer certainty than  
 446 their counterparts trained on balanced data. More details are provided in App. I and J.

### 447 5.3 OPENREVIEWER HELPS IDENTIFY “OBVIOUS” GOOD/BAD PAPERS

449 In this section, we examine whether  
 450 OpenReviewer can reliably identify pa-  
 451 pers that are clear accepts or clear rejects. To  
 452 this end, we focus on predictions where the  
 453 model is extremely confident ( $c \in [0.9, 1.0]$ )  
 454 and analyze the corresponding error rates using  
 455 the case of Qwen3–4B model trained on the  
 456 balanced dataset. First, we rank them by their  
 457 confidence scores  $c$  and take the top- $K\%$  mass  
 458 within this band with  $K \in \{1, 3, 5, 7, 9\}$ , i.e.,  
 459 2% step increases. For each slice we report  
 460 per-class *error* ( $= 1 - \text{precision}$ ) and *coverage*.

461 Table 5 reveals encouraging results for work-  
 462 load reduction. When we consider only the top  
 463 1% most-confident predictions, the model cov-  
 464 ers 12.74% of all accept decisions with just  
 465 2.18% error, and 11.36% of all reject decisions  
 466 with 3.07% error. *In practical terms, if the model makes 500 accept predictions, the 64 most-*  
 467 *confident ones would contain fewer than two mistakes.*

468 As we expand to include more confident predictions, we naturally trade some accuracy for greater  
 469 coverage. The top 9% slice covers nearly half of all decisions, 45.02% of accepts and 41.12% of  
 470 rejects, while maintaining reasonably low error rates of 6.03% and 6.06% respectively, illustrating  
 471 the expected precision-coverage trade-off.

472 These results suggest that a confidence-based triage system could substantially reduce reviewer  
 473 workload. By automatically handling the most obvious cases where the model is highly confident,  
 474 conferences could focus human reviewer effort on the more nuanced submissions where expert judg-  
 475 ment is most valuable.

## 477 6 CHALLENGES AND FUTURE WORK

479 This paper presents the first work using LLM to predict AI paper acceptance. Our work opens sev-  
 480 eral promising directions for AI-assisted reviewing, including (i) assessing fairness across subfields,  
 481 (ii) monitoring evolving conference standards, (iii) effectively integrating human-in-the-loop review  
 482 pipelines, and (iv) exploring bias detection to ensure equitable outcomes.

| Top-mass<br>slice | ACCEPT  |            | REJECT  |            |
|-------------------|---------|------------|---------|------------|
|                   | Error ↓ | Coverage ↑ | Error ↓ | Coverage ↑ |
| Top 1.0%          | 2.18    | 12.74      | 3.07    | 11.36      |
| Top 3.0%          | 3.21    | 28.91      | 3.94    | 26.58      |
| Top 5.0%          | 4.12    | 36.84      | 4.83    | 33.71      |
| Top 7.0%          | 4.89    | 40.41      | 5.51    | 39.18      |
| Top 9.0%          | 6.03    | 45.02      | 6.06    | 41.12      |
| All (10%)         | 6.91    | 53.06      | 7.24    | 47.34      |

Table 5: Performance of high-confidence predictions ( $c \in [0.9, 1.0]$ ): error rates and coverage for progressively larger confidence slices. Error rate (%) lower is better ↓; coverage (%) shows the fraction of each class captured in the slice (higher is better ↑).

486 REFERENCES  
487

488 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
489 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*,  
490 2025.

491 Alina Beygelzimer, Yann N Dauphin, Percy Liang, and Jennifer Wortman Vaughan. Has the ma-  
492 chine learning review process become more arbitrary as the field has grown? the neurips 2021  
493 consistency experiment. *arXiv preprint arXiv:2306.03262*, 2023.

494 Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.  
495

496 Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. Revisiting parameter-efficient  
497 tuning: Are we really there yet? *arXiv preprint arXiv:2202.07962*, 2022.

498 Nuo Chen, Moming Duan, Andre Huikai Lin, Qian Wang, Jiaying Wu, and Bingsheng He. Po-  
499 sition: The current ai conference model is unsustainable! diagnosing the crisis of centralized ai  
500 conference. *arXiv preprint arXiv:2508.04586*, 2025.

501 Wuyang Chen, Yanqi Zhou, Nan Du, Yanping Huang, James Laudon, Zhifeng Chen, and Claire Cui.  
502 Lifelong language pretraining with distribution-specialized experts. In *International Conference  
503 on Machine Learning*, pp. 5383–5395. PMLR, 2023.

504 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu,  
505 Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*,  
506 2022.

507 Ahmed Elhady, Eneko Agirre, and Mikel Artetxe. Emergent abilities of large language models under  
508 continued pre-training for language adaptation. In Wanxiang Che, Joyce Nabende, Ekaterina  
509 Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the  
510 Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 32174–32186, Vienna,  
511 Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi:  
512 10.18653/v1/2025.acl-long.1547. URL <https://aclanthology.org/2025.acl-long.1547/>.

513 Yifan Feng, Chengwu Yang, Xingliang Hou, Shaoyi Du, Shihui Ying, Zongze Wu, and Yue  
514 Gao. Beyond graphs: Can large language models comprehend hypergraphs? *arXiv preprint  
515 arXiv:2410.10083*, 2024.

516 Stanislav Fort. Scaling laws for adversarial attacks on language model activations and tokens. In  
517 *The Thirteenth International Conference on Learning Representations*.

518 George Fragiadakis, Christos Diou, George Kousiouris, and Mara Nikolaidou. Evaluating human-  
519 ai collaboration: A review and methodological framework. *CoRR*, abs/2407.19098, 2024. URL  
520 <https://doi.org/10.48550/arXiv.2407.19098>.

521 Jerome H Friedman. Stochastic gradient boosting. *Computational statistics & data analysis*, 38(4):  
522 367–378, 2002.

523 Krishna Garg, Firoz Shaik, Sambaran Bandyopadhyay, and Cornelia Caragea. Let’s use chatgpt to  
524 write our paper! benchmarking llms to write the introduction of a research paper, 2025. URL  
525 <https://arxiv.org/abs/2508.14273>.

526 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,  
527 and Noah A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. In  
528 Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th  
529 Annual Meeting of the Association for Computational Linguistics*, pp. 8342–8360, Online, July  
530 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL  
531 <https://aclanthology.org/2020.acl-main.740/>.

532 David W Hosmer Jr, Stanley Lemeshow, and Rodney X Sturdivant. *Applied logistic regression*. John  
533 Wiley & Sons, 2013.

534 Loi Duc Huynh, Tianshi Che, Zijie Zhang, Yang Zhou, Ruoming Jin, and Dejing Dou. k-odd one  
535 clear (k-ooc), a novel gpu kernel that improves quantization accuracy and speed of gptq algorithm.  
536 2025.

540 Hyungkyu Kang and Min-hwan Oh. Adversarial policy optimization for offline preference-based  
 541 reinforcement learning. *arXiv preprint arXiv:2503.05306*, 2025.

542

543 Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. Continual pre-  
 544 training of language models. In *The Eleventh International Conference on Learning Representa-  
 545 tions*, 2023. URL [https://openreview.net/forum?id=m\\_GDIItaI3o](https://openreview.net/forum?id=m_GDIItaI3o).

546 Jaeho Kim, Yunseok Lee, and Seulki Lee. Position: The AI conference peer review crisis demands  
 547 author feedback and reviewer rewards. In *Forty-second International Conference on Machine  
 548 Learning Position Paper Track*, 2025. URL <https://openreview.net/forum?id=l8QemUZaIA>.

549

550 Kayvan Kousha and Mike Thelwall. Artificial intelligence to support publishing and peer review: A  
 551 summary and review. *Learned Publishing*, 37(1):4–12, 2024.

552

553 Giuseppe Russo Latona, Manoel Horta Ribeiro, Tim R. Davidson, Veniamin Veselovsky, and Robert  
 554 West. The ai review lottery: Widespread ai-assisted peer reviews boost paper scores and accep-  
 555 tance rates, 2024. URL <https://arxiv.org/abs/2405.02150>.

556

557 Neil D. Lawrence. The neurips experiment. [https://inverseprobability.com/talks/notes/  
 the-neurips-experiment-snsf.html](https://inverseprobability.com/talks/notes/the-neurips-experiment-snsf.html), 2022. Blog post.

558

559 Xinna Lin, Siqi Ma, Junjie Shan, Xiaojing Zhang, Shell Xu Hu, Tiannan Guo, Stan Z Li, and  
 560 Kaicheng Yu. Biokgbench: A knowledge graph checking benchmark of ai agent for biomed-  
 561 ical science. *arXiv preprint arXiv:2407.00466*, 2024.

562

563 Mikel Malagón, Josu Ceberio, and Jose A Lozano. Craftium: Bridging flexibility and efficiency  
 564 for rich 3d single-and multi-agent environments. In *Forty-second International Conference on  
 Machine Learning*.

565

566 Jiacheng Miao, Joe R. Davis, Jonathan K. Pritchard, and James Zou. Paper2agent: Reimagining  
 567 research papers as interactive and reliable ai agents, 2025. URL <https://arxiv.org/abs/2509.06917>.

568

569 Stefano Nembrini, Inke R König, and Marvin N Wright. The revival of the gini importance? *Bioin-  
 formatics*, 34(21):3711–3718, 2018.

570

571 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
 572 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
 573 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

574

575 Pat Pataranutaporn, Nattavudh Powdthavee, Chayapatr Achiwaranguprok, and Pattie Maes. Can ai  
 576 solve the peer review crisis? a large scale cross model experiment of llms' performance and biases  
 577 in evaluating over 1000 economics papers. *arXiv preprint arXiv:2502.00070*, 2025.

578

579 Hans-Peter Piepho. A coefficient of determination (r2) for generalized linear mixed models. *Bio-  
 metrical journal*, 61(4):860–872, 2019.

580

581 Chen Qian, Zihao Xie, YiFei Wang, Wei Liu, Kunlun Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du,  
 582 Weize Chen, Cheng Yang, et al. Scaling large language model-based multi-agent collaboration.  
 583 In *The Thirteenth International Conference on Learning Representations*.

584

585 Haolin Ruan, Shaohang Xu, Zhi Chen, Yining Dong, and Chin Pang Ho. Target-oriented soft-robust  
 586 inverse reinforcement learning. 2025.

587

588 Qian Ruan, Ilia Kuznetsov, and Iryna Gurevych. Are large language models good classifiers? a study  
 589 on edit intent classification in scientific document revisions. In Yaser Al-Onaizan, Mohit Bansal,  
 590 and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural  
 Language Processing*, pp. 15049–15067, Miami, Florida, USA, November 2024. Association for  
 591 Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.839. URL <https://aclanthology.org/2024.emnlp-main.839/>.

592

593 Rylan Schaeffer, Joshua Kazdan, Yegor Denisov-Blanch, Brando Miranda, Matthias Gerstgrasser,  
 594 Susan Zhang, Andreas Haupt, Isha Gupta, Elyas Obbad, Jesse Dodge, et al. Position: Ma-  
 595 chine learning conferences should establish a "refutations and critiques" track. *arXiv preprint  
 arXiv:2506.19882*, 2025.

594 Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification and  
 595 natural language inference. In Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (eds.), *Proceed-  
 596 ings of the 16th Conference of the European Chapter of the Association for Computational Lin-  
 597 guistics: Main Volume*, pp. 255–269, Online, April 2021. Association for Computational Linguis-  
 598 tics. doi: 10.18653/v1/2021.eacl-main.20. URL <https://aclanthology.org/2021.eacl-main.20/>.

599 Johannes Schneider. Generative to agentic ai: Survey, conceptualization, and challenges. *arXiv  
 600 preprint arXiv:2504.18875*, 2025.

602 Bernhard Schölkopf, Robert C Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support  
 603 vector method for novelty detection. *Advances in neural information processing systems*, 12,  
 604 1999.

605 Zhengxiang Shi and Aldo Lipani. Don’t stop pretraining? make prompt-based fine-tuning powerful  
 606 learner. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL  
 607 <https://openreview.net/forum?id=s7xWeJQACI>.

609 Hyungyu Shin, Jingyu Tang, Yoonjoo Lee, Nayoung Kim, Hyunseung Lim, Ji Yong Cho, Hwajung  
 610 Hong, Moontae Lee, and Juho Kim. Mind the blind spots: A focus-level evaluation framework  
 611 for llm reviews. *arXiv preprint arXiv:2502.17086*, 2025.

612 Yixiao Song, Yekyung Kim, and Mohit Iyyer. Veriscore: Evaluating the factuality of verifiable  
 613 claims in long-form text generation, 2024. URL <https://arxiv.org/abs/2406.19276>.

615 Purin Sukpanichnant, Anna Rapberger, and Francesca Toni. Peerarg: Argumentative peer review  
 616 with llms. *arXiv preprint arXiv:2409.16813*, 2024.

617 Pawin Taechoyotin and Daniel Acuna. Remor: Automated peer review generation with llm reasoning  
 618 and multi-objective reinforcement learning, 2025. URL <https://arxiv.org/abs/2505.11718>.

620 Derek Tam, Rakesh R. Menon, Mohit Bansal, Shashank Srivastava, and Colin Raffel. Improving  
 621 and simplifying pattern exploiting training. In Marie-Francine Moens, Xuanjing Huang, Lucia  
 622 Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods  
 623 in Natural Language Processing*, pp. 4980–4991, Online and Punta Cana, Dominican Republic,  
 624 November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.  
 625 407. URL <https://aclanthology.org/2021.emnlp-main.407/>.

626 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,  
 627 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical  
 628 report. *arXiv preprint arXiv:2503.19786*, 2025.

630 Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu,  
 631 Yuan Qu, Fukai Shang, Bo Zhang, Liqun Wei, Zhihao Sui, Wei Li, Botian Shi, Yu Qiao, Dahua  
 632 Lin, and Conghui He. Mineru: An open-source solution for precise document content extraction,  
 633 2024. URL <https://arxiv.org/abs/2409.18839>.

634 Gang Wang, Qi Peng, Yanfeng Zhang, and Mingyang Zhang. What have we learned from openre-  
 635 view? *World Wide Web*, 26(2):683–708, 2023.

636 Liming Wang, Muhammad Jehanzeb Mirza, Yishu Gong, Yuan Gong, Jiaqi Zhang, Brian H Tracey,  
 637 Katerina Placek, Marco Vilela, and James R Glass. Can diffusion models disentangle? a theore-  
 638 tical perspective. *arXiv preprint arXiv:2504.00220*, 2025.

640 Xiaoyu Xiong, Changyu Hu, Chunru Lin, Pingchuan Ma, Chuang Gan, and Tao Du. Topogaussian:  
 641 Inferring internal topology structures from visual clues. *arXiv preprint arXiv:2503.12343*, 2025.

642 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
 643 Chang Gao, Chengan Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint  
 644 arXiv:2505.09388*, 2025.

646 Jing Yang. Position: The artificial intelligence and machine learning community should adopt a more  
 647 transparent and regulated peer review process. In *Forty-second International Conference on Ma-  
 chine Learning Position Paper Track*, 2025. URL <https://openreview.net/forum?id=gnyqRarPzW>.

648 Yuqing Yang and Robin Jia. When do llms admit their mistakes? understanding the role of model  
649 belief in retraction. *arXiv preprint arXiv:2505.16170*, 2025.  
650

651 Rui Ye, Xianghe Pang, Jingyi Chai, Jiaao Chen, Zhenfei Yin, Zhen Xiang, Xiaowen Dong, Jing  
652 Shao, and Siheng Chen. Are we there yet? revealing the risks of utilizing large language models  
653 in scholarly peer review. *arXiv preprint arXiv:2412.01708*, 2024.

654 Chen Bo Calvin Zhang, Zhang-Wei Hong, Aldo Pacchiano, and Pulkit Agrawal. Orso: Accelerating  
655 reward design via online reward selection and policy optimization. *arXiv preprint  
656 arXiv:2410.13837*, 2024.

657

658 Qihang Zhao and Xiaodong Feng. Utilizing citation network structure to predict paper citation  
659 counts: A deep learning approach. *Journal of Informetrics*, 16(1):101235, 2022.

660

661 Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate Before Use: Im-  
662 proving Few-Shot Performance of Language Models. In Marina Meila and Tong Zhang (eds.),  
663 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-  
664 ceedings of Machine Learning Research*, pp. 12697–12706. PMLR, 18–24 Jul 2021. URL  
<https://proceedings.mlr.press/v139/zhao21c.html>.

665

666 Wuqiang Zheng, Yiyan Xu, Xinyu Lin, Chongming Gao, Wenjie Wang, and Fuli Feng. Navigating  
667 through paper flood: Advancing llm-based paper evaluation through domain-aware retrieval and  
668 latent reasoning, 2025. URL <https://arxiv.org/abs/2508.05129>.

669

670 Yuan Zhou, Peng Zhang, Mengya Song, Alice Zheng, Yiwen Lu, Zhiheng Liu, Yong Chen, and  
671 Zhaohan Xi. Zodiac: A cardiologist-level llm framework for multi-agent diagnostics. *arXiv  
672 preprint arXiv:2410.02026*, 2024.

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

## 702 A REPRODUCIBILITY STATEMENT

704 We take several steps to enable full replication of our results.

706 **Data.** We use ICLR 2024–2025 submissions and final decisions obtained via the OpenReview  
 707 API–V2 under CC BY 4.0; our crawl, de-identification, and parsing pipeline and the rules for do-  
 708 main labeling and class balancing are described in App. C and summarized in Table 6.

710 **Models & training.** Exact model checkpoints and modalities appear in Sec. 4.2. The prompt tem-  
 711 plate and label verbalizers are given in App. D; the continual pre-training corpus construction, pack-  
 712 ing block size, and optimization details are in Sec. 3.1 and App. E. All experiments were run on a  
 713 two NVIDIA A100 80GB GPUs; precision and optimizer choices match App. E.4.

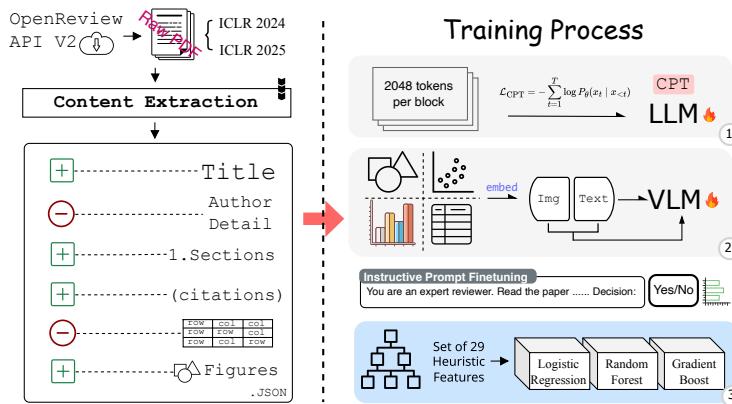
714 **Baselines & features.** The 29 engineered features and model choices are documented in Sec. 4.5  
 715 and App. H. Evaluation. We report Accuracy, Macro-Precision, Macro-Recall and F1 with results in  
 716 Tables 1–2, 8–10; the out-of-distribution tests is detailed in Apps. G. The confidence-stratified utility  
 717 analysis includes formulas and binning definitions in Sec. 5, Apps. I–J.

718 Upon publication, we will release our complete codebase and processed datasets, with rebuild  
 719 scripts, to facilitate replication and extension of this work.

## 722 B USE OF LARGE LANGUAGE MODELS

724 In this work, we used large language models (LLMs) for two distinct purposes. First, we employed  
 725 OpenAI’s ChatGPT (GPT-5) exclusively for grammar correction and improving the fluency of the  
 726 manuscript. Second, we evaluated ChatGPT’s performance on our prediction task as part of the ex-  
 727 perimental analysis. Significantly, the model did not contribute to the research design, methodology,  
 728 or interpretation of results; its role in writing was strictly limited to polishing sentence structure and  
 729 enhancing readability. All technical contributions remain the sole work of the authors.

## 731 C DATA COLLECTION AND PRE-PROCESSING



746 Figure 5: Data collection and preprocessing workflow and training pipeline.

749 We collect all ICLR 2025 and 2024 submissions and their corresponding final decisions (*accepted* or  
 750 *rejected*) via the OpenReview API–V2. All acquired data complies with the Creative Commons  
 751 Attribution 4.0 International (CC BY 4.0) license. The papers were further partitioned into four  
 752 main subfields based on title keywords: Large Language Models (LLM), Computer Vision (CV),  
 753 Reinforcement Learning (RL), and Theoretical (Theory). Papers that do not fall into these categories  
 754 are left for future discussion. Summary counts for each subfield are reported in Table 6.

755 From the collected submissions, we construct two distinct datasets for our analysis: a complete  
 ICLR2025 dataset as well as balanced domain-specific datasets by combining papers from both

756 ICLR 2024 and 2025 to ensure equal representation of accepted and rejected papers, addressing  
 757 potential class imbalance issues that could bias our analysis.  
 758

759 We employ MINERU (Wang et al.,  
 760 2024), an OCR-based tool, to ex-  
 761 tract structured content from the  
 762 collected PDFs. As shown in Fig-  
 763 ure 5.

764 MINERU processes each document  
 765 by separating text, images, ta-  
 766 bles, and equations, and generates  
 767 a structured JSON representation.  
 768 From this output, we retain only  
 769 elements labeled as figures, tables,  
 770 or equations, and restricted text ex-  
 771 traction to the title, abstract, and in-  
 772 troduction sections for use in our prediction model.<sup>4</sup> The final representation for each paper con-  
 773 sisted of clean text files for the targeted sections, alongside organized visual elements paired with  
 774 their original captions.

## 775 D PROMPT TEMPLATE

776 We design an instructive prompt template that presents the paper’s features within a natural-language  
 777 query and guides the model to generate a decision token corresponding to one of the two target  
 778 classes: *accept* or *reject*.

### 779 Input Example

#### 780 Template $\mathcal{T}(x)$ :

781 You are an expert reviewer. Read the paper content and decide if it  
 782 should be accepted.

783 Paper content:  $\langle x \rangle$

784 Decision:

785 Given a paper input instance  $x$  and prompt template  $\mathcal{T}(x)$ , the model defines a conditional proba-  
 786 bility over the *label verbalizer* Tam et al. (2021).

## 787 E CONTINUAL PRE-TRAINING

### 788 E.1 MOTIVATION

789 Continual pre-training (CPT) adapts a strong general-purpose language model to the peer-review do-  
 790 main by further training on large-scale, unlabeled scientific corpora. Unlike supervised fine-tuning,  
 791 CPT retains the original causal language modeling objective, thereby aligning the model’s genera-  
 792 tive priors with the linguistic and structural regularities of academic manuscripts. This is particularly  
 793 important in OpenReviewer, where downstream tasks rely on prompt-conditioned generation rather  
 794 than explicit classification heads. This section will describe the training detail used for CPT.

### 801 E.2 INPUT SETTING

802 We construct the CPT corpus by aggregating unlabeled texts from academic paper PDFs pro-  
 803 cessed with MinerU. Each document is concatenated with an EOS separator, tokenized using the  
 804 model’s native tokenizer, and packed into fixed-length blocks of size  $B$  (default  $B = 2048$ ). This  
 805 block-packing strategy eliminates under-filled sequences and ensures efficient utilization of training  
 806 batches. The input IDs and labels are identical, enabling pure causal next-token prediction.

807  
 808  
 809 <sup>4</sup>Manual spot-checking confirmed high quality of the extracted content.

This preprocessing not only exposes the model to scientific writing styles, rhetorical markers, and citation format, etc. but also reduces the domain gap between generic pre-training corpora and the specialized peer-review domain.

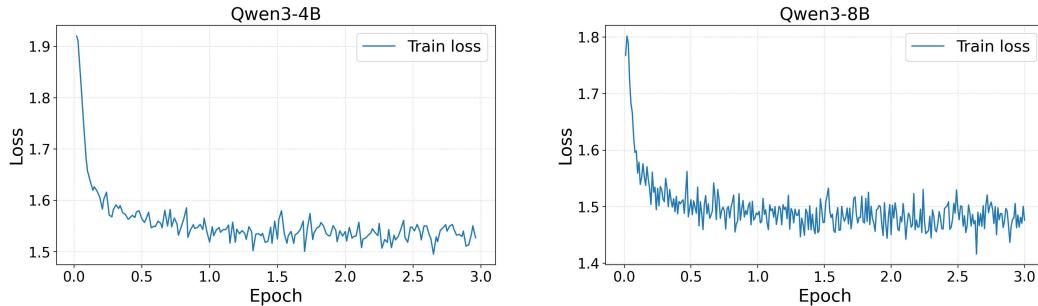


Figure 6: Continual pre-training loss results of Qwen3-4B and Qwen3-8B on the balanced dataset corpus.

### E.3 RESULT

During continual pre-training on the balanced corpus, the training loss decreases steadily for both Qwen3-4B and Qwen3-8B, indicating stable optimization. The 8B model converges slightly faster and to a lower final loss than the 4B model, consistent with its larger capacity. We observed no signs of divergence or instability across the three epochs, suggesting CPT effectively adapts the base models to scientific writing before downstream fine-tuning.

### E.4 HYPERPARAMETER SETTINGS

To maintain training stability, we adopt AdamW optimization with cosine learning rate decay, gradient checkpointing, and norm clipping. CPT is performed *prior* to prompt-based fine-tuning so that the updated parameters  $\theta$  encode domain knowledge without introducing task-specific biases. Training was conducted on a single NVIDIA A100 80GB GPU.

| Hyperparameter          | Value  |
|-------------------------|--|
| Backbone Model          | Qwen3-4B(8B)   |
| Sequence Length ( $B$ ) | 2048   |
| Batch Size (per device) | 2  |
| Gradient Accumulation   | 8 (effective batch = $2 \times 8 \times \text{GPUs}$ ) |
| Epochs                  | 3  |
| Learning Rate           | $1(2) \times 10^{-5}$                                  |
| Warmup Ratio            | 0.1  |
| Weight Decay            | 0.1  |
| Optimizer               | AdamW  |
| Scheduler               | Cosine decay   |
| Precision               | bfloat16 (default)                                     |
| Attention Backend       | SDPA (FlashAttention-2 optional)                       |
| Gradient Checkpointing  | Enabled  |
| Max Grad Norm           | 1.0  |

Table 7: Hyperparameter settings for continual pre-training in OpenReview.

## 864 F ADDITIONAL RESULTS ON VL-MODEL

### 866 F.1 RESULTS ON IMBALANCE DATASET (IN-DISTRIBUTION)

868 Table 8 shows results on the imbalanced dataset. The models exhibit base-rate and threshold bias:  
869 minimizing loss encourages predicting the majority class. The prediction becomes more sensitive to  
870 textreject patterns while under-covering the minority.

871 *Imbalanced In-Distribution Test*

| 872 SUB-<br>873 DOMAIN | 874 LLM |           |           |        | 875 CV  |           |           |        | 876 RL  |           |           |        | 877 THEORY |           |           |        | 878 ALL |           |           |        |
|------------------------|---------|-----------|-----------|--------|---------|-----------|-----------|--------|---------|-----------|-----------|--------|------------|-----------|-----------|--------|---------|-----------|-----------|--------|
|                        | 879 ACC | 879 MAC-P | 879 MAC-R | 879 F1 | 879 ACC | 879 MAC-P | 879 MAC-R | 879 F1 | 879 ACC | 879 MAC-P | 879 MAC-R | 879 F1 | 879 ACC    | 879 MAC-P | 879 MAC-R | 879 F1 | 879 ACC | 879 MAC-P | 879 MAC-R | 879 F1 |
| txt&img                | 35.7    | 17.9      | 50.0      | 26.3   | 35.6    | 17.8      | 50.0      | 26.3   | 32.9    | 16.5      | 50.0      | 24.8   | 33.1       | 16.6      | 50.0      | 24.9   | 34.8    | 17.4      | 50.0      | 25.8   |
| txt*                   | 35.7    | 17.9      | 50.0      | 26.3   | 35.7    | 17.8      | 50.0      | 26.3   | 32.9    | 16.5      | 50.0      | 24.8   | 33.1       | 16.6      | 50.0      | 24.9   | 34.8    | 17.4      | 50.0      | 25.8   |
| txt&king               | 73.2    | 70.9      | 71.2      | 71.0   | 63.1    | 63.5      | 64.0      | 59.9   | 59.6    | 55.4      | 55.7      | 55.5   | 68.5       | 63.1      | 59.3      | 59.4   | 74.8    | 65.8      | 76.9      | 66.4   |
| txt                    | 74.2    | 68.2      | 69.5      | 68.2   | 69.5    | 68.1      | 69.5      | 68.2   | 66.3    | 57.9      | 53.9      | 51.9   | 67.2       | 62.5      | 62.0      | 62.2   | 75.5    | 71.9      | 69.4      | 67.1   |
| Gemma-3-4B             | 35.6    | 17.8      | 49.8      | 26.2   | 35.6    | 17.8      | 50.0      | 26.3   | 32.5    | 16.3      | 49.4      | 24.5   | 33.1       | 16.6      | 50.0      | 24.9   | 34.4    | 17.2      | 50.0      | 25.6   |
| txt&img                | 34.0    | 17.0      | 50.0      | 25.4   | 34.6    | 17.3      | 50.0      | 25.7   | 34.8    | 17.4      | 50.0      | 25.8   | 37.1       | 18.6      | 50.0      | 27.1   | 34.4    | 17.2      | 50.0      | 25.6   |
| txt&king               | 64.3    | 32.2      | 50.0      | 39.1   | 63.8    | 54.4      | 50.9      | 43.7   | 67.9    | 83.8      | 51.3      | 42.8   | 66.0       | 55.9      | 55.0      | 45.6   | 76.8    | 91.5      | 35.7      | 51.4   |
| txt                    | 70.7    | 71.4      | 59.1      | 57.8   | 73.7    | 72.2      | 66.5      | 67.4   | 69.6    | 66.3      | 65.7      | 65.9   | 72.5       | 72.5      | 66.1      | 66.7   | 72.4    | 70.2      | 64.4      | 65.1   |

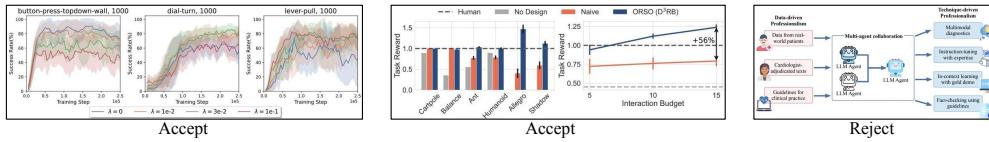
880 \* Baseline models.

881 Table 8: Accuracy performance (%) of Qwen2.5-VL-3B-Instruct and Gemma-3-4B-it  
882 across four broad domains on imbalanced dataset. Mac-P and Mac-R denote Macro Precision and  
883 Macro Recall, respectively.

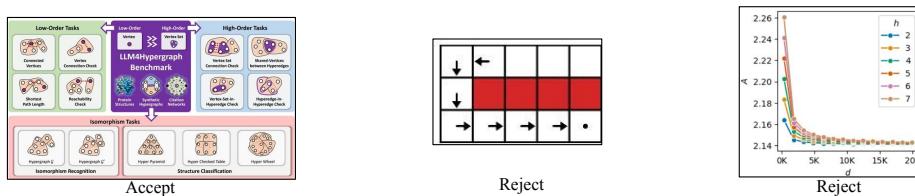
## 885 F.2 MORE QUALITY ANALYSIS

888 Figure 7 shows additional examples of these two patterns from prior submissions (Kang & Oh, 2025;  
889 Zhang et al., 2024; Zhou et al., 2024; Feng et al., 2024; Ruan et al., 2025; Huynh et al., 2025). Most  
890 image-help cases are teaser images, which usually contain clear text and visual cues that support the  
891 model’s judgment. In contrast, many image-mislead cases come from result analysis figures rather  
892 than teaser images, and thus contain little or no explicit textual guidance, making them harder for  
893 the model to interpret correctly.

### 894 Case 1 Image Helps: Text Incorrect, Text + Image Correct.



### 900 Case 2 Image Misleads: Text Correct, Text + Image Incorrect.



901 Figure 7: More examples of vision–language model predictions on previous submission.

## 910 F.3 IMAGE-ONLY FOR PREDICTION

912 We further evaluate models using images as the sole input modality. First, we employ  
913 DINO-v2 (Oquab et al., 2023)<sup>5</sup> as a classifier, where the inputs are the first two main figures from  
914 each paper. This setting yields an accuracy of 39.5% and an F1 score of 49.8%. In addition, we  
915 experiment with converting the first two pages of each PDF into images and training Qwen-VL  
916 with these image-only inputs. However, the performance in this setting remains close to that of the  
917 untrained baseline.

918 <sup>5</sup><https://huggingface.co/facebook/dinov2-base>

## 918 G ABLATION STUDIES

920 We train on the imbalanced ICLR-2025 split and evaluate on a balanced 50/50 test to probe robustness. Across sizes, OOD accuracy hovers around 67–69%, with noticeable drops in macro-recall/F1  
 921 versus in-distribution, reflecting a reject-majority bias learned from imbalanced training. CPT offers  
 922 modest, inconsistent gains (slightly higher macro-recall/F1 in some domains) but does not eliminate  
 923 the bias; larger models (8B) do not guarantee better OOD generalization than 4B. Overall, results  
 924 show that class balance during training matters more than scale, and that simple fine-tuning on im-  
 925 balanced data leads to systematic under-coverage of ACCEPT, suggesting the need for rebalancing,  
 926 threshold calibration, or post-hoc confidence conditioning for reliable deployment.

| LLMs Out-of-Distribution Test |     |       |       |      |      |       |       |      |      |       |       |      |      |       |       |      |        |       |       |      |      |
|-------------------------------|-----|-------|-------|------|------|-------|-------|------|------|-------|-------|------|------|-------|-------|------|--------|-------|-------|------|------|
| SUB-DOMAIN                    | ALL |       |       |      | LLM  |       |       |      | CV   |       |       |      | RL   |       |       |      | THEORY |       |       |      |      |
|                               | ACC | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   | ACC    | MAC-P | MAC-R | F1   |      |
| CPT                           | 4B  | 67.8  | 72.5  | 67.0 | 68.4 | 65.6  | 71.4  | 64.9 | 65.6 | 68.2  | 70.8  | 69.9 | 71.6 | 66.4  | 68.1  | 53.1 | 46.7   | 54.8  | 54.3  | 51.9 | 44.2 |
|                               | 8B  | 68.5  | 68.2  | 70.5 | 70.2 | 67.9  | 68.3  | 70.3 | 69.1 | 68.7  | 72.9  | 70.1 | 71.8 | 65.6  | 62.2  | 52.4 | 54.3   | 53.7  | 48.8  | 51.2 | 44.2 |
| Orig                          | 4B  | 68.6  | 70.9  | 69.0 | 70.8 | 65.9  | 67.4  | 66.8 | 68.0 | 67.5  | 71.5  | 68.4 | 69.8 | 62.8  | 67.4  | 62.7 | 62.3   | 59.5  | 65.5  | 58.5 | 55.8 |
|                               | 8B  | 67.0  | 70.5  | 66.1 | 67.4 | 66.6  | 67.2  | 68.3 | 66.7 | 67.1  | 70.8  | 67.9 | 69.3 | 61.6  | 64.8  | 59.7 | 57.2   | 52.3  | 51.3  | 50.4 | 44.6 |

935 Table 9: Ablation results on the imbalanced ICLR 2025 dataset. Models are trained with the original  
 936 accept/reject ratio (31.7% / 68.3%) and evaluated on the balanced 50/50 Out-of-Distribution test set.

| VLMs Out-of-Distribution Test |         |       |       |      |      |       |       |      |      |       |       |      |        |       |       |      |      |       |       |      |      |
|-------------------------------|---------|-------|-------|------|------|-------|-------|------|------|-------|-------|------|--------|-------|-------|------|------|-------|-------|------|------|
| SUB-DOMAIN                    | ALL     |       |       |      | CV   |       |       |      | RL   |       |       |      | THEORY |       |       |      | LLM  |       |       |      |      |
|                               | ACC     | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   | ACC    | MAC-P | MAC-R | F1   | ACC  | MAC-P | MAC-R | F1   |      |
| Qwen                          | txt&img | 75.2  | 89.8  | 55.4 | 68.5 | 66.7  | 63.7  | 75.6 | 69.2 | 59.9  | 67.8  | 45.5 | 54.4   | 59.1  | 73.8  | 44.5 | 47.0 | 65.0  | 83.2  | 45.1 | 58.5 |
|                               | txt     | 76.2  | 89.4  | 58.0 | 70.3 | 75.2  | 80.7  | 65.7 | 72.4 | 55.7  | 88.9  | 18.2 | 30.2   | 66.1  | 82.0  | 45.6 | 58.6 | 59.3  | 87.5  | 29.9 | 44.6 |
| Gemma                         | txt&img | 70.2  | 80.2  | 69.5 | 67.0 | 53.2  | 73.7  | 8.1  | 14.7 | 50.3  | 58.1  | 20.5 | 30.3   | 60.2  | 58.7  | 82.2 | 68.5 | 62.0  | 64.4  | 68.3 | 66.3 |
|                               | txt     | 76.2  | 81.2  | 75.7 | 75.0 | 53.7  | 58.0  | 23.3 | 33.2 | 57.5  | 63.5  | 45.5 | 53.0   | 54.4  | 73.1  | 21.1 | 32.8 | 55.7  | 75.4  | 28.1 | 40.9 |

945 Table 10: Accuracy performance (%) of Qwen2.5–VL–3B–Instruct and Gemma–3–4B–it  
 946 across four broad domains under the Out-of-Distribution Test setting

## 949 H COMPLETE TABLE OF WHITE BOX FEATURES

951 The complete 29 white-box features importance are reported in Table 11.

## 953 I $R_2$ DERIVATION

955 We quantify how reliability scales with certainty by *separately for each predicted class* (ACCEPT,  
 956 REJECT) fitting an ordinary least squares line to bin-level precision vs. confidence (using the *filtered*  
 957 *bin midpoints* as  $x$ ):

959 Given paired points  $\{(x_i, y_i)\}$  where  $x_i$  is the confidence-bin midpoint and  $y_i$  the corresponding  
 960 precision:

$$961 \text{Fit: } y = mx + b, \quad (8)$$

$$963 \text{Residual sum of squares: } SS_{\text{res}} = \sum_i (y_i - \hat{y}_i)^2, \quad (9)$$

$$965 \text{Total sum of squares: } SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2, \quad (10)$$

$$967 \text{Coefficient of determination: } R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}. \quad (11)$$

## 970 Interpretation.

- 971 •  $R^2 = 1$ : perfect linear fit.

| 972                            | Feature | Bal    | 973 | Imb |
|--------------------------------|---------|--------|-----|-----|
| <b>974 Structure</b>           |         |        |     |     |
| 975 total words                | 0.0792  | 0.0739 | 976 |     |
| 977 total pages                | 0.0658  | 0.0575 | 978 |     |
| header count                   | 0.0569  | 0.0587 |     |     |
| section balance variance       | 0.0474  | 0.0492 |     |     |
| words/page                     | 0.0446  | 0.0437 |     |     |
| <b>979 Visual Content</b>      |         |        |     |     |
| 980 avg caption length         | 0.0439  | 0.0465 | 981 |     |
| image density                  | 0.0374  | 0.0378 | 982 |     |
| table density                  | 0.0363  | 0.0370 | 983 |     |
| image count                    | 0.0362  | 0.0332 | 984 |     |
| equation density               | 0.0360  | 0.0372 | 985 |     |
| table count                    | 0.0355  | 0.0356 |     |     |
| equation count                 | 0.0315  | 0.0341 |     |     |
| <b>986 Citation Engagement</b> |         |        |     |     |
| 987 citations in text          | 0.0448  | 0.0403 | 988 |     |
| citation density               | 0.0408  | 0.0395 |     |     |

(a) Structural, visual, and citation features

| 990 Feature                           | Bal    | 991 Imb |
|---------------------------------------|--------|---------|
| <b>992 Methodological Rigor</b>       |        |         |
| dataset mentions                      | 0.0389 | 0.0374  |
| metrics mentions                      | 0.0355 | 0.0356  |
| baseline mentions                     | 0.0241 | 0.0240  |
| statistical tests                     | 0.0061 | 0.0058  |
| experiment count                      | 0.0030 | 0.0031  |
| <b>993 Writing Quality</b>            |        |         |
| abstract word count                   | 0.0430 | 0.0446  |
| avg sentence length                   | 0.0423 | 0.0432  |
| <b>994 Novelty &amp; Contribution</b> |        |         |
| novel method claims                   | 0.0346 | 0.0365  |
| comparison studies                    | 0.0327 | 0.0356  |
| contribution statements               | 0.0318 | 0.0314  |
| <b>995 App. Material</b>              |        |         |
| word count (appendix)                 | 0.0188 | 0.0203  |
| header count (appendix)               | 0.0179 | 0.0177  |
| images count (appendix)               | 0.0125 | 0.0134  |
| table count (appendix)                | 0.0121 | 0.0161  |
| equation count (appendix)             | 0.0088 | 0.0095  |

(b) Methodological, writing, novelty, and appendix features

*Bal* = balanced dataset importance; *Imb* = imbalanced dataset importance.

Table 11: Feature importance across balanced (Bal) and imbalanced (Imb) datasets using a Random Forest. Values are normalized importances.

- $R^2 = 0$ : no better than predicting the mean  $\bar{y}$ .
- $R^2 < 0$ : worse than predicting the mean.
- $R^2 = \text{NaN}$ : not enough points, constant  $x$ , or zero variance in  $y$  ( $SS_{\text{tot}} = 0$ ).

## J LINEAR PRECISION-CONFIDENCE RELATIONSHIP

We further quantify how reliability scales with certainty by fitting, separately for ACCEPT and REJECT, an least squares model of precision against confidence, and we report the linear coefficient of determination  $R^2$  in the figures to characterize the strength of the linearity.

On the **balanced** dataset (Fig. 4a, 4b), two classes exhibit similar  $R^2$  values, indicating that increases in confidence translate into nearly equivalent gains in precision for both ACCEPT and REJECT. Moreover, the Qwen3-4B model exhibits a stronger linear relationship than the Qwen3-8B model on this dataset, with the highest fit  $R^2 = 0.85$ .

On the **imbalanced** dataset (Fig. 4c, 4d), by contrast, the precision–confidence relationship diverges across classes: the minority class (ACCEPT) typically shows a lower  $R^2$ , reflecting weaker separability than under balanced training.