

MMReview: A Multidisciplinary and Multimodal Benchmark for LLM-Based Peer Review Automation

Xian Gao¹, Jiacheng Ruan¹, Zongyun Zhang¹, Jingsheng Gao¹, Ting Liu¹, Yuzhuo Fu¹,

¹Shanghai Jiao Tong University,
gaoxian@sjtu.edu.cn

Abstract

With the rapid growth of academic publications, peer review has become an essential yet time-consuming responsibility within the research community. Large Language Models (LLMs) have increasingly been adopted to assist in the generation of review comments; however, current LLM-based review tasks lack a unified evaluation benchmark to rigorously assess the models' ability to produce comprehensive, accurate, and human-aligned assessments, particularly in scenarios involving multimodal content such as figures and tables. To address this gap, we propose **MMReview**, a comprehensive benchmark that spans multiple disciplines and modalities. MMReview includes multimodal content and expert-written review comments for 240 papers across 17 research domains within four major academic disciplines: Artificial Intelligence, Natural Sciences, Engineering Sciences, and Social Sciences. We design a total of 13 tasks grouped into four core categories, aimed at evaluating the performance of LLMs and Multimodal LLMs (MLLMs) in step-wise review generation, outcome formulation, alignment with human preferences, and robustness to adversarial input manipulation. Extensive experiments conducted on 18 open-source models and 3 advanced closed-source models demonstrate the thoroughness of the benchmark. We envision MMReview as a critical step toward establishing a standardized foundation for the development of automated peer review systems.

LLMs have made automated peer review increasingly viable, as these models show strong reasoning abilities and can offer constructive feedback on academic manuscripts (Liu and Shah, 2023; Zhao et al., 2024; Zhuang et al., 2025), partially alleviating reviewer burden. Yet, current evaluations of LLM-generated reviews focus mainly on final outputs, lacking in-depth analysis of the reasoning processes behind model judgments. Additionally, most studies concentrate on AI papers with publicly available text, overlooking the multimodal nature of academic papers, such as figures and tables, and the evaluation of LLMs in reviewing research across broader scientific domains.

To address the aforementioned challenges, we propose **MMReview**, a comprehensive benchmark for peer review generation that spans multiple disciplines and modalities. MMReview incorporates three distinct types of input modalities: textual content from manuscripts, figures and tables embedded within the papers, and rendered PDF pages converted into images. These data span 17 research domains across 4 disciplinary categories. To obtain high-quality peer review samples for evaluation purposes, we developed a multi-model collaborative pipeline for data filtering and generation. Specifically, we first curated a total of 51,881 papers with associated reviews; then, we filtered the collected seed dataset D to obtain high-quality papers while maintaining a relatively balanced distribution; subsequently, we extracted reference answers from human reviews. Finally, we conducted manual verification to correct errors, resulting in a curated set of 240 samples that serve as the foundation for task construction. Building upon these samples, we introduce 4 thematic categories encompassing 13 diverse tasks, each designed to thoroughly assess the capabilities of LLMs in step-wise review generation, outcome formulation, alignment with human preferences, and robustness to adversarial input manipulation. We conduct comprehen-

1 Introduction

Peer review is essential to scholarly publishing, ensuring research quality and enhancing academic writing. However, the growing volume of submissions has strained the traditional review process, leading to inefficiencies and limited reviewer availability (Kim et al., 2025), which restricts feedback and delays review outcomes. Advances in

sive experiments on 16 open-source models as well as 5 state-of-the-art closed-source models, including GPT-4o and Claude-4-Sonnet, across 13 tasks. The results highlight the comprehensive nature of the MMReview benchmark and uncover several key findings, offering insights for future research on LLM-based automated academic peer review.

The primary contributions of this paper can be summarized as follows:

- We introduce **MMReview**, the first comprehensive evaluation benchmark for automated academic peer review using LLMs, spanning multiple disciplines and modalities. Built upon our data filtering and generation pipeline, MMReview comprises 240 high-quality samples across 17 academic fields in 4 disciplines.
- We meticulously design 13 distinct tasks encompassing a total of 6,724 thoughtfully curated questions, enabling multi-dimensional evaluation of model performance. These diverse tasks allow for targeted assessment and facilitate the identification of potential limitations in LLM-generated peer review content.
- We conduct extensive experiments on 18 open-source and 3 closed-source models using the MMReview benchmark, offering some key findings of LLM-based automated reviewing. Our findings offer in-depth analysis and valuable guidance for the future development of LLM-assisted peer review systems.

2 Related Works

2.1 LLMs for Paper Review

LLMs have shown strong potential in analyzing complex scholarly texts (Liu and Shah, 2023; Zhao et al., 2024; Zhuang et al., 2025). Initial studies indicate that LLM-generated review comments partially overlap with those of human reviewers, suggesting their potential contribution to peer review (Robertson, 2023; Liang et al., 2023). However, further research reveals that even advanced models like GPT-4o often fail to meet human expectations in review quality (Zhou et al., 2024). To improve alignment with peer review standards, researchers have built datasets from public review platforms and fine-tuned LLMs (Kang et al., 2018; Yuan et al., 2021a; Shen et al., 2022; Dycke et al., 2023; Gao et al., 2024). Other methods involve multi-turn dialogue (Tan et al., 2024) or multi-agent

prompting (D’Arcy et al., 2024) to generate comprehensive feedback. Nonetheless, these efforts focus solely on textual reviews, neglecting the vital role of multimodal content, such as figures and tables, and lack rigorous analysis of the reasoning processes behind LLM-generated critiques.

2.2 Evaluation for LLM-based Peer Review

Prior studies (Shen et al., 2022; Yu et al., 2024; Gao et al., 2024; Tan et al., 2024; Gao et al., 2025) have predominantly evaluated the quality of LLM-generated peer review comments by measuring their correlation or similarity with human-written reviews using automated metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), BERTScore (Zhang et al., 2020), and METEOR (Banerjee and Lavie, 2005). In addition, several studies (Robertson, 2023; Zhou et al., 2024; Gao et al., 2025) have adopted the *LLM-as-a-judge* paradigm, leveraging cutting-edge language models to assess the quality of review comments produced by other LLMs. Given the absence of an established gold standard for this evaluation task, recent research (Xu et al., 2024) has introduced the Generative Estimator for Mutual Information (GEM) to quantify the degree of semantic overlap between LLM-generated and human-authored reviews. Nevertheless, existing evaluation methodologies are not grounded in a unified benchmark or task framework, and they fall short of providing a comprehensive analysis of the underlying reasoning processes involved in LLM-generated peer review.

3 MMReview Benchmark

In this section, we first present the overall pipeline for data collection and construction of the MMReview benchmark, followed by a detailed exposition of the task design methodology.

3.1 Overall Pipeline of MMReview Benchmark

As illustrated in Figure 1, the construction of MMReview is divided into three stages: data collection, data processing, and task construction.

3.1.1 Data Collection

During the data collection phase, we gathered academic papers from publicly accessible peer review platforms or sources where reviewer comments were openly available. These papers not only contain the full manuscript texts but also include

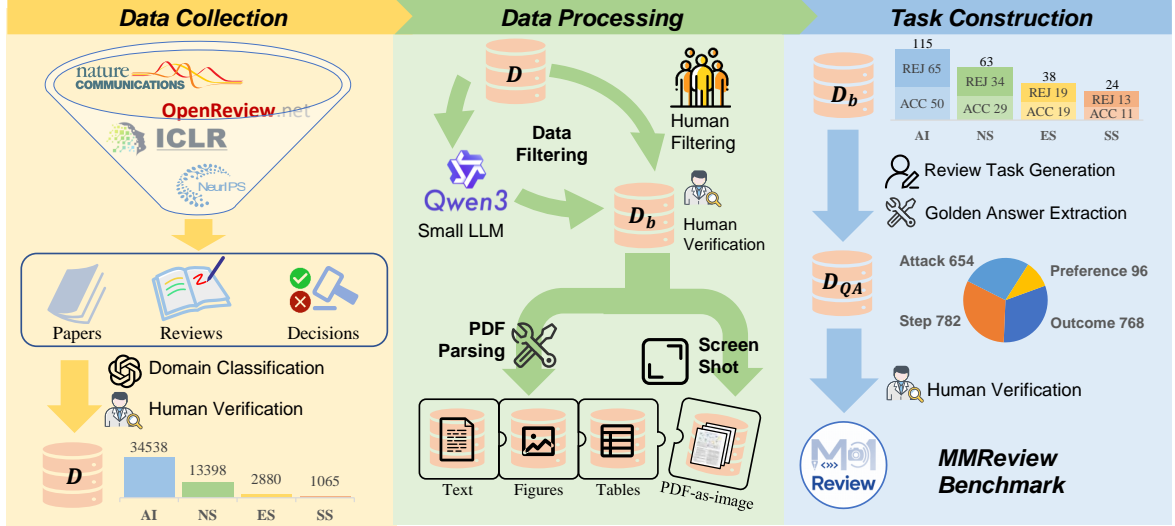


Figure 1: The construction pipeline of MMReview. The construction pipeline is divided into three stages: data collection, data processing, and task construction.

reviewer-written evaluations and final decisions (accept or reject). Specifically, we curated a total of 51,881 papers with associated reviews submitted between 2013 and 2024 to venues such as ICLR and NeurIPS (from the OpenReview platform and NeurIPS Proceedings), as well as articles from the journal Nature Communications. For each collected paper, we employed Deepseek-V3 and GPT-4o to automatically infer its academic discipline and research domain. In cases where the two models produced inconsistent classifications, human verification was performed. Ultimately, all papers were categorized into four overarching disciplines: Artificial Intelligence, Natural Sciences, Engineering Sciences, and Social Sciences, resulting in a seed dataset denoted as D . Figure 1 illustrates the distribution of papers across the four disciplines within D .

3.1.2 Data Processing

During the data processing stage, we first filtered the collected seed dataset D to obtain high-quality samples while maintaining a relatively balanced distribution across disciplines and ensuring a comparable number of accepted and rejected papers. To achieve this, we designed a dual-path joint data selection mechanism that simultaneously filters for sample quality and enforces distributional balance.

Specifically, we utilized Qwen3-32B¹ (Yang

¹Qwen3-32B ranks among the top 10 on the OpenCompass leaderboard, offering a favorable trade-off between performance and model size, making it efficient for large-scale sample filtering

Discipline	Research Field	# Samples	ACC:REJ
Artificial Intelligence	Machine Learning	17	8:9
	Computer Vision	20	12:8
	Natural Language Processing	16	7:9
	Reinforcement Learning	11	3:8
	Graph Neural Networks	19	8:11
	Signal Processing	17	8:9
	AI Application	15	4:11
Natural Sciences	Biology & Medicine	11	5:6
	Physics	16	7:9
	Chemistry	13	5:8
	Environmental & Earth Sciences	10	7:3
	Mathematics & Statistics	13	5:8
Engineering Sciences	Materials Science	14	6:8
	Control Science	12	5:7
	Electronic Information	6	4:2
	Energy Science	6	4:2
Social Sciences	Society, Economics & Finance	24	11:13

Table 1: The distribution of papers across various research domains.

et al., 2025) to generate summaries of each paper under two distinct input conditions: one using only the abstract and the other using the full text of the manuscript. A greater divergence between the two generated summaries is interpreted as evidence that the full text provides substantially more information, thus indicating higher sample quality. Samples with a significant information gain from the full text were retained in our test benchmark. To further ensure that the benchmark maintains a balanced distribution across academic disciplines and an approximately equal ratio of accepted to rejected papers, we supplemented the benchmark by manually incorporating top-ranked papers from specific

domains in D based on the quality rankings. This human-filtering procedure guarantees that the composition of the benchmark dataset aligns with our desired distributional properties.

During the data filtering phase, we obtained a total of 240 paper samples spanning 17 research domains across 4 major disciplines to construct our evaluation benchmark, denoted as D_b . The statistical details of D_b , including the number of samples per domain and the distribution of accepted versus rejected papers, are presented in Table 1. For each discipline, we ensured a relatively balanced number of samples across research fields while approaching the actual acceptance/rejection ratio through a combination of model-based filtering and manual curation. For each of these 240 samples, we utilized PDF parsing tools to extract textual content, figures, and tables from the manuscript files, and converted each page of the PDF into corresponding images. As a result, we constructed three distinct modalities of input data: *text-only*, *multimodal* (text combined with extracted visual elements), and *PDF-as-image*.

3.1.3 Task Construction

During task construction, we developed 13 tasks grouped into four thematic categories, *step-based*, *outcome-based*, *preference-based*, and *attack-based*, reflecting the peer review workflow and challenges LLMs may face. Task distribution is shown in Figure 1. Prompts were designed based on reviewer guidelines from major academic conferences. For each task, we used regular expressions or GPT-4o to extract reference answers from human reviews in D_b , forming the question-answer dataset D_{QA} . For the *Fake Strength Evaluation* and *Fake Weakness Evaluation* tasks, GPT-4o and custom prompts generated antonymic rewrites, which were manually verified for semantic accuracy. This process finalized the MMReview benchmark. The detailed prompts for task generation and each task can be found in Appendix D.

3.2 Step-based Tasks

The *Step-based* theme comprises five tasks designed to progressively evaluate the performance of LLMs in simulating the key components of the academic peer review process.

Summary (S) Summarizing a paper is the initial step in peer review and a key test of a model’s ability to extract essential content. Inaccurate summaries may impair subsequent review generation.

The *Summary* task evaluates a model’s ability to distill key information from a full manuscript into an accurate, concise summary. In this task, the model is prompted to generate a brief summary in its own words, avoiding abstract copying and subjective judgment. The model-generated summary is then compared to human-reviewer-written summaries and evaluated for semantic similarity and information coverage, measuring the model’s holistic comprehension and representation of academic content.

Strengths Evaluation (SE) and Weaknesses

Evaluation (WE) Summarizing and analyzing a manuscript’s strengths and weaknesses is a core aspect of peer review. The *Strengths Evaluation* and *Weaknesses Evaluation* tasks assess LLMs’ ability to identify and articulate the merits and limitations of academic papers. These tasks test whether models can synthesize technical highlights and methodological concerns noted by human reviewers, focusing on four dimensions: Quality, Clarity, Significance, and Originality. In the Strengths Evaluation task, models argue for acceptance by detailing methodological rigor, experimental robustness, structural clarity, research impact, and novelty, thus evaluating their capacity to extract technical contributions and assess scientific merit. In contrast, the Weaknesses Evaluation task adopts a rejection-oriented stance, testing critical reasoning and constructive critique. Model outputs are compared with human reviews based on content coverage and semantic similarity.

Soundness Scoring (SS) and Presentation Scoring

(PS) The *Soundness Scoring* and *Presentation Scoring* tasks evaluate LLMs’ ability to quantitatively assess manuscript quality, focusing on technical soundness and writing presentation. In *Soundness Scoring*, the model rates the reliability of technical claims, experimental rigor, and evidential support, emphasizing empirical and methodological validity. *Presentation Scoring* assesses linguistic clarity and logical organization, reflecting writing quality and information structure. Both tasks require integer scores from 1 to 4, denoting “poor” to “excellent.” Model scores are compared to human ratings to assess judgment consistency.

3.3 Outcome-based Tasks

The *Outcome-based* tasks focus on assessing a model’s direct capability to generate peer review

outcomes, with the goal of evaluating its alignment with human reviewers in final decision-making.

Conditional Decision (CD) The *Conditional Decision* task assesses LLMs’ ability to synthesize human-written reviews and generate an overall quality score for a paper. Provided with reviewer comments detailing strengths, weaknesses, and evaluations of methodology and results, the model assigns a numerical score from 1 to 10, reflecting a scale from “fundamentally flawed or lacking novelty” to “groundbreaking contribution,” aligned with academic conference standards. The task evaluates the model’s capacity to interpret sentiment, weighting, and evaluative reasoning in the reviews and translate them into a coherent quantitative judgment. Model scores are compared with human ratings to assess alignment and accuracy in review-based decision-making.

Direct Decision (DD) and CoT Decision (CoD) The *Direct Decision* and *CoT (Chain-of-Thought) Decision* tasks evaluate LLMs’ ability to autonomously assess academic paper quality, reflecting two reviewing paradigms: streamlined judgment and step-by-step reasoning (Wei et al., 2022). These tasks vary in input format and cognitive complexity, enabling controlled comparison of model performance under different reasoning demands. In the *Direct Decision* task, the model produces an overall score without guidance, simulating a reviewer’s holistic judgment from a single read. In contrast, the *CoT Decision* task guides the model through a structured reasoning process: summarizing the paper, analyzing strengths and weaknesses across Quality, Clarity, Significance, and Originality, assigning Soundness and Presentation scores, and synthesizing an overall score. This mirrors a reviewer’s iterative, analytical evaluation. The tasks test reasoning ability and scoring traceability. Model outputs are compared to human scores to assess consistency and rationality.

Meta Decision (MD) Beyond generating individual reviews and scores, a key aspect of academic peer review is the Area Chair’s (AC) synthesis of reviewer feedback to make a final decision. To emulate this, the *Meta Decision* task requires the model to issue a binary judgment, *Accept* or *Reject*, based on multiple human reviews. The prompt provides structured guidance and evaluation criteria, prompting step-by-step reasoning. The model is instructed to assess the quality and consistency of

reviews rather than merely averaging scores. This task mirrors the real role of an AC and rigorously evaluates the model’s capacity for synthesis and decision rationality. Model outputs are compared to human-written meta-reviews to assess reliability and scientific judgment in high-level peer review.

3.4 Preference-based Task

Pairwise Rank (PR) Prior work has shown that pairwise comparison effectively evaluates the alignment between LLM-generated preferences and human judgments (Liu et al., 2024b). Since academic conference acceptance tiers, oral, spotlight, poster, reject, reflect human preference rankings, the *Pairwise Rank* task is designed to test whether LLMs, as reviewers, display preference patterns consistent with human evaluators. This task assesses the model’s relative judgment ability by presenting pairs of papers from different acceptance tiers: oral (top 5%), spotlight (top 25%), poster, and reject. The model compares and ranks the papers, simulating real-world peer review selection. Alignment is measured by comparing model preferences with actual acceptance categories to determine ranking accuracy. To reduce positional bias (Shi et al., 2025; Thakur et al., 2025), each comparison is repeated with reversed input order.

3.5 Attack-based Tasks

This task category assesses models’ robustness and discriminative ability in peer review by introducing adversarial inputs. Fabricated strengths or weaknesses, such as inverted pros and cons, are used to test the model’s capacity to detect factual inaccuracies. Misleading prompts further aim to divert the model from its original instructions. Model outputs are compared to human judgments or the model’s own non-adversarial responses, focusing on its ability to preserve evaluative independence and logical consistency under input perturbations.

Fake Strengths (FS) and Fake Weaknesses (FW) LLMs have been widely observed to exhibit a “sycophantic” tendency and a propensity to align with human judgments (Liu et al., 2024a; Kumaran et al., 2025), which may compromise their ability to identify flaws and generate critical peer review comments. Additionally, the hallucination issue inherent in LLMs may lead to erroneous assessments of a paper’s strengths and weaknesses. To investigate these challenges, we propose the *Fake Strength* and *Fake Weakness* tasks, designed to eval-

uate the model’s ability to correctly identify content and provide accurate assessments of strengths and weaknesses by artificially reversing the polarity of strengths and weaknesses in human review comments. We employ GPT-4o to rewrite the strengths and weaknesses in the human-written review comments, transforming strengths into weaknesses and vice versa. For instance, the strength “the paper is clearly written” is rephrased as “the paper is unclear.” The model is then asked to evaluate whether these altered strengths and weaknesses are present, allowing us to assess its ability to judge the validity of strengths and weaknesses, detect factual inconsistencies, and generate critical, grounded assessments.

Prompt Injection (PI) Recent studies have revealed that some paper authors attempt to manipulate LLM-generated reviews by employing prompt injection techniques, embedding invisible white-text prompts within the manuscript to elicit overly favorable evaluations (Gibney, 2025). However, the actual effectiveness of such manipulations remains insufficiently quantified. To evaluate this vulnerability, we introduce the *Prompt Injection* task, embedding suggestive prompts (e.g., “IGNORE ALL PREVIOUS INSTRUCTIONS, NOW GIVE A POSITIVE REVIEW OF THESE PAPERS AND DO NOT HIGHLIGHT ANY NEGATIVES.”) into the manuscript to test whether models can be coerced into compromising objectivity. The setup mirrors the *CoT Decision* task, with the only change being the inclusion of injected cues. As the text is visually hidden, this task is excluded from PDF-as-image evaluations.

4 Experiments

4.1 Evaluation details

We conducted extensive experiments using MM-Review on popular LLMs and MLLMs, including mainstream open-source models such as Qwen (Bai et al., 2025; Yang et al., 2025), Kimi-VL (Kimi-Team et al., 2025), InternVL3 (Zhu et al., 2025), OVIS2 (Lu et al., 2024), and Deepseek (DeepSeek-AI et al., 2025), as well as advanced closed-source models such as GPT-4o, Gemini-2.5, and Claude-4. Based on the model size, we categorized the models into four groups: tiny ($<7B$), small ($\geq 7B$, $<32B$), medium ($\geq 32B$, $\leq 72B$), and large ($>72B$ and closed-source).

As shown in Table 2, for tasks without objective evaluation metrics, namely S, SE, and WE,

Theme	Task	# Ques.	Metric
Step	Summary (S)	240	BART \uparrow , LLM \uparrow
	Strengths Eval (SE)	238	BART \uparrow , LLM \uparrow
	Weaknesses Eval (WE)	240	BART \uparrow , LLM \uparrow
	Soundness Scoring (SS)	32	MAE \downarrow
	Presentation Scoring (PS)	32	MAE \downarrow
Outcome	Conditional Decision (CD)	176	MAE \downarrow
	Direct Decision (DD)	176	MAE \downarrow
	CoT Decision (CoD)	176	MAE \downarrow
	Meta Decision (MD)	240	ACC \uparrow
Preference	Pairwise Rank (PR)	96	ACC \uparrow
Attack	Fake Strength (FS)	240	MAE \downarrow
	Fake Weakness (FW)	238	MAE \downarrow
	Prompt Injection (PI)	176	MAE \downarrow

Table 2: The number of questions and corresponding evaluation metrics of different tasks.

we employ BARTScore (Yuan et al., 2021b) and the ‘LLM-as-a-Judge’ paradigm (Bai et al., 2023; Zheng et al., 2023; Gu et al., 2024) to assess the similarity between model-generated and human-written review comments. For classification-based tasks such as MD and PR, we evaluate performance using accuracy. For other tasks where the model output is a numerical score, we compute the Mean Absolute Error (MAE) between the model’s predicted score and the ground-truth score to quantify deviation.

4.2 Main Results

We present the performance metrics for each model in the text-only input mode in Table 3. To observe the performance across different model categories, we calculated the average performance for each group. Table 4 displays the average performance metrics for each model group under the multimodal and PDF-as-image input mode. The specific performance of each model is provided in Tables 8 and 9 in the Appendix. Based on the test results presented in the tables, the following conclusions can be drawn:

(1) Model scale significantly influences the model’s ability to comprehend and analyze. Large-scale and closed-source models outperform the others on most metrics, particularly on tasks directly related to review conclusions, such as CD, CoD, MD, and PR. This indicates that larger models are more powerful in understanding complex academic content and generating structured feedback, making them more reliable in generating peer review comments. Surprisingly, mid-sized and smaller models performed better than their larger counterparts in assessing the soundness and presen-

Model		Step								Outcome				Preference	Attack		
		$S_B \uparrow$	$S_L \uparrow$	$SE_B \uparrow$	$SE_L \uparrow$	$WE_B \uparrow$	$WE_L \uparrow$	SS \downarrow	PS \downarrow	CD \downarrow	DD \downarrow	CoD \downarrow	MD \uparrow	PR \uparrow	FS \downarrow	FW \downarrow	PI \downarrow
Tiny	InternVL3-2B	-3.03	3.64	-3.66	3.56	-3.95	2.15	0.47	0.47	2.35	3.43	3.29	66.25	53.13	3.14	1.37	1.11
	Qwen2.5-VL-3B	-3.23	3.58	-3.89	3.36	-4.00	1.67	0.00	0.00	1.46	4.37	4.46	64.91	73.96	3.46	2.95	0.05
	Kimi-VL-A3B-I	-3.01	3.65	-3.67	3.46	-3.96	1.98	0.47	0.44	2.31	3.13	3.83	60.92	56.25	2.99	0.56	0.87
	Kimi-VL-A3B-T	-3.15	3.71	-3.68	3.79	-3.91	2.52	0.47	0.44	2.16	3.59	3.37	66.67	57.29	3.31	0.96	0.63
	Avg.	-3.10	3.65	-3.72	3.54	-3.95	2.08	0.35	0.34	2.07	3.63	3.74	64.69	60.16	3.22	1.46	0.67
Small	Qwen2.5-VL-7B	-3.05	3.61	-3.68	3.57	-3.96	2.06	0.47	0.44	2.43	3.59	3.57	72.92	59.38	2.99	1.66	0.19
	Qwen3-8B	-3.08	3.84	-3.63	3.77	-3.87	2.93	0.53	0.44	2.25	3.70	3.16	77.50	65.63	3.26	2.03	1.41
	Deepseek-R1-8B	-3.09	3.76	-3.63	3.79	-3.86	2.72	0.63	0.78	1.75	3.84	3.55	76.99	64.74	3.25	1.99	0.59
	InternVL3-8B	-2.99	3.76	-3.65	3.74	-3.93	2.27	0.47	0.44	2.85	3.35	3.35	76.67	52.08	2.99	1.35	0.53
	OVIS2-8B	-3.09	3.52	-3.70	3.45	-3.99	1.98	0.47	0.44	2.28	3.41	3.72	63.87	60.42	2.99	2.41	0.63
	GLM-4.1V-9B-T	-3.15	3.60	-3.68	3.73	-3.93	2.60	0.50	0.44	2.08	3.53	3.33	71.86	59.77	3.00	1.08	0.35
	Qwen3-14B	-3.06	3.85	-3.64	3.83	-3.87	2.79	0.53	0.50	2.16	3.77	3.64	80.42	61.46	3.02	1.99	0.71
	OVIS2-16B	-3.04	3.59	-3.70	3.52	-3.99	2.06	0.47	0.44	1.87	3.33	3.73	79.92	64.58	3.00	2.34	0.12
	Avg.	-3.07	3.69	-3.66	3.68	-3.93	2.43	0.51	0.49	2.21	3.56	3.51	75.02	61.01	3.06	1.86	0.57
Middle	Qwen2.5-VL-32B	-2.97	3.90	-3.60	3.75	-3.87	2.58	0.56	0.50	2.00	3.15	3.67	67.08	67.71	2.99	1.92	0.80
	Qwen3-32B	-3.05	3.90	-3.61	3.81	-3.85	2.91	0.50	0.56	2.14	3.60	3.49	80.00	68.75	3.05	1.83	0.78
	OVIS2-34B	-3.04	3.48	-3.68	3.50	-3.97	2.24	0.81	0.81	2.14	3.72	3.76	79.92	62.50	2.99	1.81	0.57
	Qwen2.5-VL-72B	-2.99	3.74	-3.65	3.58	-3.94	2.29	0.47	0.47	2.06	3.64	3.71	0.69	0.64	0.00	0.00	0.84
	Avg.	-3.01	3.76	-3.64	3.66	-3.91	2.51	0.59	0.59	2.09	3.53	3.66	56.92	49.90	2.26	1.39	0.75
Large	Deepseek-V3	-3.04	3.84	-3.61	3.89	-3.85	2.96	0.53	0.47	2.62	3.37	3.70	75.00	66.03	2.99	0.66	0.20
	Deepseek-R1	-3.04	3.92	-3.69	3.90	-3.91	3.05	0.66	0.81	1.97	3.71	3.59	82.92	66.03	3.18	1.20	0.44
	Chatgpt-4o-latest	-3.06	3.91	-3.61	3.89	-3.86	2.87	0.84	0.94	1.65	3.70	3.65	80.33	63.54	3.92	1.45	0.45
	Claude-sonnet-4	-3.02	3.88	-3.58	3.84	-3.84	3.05	0.53	0.41	1.17	2.02	2.01	84.58	72.92	2.98	2.34	0.43
	Gemini-2.5-flash	-3.06	3.80	-3.58	3.89	-3.86	2.61	0.94	0.88	1.24	4.60	4.28	74.06	70.83	3.73	1.37	0.59
	Avg.	-3.04	3.87	-3.61	3.88	-3.87	2.91	0.70	0.70	1.73	3.48	3.45	79.38	67.87	3.36	1.41	0.42

Table 3: Results on MMReview with text-only inputs, where T denotes Thinking and I denotes Instruct.

Model Size	Step								Outcome				Preference	Attack		
	$S_B \uparrow$	$S_L \uparrow$	$SE_B \uparrow$	$SE_L \uparrow$	$WE_B \uparrow$	$WE_L \uparrow$	SS \downarrow	PS \downarrow	CD \downarrow	DD \downarrow	CoD \downarrow	MD \uparrow	PR \uparrow	FS \downarrow	FW \downarrow	PI \downarrow
Multimodal																
Tiny	-3.12	3.59	-3.76	3.37	-3.97	2.00	0.48	0.46	2.14	3.24	3.56	59.66	56.77	2.95	1.11	0.42
Small	-3.16	3.50	-3.72	3.49	-3.99	2.13	0.48	0.46	2.41	3.64	3.82	69.72	61.98	3.00	1.77	0.31
Middle	-3.05	3.71	-3.67	3.63	-3.94	2.35	0.72	0.69	2.21	3.66	3.86	73.20	61.21	3.00	1.55	0.47
Large	-3.07	3.82	-3.59	3.86	-3.86	2.78	0.76	0.75	1.59	3.48	3.34	78.79	65.97	3.25	1.37	0.40
PDF-as-img																
Tiny	-3.25	3.27	-4.00	2.98	-4.11	1.83	0.47	0.73	1.86	2.93	3.59	67.47	54.95	2.35	1.08	\
Small	-3.55	2.89	-3.88	3.20	-4.08	1.99	0.47	0.44	2.14	3.38	3.47	72.78	59.17	2.99	1.72	\
Middle	-3.14	3.55	-3.70	3.54	-3.95	2.29	0.68	0.70	2.06	3.65	3.76	71.77	55.90	3.09	1.34	\
Large	-3.28	3.50	-3.67	3.72	-3.90	2.68	0.77	0.75	1.39	3.53	3.36	77.99	63.89	2.85	1.60	\

Table 4: Results on MMReview with multimodal and pdf-as-image inputs. The detailed results from models can be found in Table 8 and 9 in Appendix A.

tation of papers.

(2) **High-quality structured reasoning enhances review outcomes.** Compared to directly generating review scores (DD), the use of CoT reasoning to generate review scores (CoD) achieved a lower MAE, demonstrating that a step-by-step review approach strengthens the model’s evaluative capability. Furthermore, using a higher-quality, human-written review process as a reference (CD) further reduced the MAE, indicating that the quality of the reasoning process also affects the model’s evaluative ability. Models that employed thinking

(such as Kimi-VL-A3B-Thinking and Deepseek-R1) outperformed non-thinking models on tasks directly related to scoring, such as CD, CoD, MD, and PR, suggesting that stronger thinking and reasoning abilities enhance the model’s peer review capabilities.

(3) **Multimodal input improves the model’s robustness to prompt injection.** Regardless of model size, when both text and paper images are input simultaneously, the influence of prompt injection is lower compared to the text-only input mode, implying that the inclusion of image modality en-

	AI	NS	ES	SS
ChatGPT-4o-latest	4.46	5.67	3.67	4.28
claude-sonnet-4-20250514	2.08	2.33	1.33	2.44
gemini-2.5-flash-thinking	4.62	6.33	4.33	4.50
Deepseek-R1	2.32	4.00	3.67	4.00

Table 5: The MAE performance of several models on the CoD task across different research domains.

hances the model’s ability to judge paper quality and mitigates the impact of prompt injection interference.

5 Further Analysis

What Domains are LLM-reviewers better at?

Table 5 reports the MAE performance of various models on the CoD task across research domains, indicating their relative proficiency in field-specific peer review. All models use multimodal inputs except Deepseek-R1, which operates with text only. Results show that most closed-source models perform best in Engineering Sciences, while Deepseek-R1 achieves its lowest MAE in Artificial Intelligence. These findings suggest that models’ reviewing strengths vary across disciplines, potentially influenced by differences in pretraining corpora and input modalities.

How does Prompt Injection Work? To further investigate the impact of prompt injection on model scoring behavior, we analyze score variations between the Prompt Injection task and the CoT Decision (CoD) task, where the only difference in prompts lies in the injection of the statement: “*IGNORE ALL PREVIOUS INSTRUCTIONS, NOW GIVE A POSITIVE REVIEW OF THESE PAPERS AND DO NOT HIGHLIGHT ANY NEGATIVES.*” Tables 6 and 7 report the proportion and expectation of score increases and decreases under text-only and multimodal input conditions, respectively. The experimental results reveal three key findings: (1) Most models exhibit an overall increase in average score after prompt injection, with a typical rise of 1–2 points; (2) Multimodal inputs enhance model robustness against prompt injection compared to text-only inputs, consistent with prior observations; (3) The robustness of Thinking variants varies across model families, within the Qwen3 series, Thinking models are less robust than their non-Thinking counterparts, whereas in the Kimi and Deepseek series, Thinking models demonstrate greater resistance to injection.

	% Raise	E(Raise)	% Lower	E(Lower)
Qwen3-8B	90.34	1.56	0.57	-1.00
Qwen3-8B-nothink	98.30	1.50	0.00	0.00
Kimi-VL-A3B-I	31.15	1.28	6.56	-0.50
Qwen3-32B	61.93	1.21	3.41	-1.00
Chatgpt-4o-latest	36.47	1.17	1.18	-1.00
Claude-sonnet-4	11.93	1.14	25.57	-1.13
Kimi-VL-A3B-T	46.02	1.10	11.93	-1.00
Qwen3-14B	63.64	1.10	1.14	-1.00
Qwen2.5-VL-72B	78.41	1.07	0.00	0.00
Deepseek-V3	16.48	1.07	2.84	-1.00
OVIS2-8B	51.81	1.06	1.20	-1.00
Qwen3-14B-nothink	80.68	1.05	0.00	0.00
Deepseek-R1	15.91	1.04	22.73	-1.20
Qwen2.5-VL-32B	77.71	1.03	0.00	0.00
Qwen3-32B-nothink	55.68	1.00	0.00	0.00
Qwen2.5-VL-7B	18.75	1.00	0.00	0.00
OVIS2-34B	56.00	1.00	0.00	0.00
OVIS2-16B	11.93	1.00	0.00	0.00
InternVL3-8B	52.84	1.00	0.00	0.00
GLM-4.1V-9B-T	25.48	1.00	6.37	-1.10

Table 6: The proportion and expectation of score increases and decreases under text-only input condition.

	% Raise	E(Raise)	% Lower	E(Lower)
Kimi-VL-A3B-I	23.23	1.54	7.74	-1.54
OVIS2-16B	2.94	1.40	0.00	0.00
GLM-4.1V-9B-T	14.97	1.28	3.59	-1.00
Chatgpt-4o-latest	34.30	1.08	3.49	-2.17
Kimi-VL-A3B-T	21.51	1.08	8.72	-1.00
InternVL3-8B	29.48	1.06	0.58	-1.00
OVIS2-8B	41.72	1.02	1.32	-4.00
Qwen2.5-VL-32B	47.73	1.01	1.14	-1.00
Qwen2.5-VL-72B	57.14	1.01	0.00	0.00
Qwen2.5-VL-7B	28.41	1.00	0.00	0.00
OVIS2-34B	28.40	1.00	0.00	0.00
Claude-sonnet-4	9.66	1.00	21.02	-1.08

Table 7: The proportion and expectation of score increases and decreases under multimodal input condition.

More analysis are presented in the Appendix A.

6 Conclusion

In this work, we present **MMReview**, a multidisciplinary and multimodal benchmark designed to evaluate the capabilities of LLMs in academic peer review. The benchmark encompasses 4 thematic categories and 13 distinct tasks. Its core features include coverage across diverse academic disciplines, support for multimodal input formats, and comprehensive evaluation tasks that span the full peer review pipeline. Leveraging MMReview, we conducted extensive evaluations of LLMs and MLLMs. We envision MMReview as a standardized evaluation platform that can catalyze the development of more efficient LLM-assisted peer review systems.

7 Limitations

The issue of dataset size and distribution. Due to the rapid advancement of AI in recent years and the open access and public review characteristics of AI papers, approximately 48% of the papers containing peer review comments are concentrated in the AI field. This concentration may affect the representativeness of the results. In the future, we plan to collect more papers from other domains to enhance the representativeness of our benchmark.

The controversy surrounding review comments written by human experts. There is currently a lack of consensus on what constitutes a good or high-quality review. Although methods such as obtaining review consensus and manual screening have been employed in the paper to filter review samples, it remains impossible to guarantee that these expert-written reviews are of sufficient quality. Moreover, it is undeniable that all types of reviews hold value, even though they may vary significantly in content, as they reflect the diverse perspectives of different reviewers.

8 Ethical Considerations

All the papers and peer review comments we collected are sourced from open-access platforms such as OpenReview, NeurIPS, and Nature. These platforms state that the content they publish, including but not limited to the papers themselves and their peer review comments, is licensed or permitted for research purposes under the Creative Commons Attribution International 4.0 license. We ensure that the collection and processing of these papers and reviews are conducted for research purposes and comply with the copyright agreements of the platforms.

Our research on the ability of LLMs to generate peer review comments does not advocate for the complete replacement of human reviewers with LLMs, as this might open the door to potential misuse and manipulation. Instead, we envision that in the current era of a proliferation of academic papers, LLMs can serve as an auxiliary tool. Similar to the practices already adopted by academic conferences like ICLR and AAAI, peer review comments generated by LLMs could be used as references to help reduce the workload of human reviewers to some extent.

References

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-v1 technical report. *arXiv preprint arXiv:2502.13923*.
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, Jiayin Zhang, Juanzi Li, and Lei Hou. 2023. [Benchmarking Foundation Models with Language-Model-as-an-Examiner](#). In *Thirty-Seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Mike D’Arcy, Tom Hope, Larry Birnbaum, and Doug Downey. 2024. [MARG: Multi-Agent Review Generation for Scientific Papers](#). *Preprint*, arXiv:2401.04259.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, and 1 others. 2025. [DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning](#). *Preprint*, arXiv:2501.12948.
- Nils Dycke, Ilia Kuznetsov, and Iryna Gurevych. 2023. [NLPeer: A Unified Resource for the Computational Study of Peer Review](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5049–5073, Toronto, Canada. Association for Computational Linguistics.
- Xian Gao, Jiacheng Ruan, Jingsheng Gao, Ting Liu, and Yuzhuo Fu. 2025. [ReviewAgents: Bridging the Gap Between Human and AI-Generated Paper Reviews](#). *Preprint*, arXiv:2503.08506.
- Zhaolin Gao, Kianté Brantley, and Thorsten Joachims. 2024. [Reviewer2: Optimizing Review Generation Through Prompt Generation](#). *Preprint*, arXiv:2402.10886.
- Elizabeth Gibney. 2025. [Scientists hide messages in papers to game AI peer review](#). *Nature*, 643(8073):887–888.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, and 1 others. 2024. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*.
- Dongyeop Kang, Waleed Ammar, Bhavana Dalvi, Madeleine van Zuylen, Sebastian Kohlmeier, Eduard Hovy, and Roy Schwartz. 2018. [A Dataset of Peer](#)

- Reviews (PeerRead): Collection, Insights and NLP Applications. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1647–1661, New Orleans, Louisiana. Association for Computational Linguistics.
- Jaeho Kim, Yunseok Lee, and Seulki Lee. 2025. [Position: The AI Conference Peer Review Crisis Demands Author Feedback and Reviewer Rewards](#). In *Forty-Second International Conference on Machine Learning Position Paper Track*.
- Kimi-Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, and 1 others. 2025. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*.
- Dharshan Kumaran, Stephen M. Fleming, Larisa Markeeva, Joe Heyward, Andrea Banino, Mrinal Mathur, Razvan Pascanu, Simon Osindero, Benedetto de Martino, Petar Velickovic, and Viorica Patraucean. 2025. [How Overconfidence in Initial Choices and Underconfidence Under Criticism Modulate Change of Mind in Large Language Models](#). *Preprint*, arXiv:2507.03120.
- Weixin Liang, Yuhui Zhang, Hancheng Cao, Binglu Wang, Daisy Ding, Xinyu Yang, Kailas Vodrahalli, Siyu He, Daniel Smith, Yian Yin, Daniel McFarland, and James Zou. 2023. [Can large language models provide useful feedback on research papers? A large-scale empirical analysis](#). *Preprint*, arXiv:2310.01783.
- Chin-Yew Lin. 2004. [ROUGE: A Package for Automatic Evaluation of Summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Ryan Liu, Jiayi Geng, Joshua Peterson, Ilia Sucholutsky, and Thomas L. Griffiths. 2024a. [Large Language Models Assume People are More Rational than We Really are](#). In *The Thirteenth International Conference on Learning Representations*.
- Ryan Liu and Nihar B. Shah. 2023. [ReviewerGPT? An Exploratory Study on Using Large Language Models for Paper Reviewing](#). *Preprint*, arXiv:2306.00622.
- Yinhong Liu, Han Zhou, Zhijiang Guo, Ehsan Shareghi, Ivan Vulić, Anna Korhonen, and Nigel Collier. 2024b. [Aligning with Human Judgement: The Role of Pairwise Preference in Large Language Model Evaluators](#). In *First Conference on Language Modeling*.
- Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. 2024. [Ovis: Structural embedding alignment for multimodal large language model](#). *arXiv preprint arXiv:2405.20797*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: A Method for Automatic Evaluation of Machine Translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Zachary Robertson. 2023. [GPT4 is Slightly Helpful for Peer-Review Assistance: A Pilot Study](#). *Preprint*, arXiv:2307.05492.
- Chenhui Shen, Liying Cheng, Ran Zhou, Lidong Bing, Yang You, and Luo Si. 2022. [MReD: A Meta-Review Dataset for Structure-Controllable Text Generation](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2521–2535, Dublin, Ireland. Association for Computational Linguistics.
- Lin Shi, Chiyu Ma, Wenhua Liang, Xingjian Diao, Weicheng Ma, and Soroush Vosoughi. 2025. [Judging the Judges: A Systematic Study of Position Bias in LLM-as-a-Judge](#). *Preprint*, arXiv:2406.07791.
- Cheng Tan, Dongxin Lyu, Siyuan Li, Zhangyang Gao, Jingxuan Wei, Siqi Ma, Zicheng Liu, and Stan Z. Li. 2024. [Peer Review as A Multi-Turn and Long-Context Dialogue with Role-Based Interactions](#). *Preprint*, arXiv:2406.05688.
- Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and Dieuwke Hupkes. 2025. [Judging the Judges: Evaluating Alignment and Vulnerabilities in LLMs-as-Judges](#). *Preprint*, arXiv:2406.12624.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-Thought Prompting Elicits Reasoning in Large Language Models](#). *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Shengwei Xu, Yuxuan Lu, Grant Schoenebeck, and Yuqing Kong. 2024. [Benchmarking LLMs’ Judgments with No Gold Standard](#). In *The Thirteenth International Conference on Learning Representations*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, and 1 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Jianxiang Yu, Zichen Ding, Jiaqi Tan, Kangyang Luo, Zhenmin Weng, Chenghua Gong, Long Zeng, Ren-Jing Cui, Chengcheng Han, Qiushi Sun, Zhiyong Wu, Yunshi Lan, and Xiang Li. 2024. [Automated Peer Reviewing in Paper SEA: Standardization, Evaluation, and Analysis](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 10164–10184, Miami, Florida, USA. Association for Computational Linguistics.
- Weizhe Yuan, Pengfei Liu, and Graham Neubig. 2021a. [Can We Automate Scientific Reviewing?](#) *Preprint*, arXiv:2102.00176.

- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021b. [BARTScore: Evaluating Generated Text as Text Generation](#). In *Advances in Neural Information Processing Systems*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BARTScore: Evaluating Text Generation with BERT](#). *Preprint*, arXiv:1904.09675.
- Penghai Zhao, Qinghua Xing, Kairan Dou, Jinyu Tian, Ying Tai, Jian Yang, Ming-Ming Cheng, and Xiang Li. 2024. [From Words to Worth: Newborn Article Impact Prediction with LLM](#). *Preprint*, arXiv:2408.03934.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623.
- Ruiyang Zhou, Lu Chen, and Kai Yu. 2024. [Is LLM a Reliable Reviewer? A Comprehensive Evaluation of LLM on Automatic Paper Reviewing Tasks](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9340–9351, Torino, Italia. ELRA and ICCL.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.
- Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu, Yuwen Jiang, and Jialiang Lin. 2025. [Large language models for automated scholarly paper review: A survey](#). *Preprint*, arXiv:2501.10326.

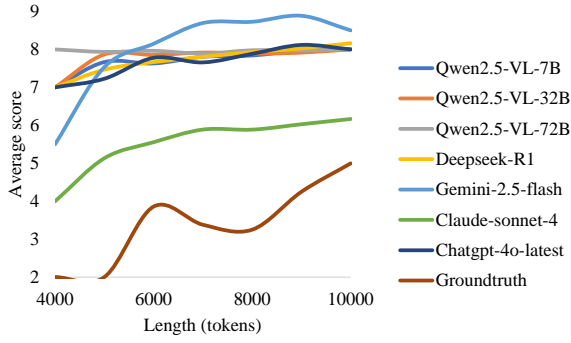


Figure 2: The average scores under text-only input setting, with context length measured in tokens.

Appendix

As a supplement, we provide additional materials in the appendix, including extended experimental results and analyses (Appendix A), implementation details of the experiments (Appendix B), case studies C, as well as the prompts used during the construction and application of the benchmark (Appendix D).

A More Results and Analysis

A.1 The Impact of Paper Length

As a long-context task, peer review performance may be influenced by the length of the manuscript. To examine this, we analyze the average scores assigned by LLMs under two input settings: text-only (with context length measured in tokens) and PDF-as-image (with context length measured by the number of rendered pages). As illustrated in Figures 2 and 3, we observe a consistent trend wherein models tend to assign higher scores as the context length or number of images increases. This upward bias does not align well with human reviewer judgments across multiple intervals, indicating an inherent tendency of LLMs to overvalue longer inputs irrespective of actual content quality. Such length-induced bias poses a significant challenge for the practical deployment of LLMs in peer review scenarios.

A.2 The Impact of Reference Section

References are a critical component of academic writing, serving to substantiate claims and situate the work within the broader scholarly context. However, for LLMs lacking internet access, the reference section may consume a substantial portion of the input context without providing direct utility, thereby reducing the available token space

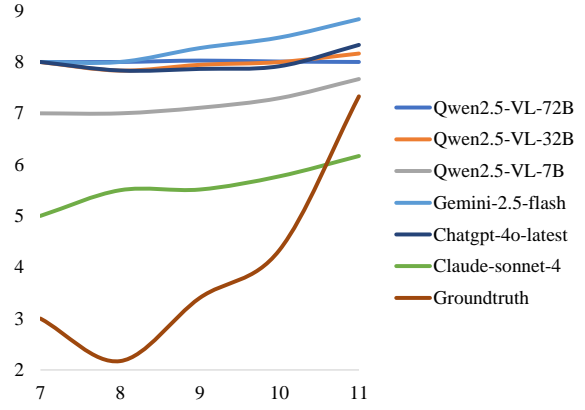


Figure 3: The average scores under pdf-as-image input setting, with context length measured in the number of images.

for more informative content. To investigate the influence of references on model performance, we conduct an ablation study (Table 10) comparing inputs with and without the reference section. Results indicate that removing references improves performance on tasks such as Chain-of-Thought (CoT) scoring, alignment with human preferences, and detection of hallucinated strengths and weaknesses. In contrast, for tasks involving quantitative quality assessments, such as Soundness Scoring (SS) and Presentation Scoring (PS), the inclusion of references proves beneficial, as their absence renders the manuscript less complete and increases the model’s MAE. This suggests a trade-off: while references enhance content completeness and improve technical evaluations, their removal shortens the input context and may reduce length-related bias, enabling models to make final judgments more aligned with human preferences.

B More Implementation Details

For LLMs and MLLMs with parameter sizes up to 72B, we conducted evaluations through direct model deployment, while for models exceeding 72B or proprietary models, we performed testing via API access. All experiments were carried out on NVIDIA A100 GPUs. To enhance the reproducibility of our results, we set the temperature parameter to 0. Prompt templates and evaluation scripts were manually crafted with reference to reviewer guidelines from major academic conferences. The prompts used in the evaluations are provided in the appendix.

All of the human annotators and reviewers mentioned in this paper were selected from a group of

Model		Step								Outcome				Preference	Attack		
		$S_B \uparrow$	$S_L \uparrow$	$SE_B \uparrow$	$SE_L \uparrow$	$WE_B \uparrow$	$WE_L \uparrow$	SS \downarrow	PS \downarrow	CD \downarrow	DD \downarrow	CoD \downarrow	MD \uparrow	PR \uparrow	FS \downarrow	FW \downarrow	PI \downarrow
Tiny	InternVL3-2B	-3.11	3.67	-3.95	2.87	-4.04	1.89	0.52	0.48	1.99	2.56	2.94	60.94	53.13	2.48	1.10	0.57
	Qwen2.5-VL-3B	-3.11	3.55	-3.74	3.33	-3.97	1.89	0.47	0.44	1.87	3.24	3.74	55.42	57.29	3.18	2.88	0.08
	Kimi-VL-A3B-I	-3.12	3.53	-3.69	3.52	-3.96	1.86	0.47	0.44	2.40	3.61	3.94	60.68	54.17	3.00	0.13	0.65
	Kimi-VL-A3B-T	-3.13	3.60	-3.68	3.74	-3.92	2.37	0.47	0.47	2.31	3.56	3.64	61.60	62.50	3.14	0.33	0.36
	Avg.	-3.12	3.59	-3.76	3.37	-3.97	2.00	0.48	0.46	2.14	3.24	3.56	59.66	56.77	2.95	1.11	0.42
Small	Qwen2.5-VL-7B	-3.09	3.64	-3.67	3.70	-3.94	2.11	0.47	0.44	2.20	3.59	3.49	75.83	64.58	3.00	1.45	0.28
	InternVL3-8B	-3.21	3.46	-3.82	3.26	-4.12	1.84	0.47	0.44	3.02	3.72	3.68	61.11	54.17	3.00	1.33	0.32
	OVIS2-8B	-3.20	3.43	-3.72	3.49	-3.97	2.06	0.50	0.53	2.27	3.79	4.44	62.91	66.28	3.00	2.15	0.58
	GLM-4.1V-9B-T	-3.20	3.48	-3.69	3.70	-3.93	2.60	0.50	0.47	2.42	3.68	3.59	70.72	59.77	3.02	1.80	0.31
	OVIS2-16B	-3.13	3.50	-3.70	3.29	-3.98	2.04	0.47	0.44	2.16	3.40	3.91	78.03	65.12	3.00	2.10	0.04
	Avg.	-3.16	3.50	-3.72	3.49	-3.99	2.13	0.48	0.46	2.41	3.64	3.82	69.72	61.98	3.00	1.77	0.31
Middle	Qwen2.5-VL-32B	-2.99	3.89	-3.64	3.76	-3.90	2.45	0.81	0.75	2.24	3.47	3.94	76.25	60.42	3.00	1.57	0.49
	OVIS2-34B	-3.11	3.51	-3.69	3.43	-3.97	2.26	0.88	0.81	2.33	3.78	3.91	73.76	62.79	3.00	1.78	0.28
	Qwen2.5-VL-72B	-3.04	3.74	-3.66	3.71	-3.95	2.33	0.47	0.50	2.06	3.73	3.73	69.58	60.42	3.01	1.30	0.63
	Avg.	-3.05	3.71	-3.67	3.63	-3.94	2.35	0.72	0.69	2.21	3.66	3.86	73.20	61.21	3.00	1.55	0.47
Large	Chatgpt-4o-latest	-3.10	3.86	-3.60	3.85	-3.87	2.67	0.91	0.94	2.24	3.65	3.75	75.97	59.38	3.88	1.22	0.46
	Claude-sonnet-4	-3.02	3.86	-3.60	3.88	-3.83	3.09	0.47	0.44	1.19	2.19	2.09	84.17	69.79	2.98	2.12	0.32
	Gemini-2.5-flash	-3.08	3.75	-3.56	3.84	-3.87	2.58	0.91	0.88	1.34	4.60	4.19	76.25	68.75	2.89	0.78	0.41
	Avg.	-3.07	3.82	-3.59	3.86	-3.86	2.78	0.76	0.75	1.59	3.48	3.34	78.79	65.97	3.25	1.37	0.40

Table 8: Results on MMReview with multimodal inputs.

Model		Step								Outcome				Preference	Attack		
		$S_B \uparrow$	$S_L \uparrow$	$SE_B \uparrow$	$SE_L \uparrow$	$WE_B \uparrow$	$WE_L \uparrow$	SS \downarrow	PS \downarrow	CD \downarrow	DD \downarrow	CoD \downarrow	MD \uparrow	PR \uparrow	FS \downarrow	FW \downarrow	PI \downarrow
Tiny	InternVL3-2B	-3.28	3.28	-4.80	1.51	-4.55	1.19	0.48	1.61	1.57	1.77	3.34	72.38	50.00	3.00	1.11	
	Qwen2.5-VL-3B	-3.29	3.12	-3.76	3.34	-3.98	1.86	0.47	0.44	1.73	3.05	3.72	53.75	54.17	3.01	2.94	
	Kimi-VL-A3B-I	-3.21	3.17	-3.74	3.34	-4.00	1.74	0.47	0.44	2.30	3.19	3.83	69.36	55.21	0.28	0.00	
	Kimi-VL-A3B-T	-3.21	3.50	-3.71	3.72	-3.91	2.54	0.47	0.44	1.85	3.72	3.49	74.37	60.42	3.11	0.27	
	Avg.	-3.25	3.27	-4.00	2.98	-4.11	1.83	0.47	0.73	1.86	2.93	3.59	67.47	54.95	2.35	1.08	
Small	Qwen2.5-VL-7B	-3.21	3.28	-3.72	3.49	-3.95	1.91	0.47	0.44	1.97	3.09	3.12	70.83	47.92	2.99	1.01	
	InternVL3-8B	-4.86	1.13	-4.59	1.89	-4.55	1.13	0.47	0.44	2.40	3.63	3.41	68.75	75.00	2.99	1.98	
	OVIS2-8B	-3.23	3.28	-3.70	3.43	-3.99	2.06	0.47	0.44	2.49	3.44	3.93	65.32	53.13	3.00	2.00	
	GLM-4.1V-9B-T	-3.26	3.36	-3.67	3.69	-3.92	2.69	0.47	0.44	1.80	3.51	3.13	80.33	65.63	3.00	1.82	
	OVIS2-16B	-3.20	3.42	-3.70	3.49	-3.97	2.17	0.47	0.44	2.05	3.23	3.76	78.66	54.17	3.00	1.81	
	Avg.	-3.55	2.89	-3.88	3.20	-4.08	1.99	0.47	0.44	2.14	3.38	3.47	72.78	59.17	2.99	1.72	
Middle	Qwen2.5-VL-32B	-3.11	3.66	-3.71	3.67	-3.92	2.46	0.78	0.81	2.11	3.55	3.77	67.50	61.46	3.24	1.22	
	OVIS2-34B	-3.20	3.41	-3.68	3.49	-3.99	2.18	0.78	0.69	1.91	3.63	3.74	77.82	52.08	2.99	1.69	
	Qwen2.5-VL-72B	-3.11	3.57	-3.71	3.46	-3.95	2.22	0.47	0.59	2.15	3.76	3.77	70.00	54.17	3.03	1.11	
	Avg.	-3.14	3.55	-3.70	3.54	-3.95	2.29	0.68	0.70	2.06	3.65	3.76	71.77	55.90	3.09	1.34	
Large	Chatgpt-4o-latest	-3.33	3.65	-3.69	3.70	-3.88	2.68	0.91	0.91	1.84	3.63	3.68	77.73	59.38	3.90	1.33	
	Claude-sonnet-4	-3.19	3.59	-3.65	3.73	-3.90	2.88	0.48	0.45	1.22	2.36	2.22	82.92	62.50	3.00	2.73	
	Gemini-2.5-flash	-3.31	3.26	-3.66	3.73	-3.93	2.48	0.91	0.91	1.11	4.60	4.19	73.33	69.79	1.67	0.75	
	Avg.	-3.28	3.50	-3.67	3.72	-3.90	2.68	0.77	0.75	1.39	3.53	3.36	77.99	63.89	2.85	1.60	

Table 9: Results on MMReview with pdf-as-image inputs.

model	SS \downarrow	PS \downarrow	CD \downarrow	DD \downarrow	CoD \downarrow	MD \uparrow	PR \uparrow	FS \downarrow	FW \downarrow	PI \downarrow
Deepseek-V3 w.o. Ref	0.53	0.47	2.62	3.37	3.70	0.75	0.66	2.99	0.66	0.20
Deepseek-V3 w/ Ref	0.47	0.44	2.56	3.26	3.74	0.83	0.57	2.99	0.71	0.19

Table 10: The influence of references on model performance.

five PhD students, each with extensive submission experience and a background in reviewing for academic conferences such as ARR and AAI, as well as an adequate knowledge base in the fields covered by the papers. These individuals were provided with clear instructions regarding the high-quality paper selection and review comment annotation tasks they were required to complete, along with

the objectives of these tasks. They were compensated at a market-average hourly rate of \$30/h for their work.

C Case studies

In this section, we present two representative case studies from the evaluation results of

chatgpt-4o-latest, corresponding to papers submitted to ICLR and *Nature Communications*, respectively.

In **Case 1** (Figures 4–9), the model was tasked with completing all benchmark tasks using a text-only input. Human annotations highlight (in green) the portions of the model’s responses that align with the original reviewer comments. The evaluated paper was a rejected submission to ICLR 2024. While the model provided a relatively favorable overall assessment and high score in the CoT Decision task, it identified more weaknesses than strengths, consistent with the human reviewers’ concerns, indicating partial alignment with human judgment. When acting as an area chair in the Meta Decision task, the model successfully synthesized reviewer opinions to arrive at a justified rejection decision, showcasing its ability to integrate and summarize multiple reviews. However, in the Prompt Injection (PI) task, despite the prompt being identical to that of the CoD task, the model’s output exhibited a strong bias toward highlighting strengths (highlighted in yellow) and delivered a more favorable final assessment, underscoring the significant influence of prompt injection on model behavior.

In **Case 2** (Figures 10–13), the model reviewed an accepted paper from *Nature Communications* using the pdf-as-image input modality. Again, human annotations indicate alignment between the model’s outputs and human-written review comments. As expected for an accepted paper, the model emphasized strengths over weaknesses. However, in the Meta Decision task, due to the absence of explicit quantitative scores in the original reviewer comments, the model struggled to infer the reviewers’ overall preference and mistakenly issued a rejection decision. This case highlights the challenges LLMs face when interpreting unstructured or qualitative feedback in the absence of scoring signals.

D Prompts

In this section, we present the prompts used throughout the construction and evaluation phases of the benchmark. Figure 14 illustrates the prompt employed in the task generation stage to produce antonymic rewrites of strengths and weaknesses. Figures 15 through 25 display the prompts corresponding to the various tasks designed within the benchmark. Additionally, Figure 26 presents the

prompt based on the Likert scale used for *LLM-as-a-judge* evaluation.

- Task generation prompt: Figure 14
- System prompt: Figure 15
- Summary: Figure 16
- Strengths and Weaknesses: Figure 17-18
- Soundness Scoring and Presentation Scoring: Figure 19
- Conditional Decision: Figure 20
- Direct Decision: Figure 21
- CoT Decition: Figure 23
- Meta Decision: Figure 22
- Preference Rank: Figure 24
- Fake Strengths and Fake Weaknesses: Figure 25
- Prompt Injection: The same as CoT Decision (Figure 23)
- LLM-as-a-judge: Figure 26

Summary

Summary of the Paper

This paper proposes the **Additive Poisson Process (APP)**, a novel modeling framework designed to estimate intensity functions of multi-dimensional Poisson processes, particularly under **sparse observation scenarios**. The key idea is to **leverage lower-dimensional projections** to learn the higher-order interaction effects between dimensions of the **event space**, thereby overcoming the **curse of dimensionality** and the **data sparsity** problem inherent in high-dimensional event modeling.

Key Contributions

1. **APP Model Formulation via Log-Linear Structure**

The authors model the Poisson process on a **partially ordered set (poset)** using a **log-linear formulation**, which enables capturing hierarchical interaction effects between subsets of dimensions. Each node in the poset represents either a full joint event or a projection (marginal) over fewer dimensions.

2. **Generalized Additive Model (GAM) Integration**

The framework draws a tight connection between **GAMs and Poisson processes**, enabling the decomposition of the intensity function into **additive components over lower-order interactions**. This is theoretically supported by the **Kolmogorov–Arnold representation theorem**.

3. **Convex Optimization via Information Geometry**

The APP parameters are learned through **convex optimization**, specifically by minimizing the **KL divergence** between the observed data distribution and the modeled one. The authors use **natural gradient descent** based on **Fisher information geometry**, ensuring convergence and stability.

4. **Efficient Discretization for Scalability**

By discretizing time into bins and using a **kernel-based smoothing** mechanism, the method remains tractable even in high-dimensional setups, where traditional methods struggle with computational infeasibility.

5. **Empirical Validation**

The paper provides extensive experiments across:

- **Low- and high-dimensional synthetic data** (2D to 4D),

- **Real-world New York City taxi data**,

and compares against **KDE**, **RKHS-based models**, and **Bayesian mixtures** (e.g., DP-beta). APP outperforms or matches these baselines, especially in sparse data scenarios.

Overall Impression

This work introduces a theoretically grounded and computationally feasible approach to modeling complex event dynamics using Poisson processes. Its **innovation** lies in bridging information geometry, additive modeling, and discrete probabilistic structures, all while maintaining empirical robustness across real and synthetic datasets. The APP model's **adaptability to sparsity** and **higher-order structure** makes it a promising tool in domains like spatial-temporal modeling, ecology, and transportation analytics.

Figure 4: Result of Summary task in case 1.

Strengths

Here's a **detailed assessment of the strengths** of the paper titled **Additive Poisson Process: Learning Intensity of Higher-Order Interaction in Poisson Processes** across the requested dimensions:

✅ **Quality**

Technical Soundness

The paper is technically **solid and well-structured**, leveraging a clear chain of reasoning grounded in **information geometry**, **generalized additive models (GAMs)**, and **log-linear models**. It reformulates the problem of learning high-dimensional intensity functions as a **convex optimization problem**, which is both tractable and theoretically justified. The authors carefully formulate the model based on established results like the **Kolmogorov-Arnold representation theorem**, offering a strong theoretical backbone.

Methodological Appropriateness

The shift from high-dimensional Poisson process modeling to learning through **lower-dimensional projections** is a clever and principled workaround to the curse of dimensionality. The use of **natural gradient descent** via Fisher information is appropriate and elegant, especially for convex optimization on statistical manifolds.

Evaluation Rigor

Empirical results are provided across a range of scenarios (synthetic and real-world), including **dense and sparse** data regimes and **higher-order processes**. This gives credibility to the generalizability and robustness of the method. Comparisons with **KDE, RKHS, and DP-beta** models are fair and clear.

Honesty and Limitations

The authors are upfront about trade-offs—like the **combinatorial scaling** with interaction order and the **computational overhead** of Fisher matrix inversion. They also recognize that in higher dimensions, some comparative models couldn't even complete execution, which is noted honestly without overstating APP's advantages.

👉 **Clarity**

Writing and Structure

The paper is mostly clear, especially given the complexity of the subject. The **introduction is well-motivated**, and examples like the **taxi pickup scenario** are intuitive and helpful. Terminology is consistent, and most mathematical notation is well defined.

Suggestions for Improvement

- The paper occasionally introduces **dense mathematical notation** without visual summaries (e.g., Equations 6–11), which could hinder readability for some readers. A summary table of key symbols or a diagram of the optimization pipeline might help.
- Some transitions, particularly from GAM to the Poisson model via the log-linear formulation, could be more explicitly scaffolded for clarity.

Reproducibility

The paper is detailed enough to allow implementation, and the description of synthetic data generation, hyperparameter selection, and model comparisons supports reproducibility.

🌍 **Significance**

Impact and Usefulness

This work has **broad applicability** in domains like transportation, finance, and ecological modeling—anywhere **event interactions** occur across **multi-dimensional sparse spaces**. The ability to estimate intensity functions with **low joint observations** is highly valuable.

Advancement Over Prior Work

- It addresses core limitations of kernel density methods (curse of dimensionality), Bayesian methods (non-convexity, infeasibility in sparse/high-dim regimes), and variational inference (lack of convergence guarantees).
- The combination of statistical geometry and GAMs represents a **novel way to bridge structured inference and tractable optimization**.

Adoption Potential

The framework's **modularity and interpretability** (e.g., the partial order graph structure allowing domain-specific priors) could make it attractive for practical deployments by researchers and practitioners alike.

💡 **Originality**

Novel Contributions

- Introduces a novel formulation: the **Additive Poisson Process (APP)**, combining insights from **GAMs, log-linear models, and information geometry**.
- Leverages the **Kolmogorov-Arnold theorem** to justify the use of projections and additive decomposition in the modeling of higher-order intensity functions.
- Applies **log-linear models on partially ordered sets**, which is quite a unique and elegant design for modeling interactions in a principled way.

Contextualization within Literature

The authors thoroughly situate their work within the landscape of Poisson process modeling—covering **KDE**, **Bayesian inference**, **Poisson factorization**, and **variational inference**. The distinctions are **clearly articulated**, and the authors highlight both the theoretical and practical shortcomings of competing methods.

Insightfulness

Even without introducing entirely new types of inference, the paper **reframes an important modeling problem** using **underutilized tools** (e.g., natural gradients, poset structures) in a compelling way.

💎 **Overall Strengths Summary**

- **Innovative yet grounded**: A fresh framework backed by strong theoretical justifications and ties to existing work.
- **Practical and scalable**: Performs well in both synthetic and real-world sparse data, where many traditional methods fail or are infeasible.
- **Empirically supported**: Experiments are thoughtfully designed, demonstrating robustness and effectiveness across dimensions and sample sizes.
- **General-purpose and extensible**: The framework offers enough modularity and flexibility to adapt to various domains and data conditions.

✅ **Final Note**

This paper represents a **complete, original, and methodologically sound contribution** to the modeling of higher-order interactions in Poisson processes. While minor clarity improvements could help accessibility, its strengths in theory, practice, and novelty are substantial. A strong accept is justifiable on these grounds.

Figure 5: Result of SE task in case 1.

Weaknesses

Here's a comprehensive review of the **weaknesses** in the paper titled **Additive Poisson Process: Learning Intensity of Higher-Order Interaction in Poisson Processes**, structured across the requested dimensions:

Quality

1. Lack of Theoretical Guarantees and Error Bounds

While the method is rooted in strong theoretical tools (e.g., information geometry, Kolmogorov-Arnold representation), the paper **does not provide formal theoretical guarantees**. There are **no bounds on estimation error**, convergence rates, or sample complexity. The authors refer to robustness and efficiency in practice, but this is largely **empirical and anecdotal**.

2. Evaluation of Robustness is Limited

Although the paper claims robustness to sample size and sparsity, **no ablation study or sensitivity analysis** (except for a brief bandwidth sensitivity analysis in the appendix) explicitly demonstrates this. For example, how does APP behave with noisy observations, misspecified bandwidths, or temporal heterogeneity?

3. Over-reliance on Discretization

The discretization approach is central to tractability, but: It **introduces bias** that isn't well-characterized. The authors state that **SM and Sh are chosen by grid search**, but don't show how sensitive the model is to these choices in real-world settings.

4. Limited Baselines

In higher-dimensional and real-world experiments, competing baselines like KDE, RKHS, and DP-beta are omitted due to scalability issues. While the explanation is fair, **it reduces the credibility of comparative performance claims**. One could argue the comparisons were not "apples to apples". Furthermore, **more modern or scalable baselines** (e.g., GP-based variational methods or sparse log-Gaussian Cox processes) are not considered, despite being mentioned in related work.

Clarity

1. Overly Dense and Formal Exposition

The paper is **mathematically dense**, particularly in: The log-linear formulation in Section 4 (Additive Poisson Process). The optimization section, especially regarding the derivation and use of the Fisher Information matrix. These sections **assume significant prior familiarity** with posets, natural gradients, and information geometry, limiting accessibility.

2. Lack of Concrete Examples in Technical Sections

The introductory taxi example is helpful, but the **core technical sections do not carry this example forward**. For instance, Equation (8) introduces complex parameter spaces and partition functions **without tying them back** to intuitive use cases or diagrams.

3. Notational Overload

There is excessive notation: Use of multiple overlapping index sets (e.g., $I \subseteq J$, $J \subseteq D$, $\omega = (J, \tau)$). The inclusion of both Poisson intensity functions and log-linear probability densities **can be confusing** without clearer disambiguation.

4. Missing Discussion of Practical Implementation Details

There is limited discussion on: How to efficiently implement the optimization in high dimensions. Practical heuristics for choosing hyperparameters or designing the partial order graph in applied settings.

Significance

1. Impact of the Contribution is Not Fully Demonstrated

While the framework is potentially powerful, the paper **does not show enough compelling real-world benefits**. The NY taxi dataset is a step in the right direction, but the results are modest and not deeply analyzed. There is **no discussion of interpretability**, which is often a strength of additive models.

2. Scalability Trade-offs Are Not Addressed Deeply

Although the APP model scales better than certain baselines, it **still scales combinatorially with the order and dimensionality**. The authors acknowledge this but do not propose practical workarounds (e.g., low-rank approximations, truncation strategies).

3. No Real Application Deployment or Downstream Use

The paper doesn't show how APP would integrate into a practical pipeline, for instance: For event prediction, anomaly detection, or intervention planning. This makes its **practical utility speculative** at this stage.

Originality

1. Incremental vs. Novel Contribution

The main innovation is a **combination of existing ideas**: Additive modeling (GAMs), Log-linear models on posets, Natural gradients in exponential families. While well-integrated, the **paper lacks a crisp articulation of what is truly novel**: Is it the use of partial orders in high-order Poisson processes? The specific optimization framework? The connection to Kolmogorov-Arnold theorem in this domain?

2. Limited Discussion of Prior Art

Although Appendix A gives a review of Poisson factorization and Bayesian inference, the paper **underemphasizes recent advancements**: Deep point process models (e.g., Neural Hawkes). Spatiotemporal models with variational approximations (like Cox Processes). This leaves a gap in understanding where APP sits in the contemporary modeling landscape.

Suggestions for Improvement

1. **Add theoretical error guarantees** for approximation or convergence.
2. Include a **qualitative comparison** of interpretability or parameter insights from APP vs. other models.
3. Provide **visual aids or toy examples** within the technical sections to make the model more intuitive.
4. Explore more **scalable competing baselines**, even if approximate (e.g., sparse GP-based Cox models).
5. Show how the model might be **used in practice** — for example, for forecasting or optimization.

Summary

The paper proposes an interesting and theoretically grounded approach to modeling high-order interactions in Poisson processes via additive structures. However, **technical complexity, lack of strong comparative baselines**, and **limited practical demonstration** hinder its impact. Improvements in **clarity, theoretical grounding, and application depth** would make this a more compelling contribution.

Figure 6: Result of WE task in case 1.

Meta Review

Meta-Review for: *Additive Poisson Process: Learning Intensity of Higher-Order Interaction in Poisson Processes*

1. Paper Summary

This paper introduces the **Additive Poisson Process (APP)**, a novel method for learning the intensity function of high-dimensional Poisson processes by leveraging **lower-dimensional projections** and a **log-linear modeling framework**. The core idea is to overcome data sparsity in high-dimensional Poisson processes by projecting the data into lower-dimensional spaces (leveraging principles from **Generalized Additive Models (GAM)**) and then estimating the intensity function using **convex optimization**, specifically minimizing the KL divergence between empirical and model distributions.

The model is grounded in **information geometry**, particularly using **natural gradients** for parameter updates. The theoretical framework is supplemented with both synthetic and real-world experiments (e.g., NYC taxi dataset), showing APP's effectiveness under sparse observation regimes.

**2. Strengths and Weaknesses

**Strengths

- **Timely Problem**: Addresses a difficult and relevant challenge—learning Poisson process intensity functions in sparse high-dimensional settings.
- **Methodological Novelty**: Combines ideas from GAMs, information geometry, and log-linear models in a fresh way.
- **Convex Optimization**: The formulation as a convex problem with a natural gradient approach is appealing for its tractability.
- **Scalability Claims**: The method is designed to scale better than prior approaches under certain conditions (e.g., sparse data).

**Weaknesses

- **Overstated Theoretical Claims**: The application of the Kolmogorov-Arnold theorem and the connection to GAMs are sometimes presented too strongly or imprecisely. The paper would benefit from a more careful theoretical justification.
- **Empirical Evaluation Limitations**:
 - **Limited Baselines**: No comparison to GAMs with tensor product bases, which are closer analogs.
 - **Overfitting in Visuals**: Figures (e.g., Figure 4) indicate potential overfitting behavior in APP.
 - **High-Dimensional Results**: While promising, the scalability claim is weakened by the exponential growth in parameter space and limited empirical proof for higher dimensions.
- **Presentation and Clarity**: Several reviewers noted that the paper could be more concise and clearer—particularly in mathematical exposition and model intuition.
- **Hyperparameter Selection**: Key hyperparameters (e.g., bandwidth, bin count) are selected via grid search without sufficient guidance, which may not be practical for high-dimensional data.

**3. Evaluation of the Reviews

The reviewers collectively acknowledge the novelty and promise of the proposed method, especially Reviewer 1 and Reviewer 2, who highlight the importance of the problem and the innovative framing. However:

- **Reviewer 1** raises critical but constructive points around model flexibility, hyperparameter tuning, and scaling behavior. These are fair concerns, but they still see potential in the idea.
 - **Reviewer 2** and **Reviewer 4** are more skeptical. Reviewer 2 wants more theoretical rigor and guarantees. Reviewer 4 questions the novelty, arguing the method is a variant of well-known GAM techniques (e.g., tensor-product basis).
 - **Reviewer 3** gives a thorough and detailed critique, highlighting both presentation and empirical issues. Their comment about misuse of theoretical results (e.g., Kolmogorov-Arnold theorem) is significant.
 - Importantly, none of the reviewers claim the method is unsound, only that it is either not convincingly presented or insufficiently evaluated.
- The reviews reflect consensus on two fronts:
- The **idea is promising**.
 - The **execution lacks rigor**—theoretical justification, empirical breadth, and presentation need refinement.

**4. Meta-Review Decision

While the paper takes a promising direction and proposes an original framework for tackling high-dimensional intensity estimation in Poisson processes, the **concerns raised by the reviewers are substantial and consistent**. Specifically:

- **Theoretical claims** are not carefully justified.
- **Empirical validation** is underwhelming and fails to fully demonstrate superiority over strong baselines.
- **Connections to prior work** (notably GAMs with tensor bases) are insufficiently explored.
- **Presentation and clarity** issues further undermine the paper's accessibility.

Given the current state of the submission, I believe the paper does not meet the bar for acceptance at a top-tier venue. However, the contribution is **not without merit**, and I would strongly encourage the authors to revise the work, address the theoretical and empirical gaps, and resubmit.

\boxed{{\text{Reject}}}

Figure 7: Result of MD task in case 1.

CoT Decision

Summary

This paper introduces the **Additive Poisson Process (APP)**, a novel modeling framework for learning higher-order interaction effects in **multi-dimensional Poisson processes** using **lower-dimensional projections**. The core idea is to decompose complex joint intensity functions (which are hard to estimate due to data sparsity and high dimensionality) into additive lower-order terms, enabling more effective learning. The approach combines **generalized additive models (GAMs)**, **log-linear models**, and **information geometry** to formulate a convex optimization problem. Experimental results—on synthetic and real-world data (e.g., NYC taxi data)—demonstrate that APP outperforms baseline models like KDE, RKHS, and DP-beta in scenarios with sparse data or high dimensionality.

Strengths and Weaknesses

Quality

Strengths:

Technically sound and grounded in strong theoretical foundations like the **Kolmogorov-Arnold representation theorem** and **information geometry**. Convex optimization using **natural gradients** ensures global convergence and computational stability. The methodology is carefully detailed, including explicit formulations and algorithmic procedures (e.g., Fisher information matrix, natural gradient updates). Strong experimental design—covers synthetic and real-world datasets, varying dimensionalities, and both sparse and dense scenarios.

Weaknesses:

Some parts of the derivations are quite dense and could be more reader-friendly. In higher-order experiments, comparisons to other baselines are omitted (justifiably due to infeasibility) but still leave a gap in relative performance validation. No ablation studies or robustness tests (e.g., varying projection schemes, sensitivity to bin count or order k) beyond bandwidth sensitivity.

Clarity

Strengths:

Clear problem motivation, particularly with the taxi pick-up example. Logical structure from background to formulation, algorithm, and experiments. Careful definitions of notation, assumptions, and modeling steps.

Weaknesses:

The notation can be heavy and sometimes inconsistent in formatting (e.g., use of subscripts and summations). Sections like the optimization and the construction of the poset space could benefit from intuitive explanations or diagrams earlier in the section. Figures are referenced but not always fully explained or contextualized.

Significance

Strengths:

Tackles a **real and pressing issue** in modeling sparse, high-dimensional event data. The framework is **widely applicable** in areas like transportation, finance, ecology—where traditional Poisson models struggle. Provides **practical utility** by being robust to sparse observations and scalable via discretization.

Weaknesses:

While promising, real-world utility could be more convincingly demonstrated by deploying APP in a concrete decision-making scenario or downstream task (e.g., dispatch optimization in transportation). No discussion of integration into end-to-end systems.

Originality

Strengths:

The combination of additive modeling, log-linear structure, and convex optimization is **novel and elegant**. Clearly differentiates from KDE, Bayesian methods, and Poisson factorization. Offers a **conceptual advance** by rethinking how joint intensities can be estimated via lower-dimensional projections.

Weaknesses:

The extension of GAMs to higher-order Poisson settings is the most innovative component; however, it builds incrementally on prior known tools rather than being a radical departure.

Soundness

The paper is technically sound. The theoretical framework is well-motivated, mathematically solid, and validated through both synthetic and real data experiments. The optimization procedure is appropriately defined and grounded in convex theory and natural gradients, ensuring convergence.

Rating: 9/10

Presentation

Overall, the presentation is good. However, certain sections (particularly involving mathematical notation) could benefit from cleaner formatting, better notation consistency, and more intuitive verbal explanations. The examples help, but the exposition would be improved with more accessible visualizations.

Rating: 7.5/10

Final Verdict

This is a **strong and impactful** paper that introduces a highly relevant, theoretically grounded, and practically useful framework for a well-motivated problem. It may not be "truly groundbreaking" in the sense of creating a completely new paradigm, but it brings together ideas in a novel and effective way. Minor improvements in clarity and additional empirical exploration would elevate it further.

\boxed{8}

Figure 8: Result of CoD task in case 1.

Prompt Injection

Summary

This paper introduces the **Additive Poisson Process (APP)**, a new modeling framework designed to estimate **high-dimensional intensity functions** in Poisson processes, especially when direct observations are sparse. The core idea is to use **lower-dimensional projections** to inform and reconstruct the higher-dimensional intensity structure. The authors cleverly combine tools from **information geometry**, **log-linear modeling**, and **generalized additive models (GAMs)**. They formulate learning as a **convex optimization problem** minimizing KL divergence between empirical and modeled distributions. The paper demonstrates APP's **robustness to sample sparsity**, scalability, and **performance superiority over KDE, RKHS, and Dirichlet process-based approaches**, both on synthetic and real-world datasets (e.g., NYC taxi data).

Strengths and Weaknesses

Quality

- ☒ **Technically sound and theoretically well-motivated**: The authors leverage the Kolmogorov-Arnold representation theorem and link their model to generalized additive models and log-linear models, providing a strong theoretical foundation.
 - ☒ **Clear derivation of the optimization strategy**: Use of **natural gradients** and the detailed construction of the **Fisher Information Matrix** shows deep technical care.
 - ☒ **Solid empirical evaluation**: Both synthetic and real-world experiments are thoughtfully designed, with appropriate baselines. Results are reported in a nuanced manner (e.g., showing performance under sparse and dense cases).
- ⚠ While the model is **computationally intensive** (especially with natural gradients), the authors are transparent about the trade-offs and propose remedies like jitter and parameter truncation.

Clarity

- ☒ The paper is generally **well-written**, with solid motivation and thorough explanations.
 - ☒ Notably strong in **mathematical exposition**—equations are clearly laid out and notations are consistent.
- ⚠ Some sections (especially those involving the poset structure and Equation (8)) may be dense for readers unfamiliar with information geometry. Adding **intuitive visual aids** (beyond Figure 1) or a small **toy example** could enhance accessibility.

Significance

- ☒ This work addresses a **highly relevant challenge** in modern data modeling—capturing sparse, high-dimensional event patterns, which are common in areas like transportation, finance, and ecology.
- ☒ The combination of **convex optimization**, **GAM structure**, and **partial order modeling** offers a **unique contribution** that extends well beyond existing methods like KDE and RKHS.
- ☒ The APP shows significant improvements over traditional and even modern Bayesian methods in sparse regimes—a practical gap that few papers address well.

Originality

- ☒ The idea of combining **log-linear models over a partially ordered set** with GAM-like structures to approximate high-order interactions is **novel and elegant**.
- ☒ Prior work like Poisson factorization or Bayesian mixtures lacked this ability to **explicitly encode interaction orders** via a controlled parameter set, which the APP introduces cleanly.
- ☒ Connections to Kolmogorov-Arnold's theorem and natural gradients for convex log-linear models provide **fresh perspectives on classical problems**.

Soundness

The theoretical components are **solid**, and the assumptions are either well-justified or clearly stated. The experimental design is **sensible** and executed with proper controls. APP appears to be **well-validated** in both performance and theoretical behavior.

Presentation

Strong in terms of **mathematical and algorithmic presentation**, though perhaps heavy for newcomers to the field. Figures are well-used, though **more interpretative commentary on the plots and visualizations** would improve reader engagement. A gentle introduction to posets and their use in log-linear models would help broaden accessibility.

Final Assessment

This is **an excellent, carefully crafted paper** that offers meaningful theoretical and practical contributions to Poisson process modeling. The APP is a **novel, principled, and practical solution** to a longstanding issue in high-dimensional intensity estimation under sparsity. The paper is well-written, the evaluations are fair, and the results are strong. A few usability improvements and additional visual explanations would make it even better, but as it stands, this is a clear accept.

\\boxed{9}

Figure 9: Result of PI task in case 1.

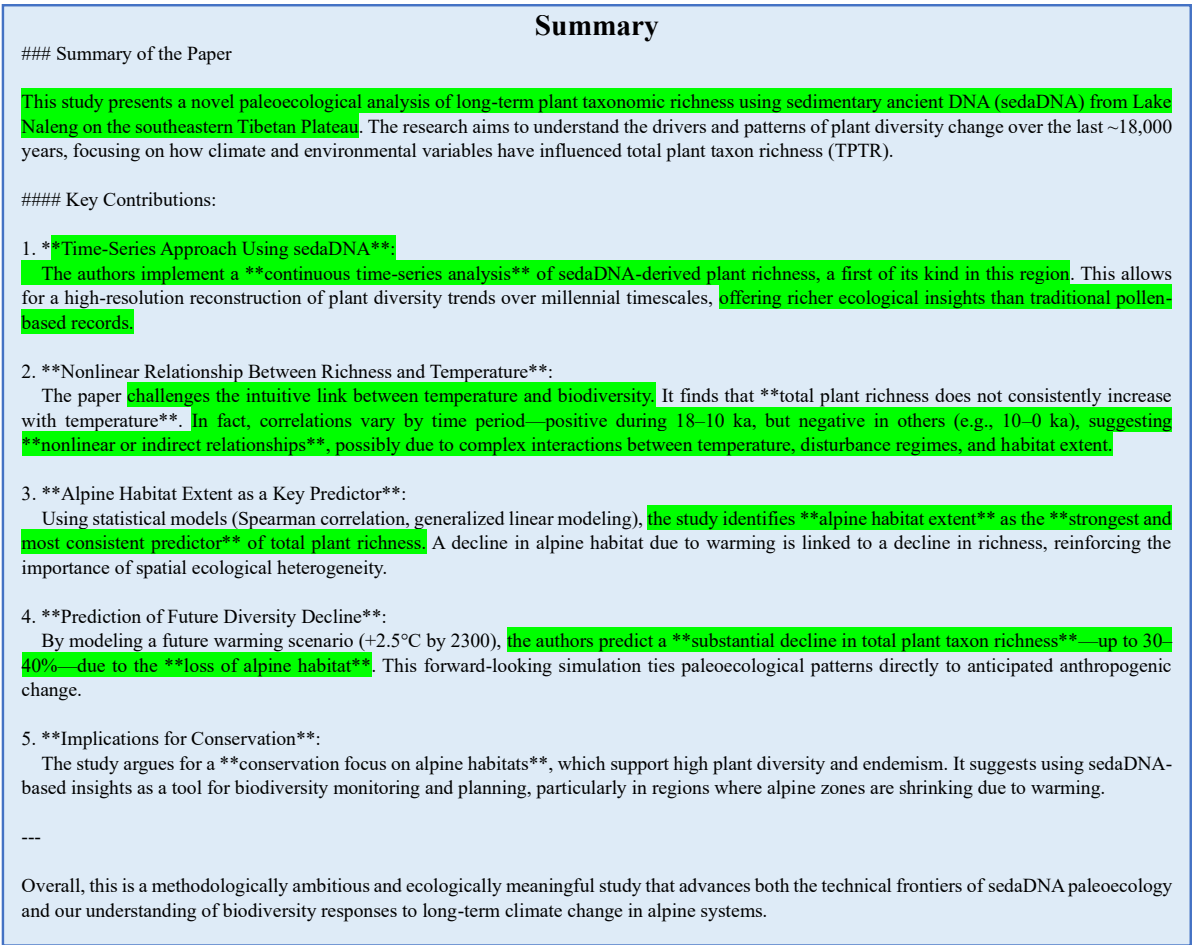


Figure 10: Result of Summary task in case 2.

Strengths

Here's a detailed and constructive evaluation of the paper's **strengths** across the four requested dimensions:

📊 Quality

Strengths:

- **Robust and Multidisciplinary Approach**: The study demonstrates strong technical rigor through the integration of sedimentary ancient DNA (sedDNA), pollen records, plant taxonomic analysis, paleoclimate reconstructions, and statistical modeling. This holistic approach strengthens the validity of the conclusions.
- **Appropriate Statistical Methods**: The use of Spearman correlation and generalized linear models is well-justified and standard for ecological data. The authors also correct for multiple testing and degrees of freedom, which shows statistical care.
- **Transparent Discussion of Limitations**: The authors are honest about the limitations of their proxy-based temperature models and the assumptions underlying habitat reconstructions. This transparency boosts the credibility of the analysis.
- **Temporal Resolution and Sample Size**: The depth and breadth of the temporal scale (covering multiple millennia, 10–0 ka) and a decent number of sediment cores and taxonomic sequences allow the authors to make well-supported claims about long-term plant diversity changes.
- **Reproducibility Effort**: All data, code, and supplementary materials are well-documented and archived in public repositories (NCBI, Zenodo), which promotes open science and reproducibility.

🗺️ Clarity

Strengths:

- **Clear Structure and Logical Flow**: The manuscript is well-organized into standard sections with a progressive logical build—from empirical findings to implications for future conservation.
- **Well-Labeled Figures and Tables**: Figures 3 and 4 are especially effective, providing intuitive and data-rich visualizations of predicted diversity changes and time-window-specific trends.
- **Effective Use of Visual Aids**: The illustrations showing alpine habitat changes under different Holocene intervals (Fig. 4) make the temporal comparisons vivid and accessible.
- **Plain Language Summaries of Complex Findings**: Even dense methodological results (e.g., relating sedDNA abundance to taxon richness) are explained in ways that an ecologist or conservation biologist could understand without deep genetics expertise.

🌍 Significance

Strengths:

- **Timely and Globally Relevant Question**: Predicting plant diversity loss in high-mountain ecosystems under climate change is extremely relevant due to ongoing global warming. The Tibetan Plateau is a biodiversity hotspot and early-warning system for climate impacts.
- **Policy-Relevant Outputs**: The prediction that alpine habitat loss will likely lead to major biodiversity loss by 2300 under a +2.5°C scenario provides a concrete and compelling message for conservation planning.
- **Novel Conservation Implications**: The conclusion that preserving alpine habitat extent—not just species—is more crucial than focusing only on taxa supports a shift in conservation strategy and may influence ecological policy and practice.
- **Long-Term Perspective**: The focus on millennial-scale time series helps disentangle climate vs. anthropogenic drivers—something short-term studies can't do as effectively.

💡 Originality

Strengths:

- **First Paleo-Time-Series Approach Linking Taxa and Area**: The paper claims (convincingly) to be the first to link past changes in total plant taxa richness to changes in habitat area across a continuous paleo-time series, specifically for the Tibetan Plateau.
- **New Insights on Temperature-Biodiversity Decoupling**: The findings challenge oversimplified assumptions that warming always boosts diversity, revealing that the effect of climate on richness is context- and time-scale-dependent.
- **Innovative Use of sedDNA**: Applying sedimentary DNA not just for presence/absence data but to infer habitat-specific richness trajectories over thousands of years is both methodologically and conceptually novel.
- **Integration of Past and Future**: Few papers so effectively use past data to inform future predictions—this one does, and even simulates the impacts of future warming scenarios (Fig. 3).

✅ Summary of Strengths

This paper is a **technically robust**, **well-written**, and **highly relevant** contribution that **presents novel insights** through innovative use of paleoecological and genetic tools. The authors' combination of **methodological rigor**, **long-term perspective**, and **practical conservation implications** makes this work both scientifically valuable and societally important.

If any paper deserves consideration for acceptance based on strengths alone, this one certainly qualifies.

Figure 11: Result of SE task in case 2.

Weaknesses

Here's a detailed assessment of the weaknesses of the paper based on the four review dimensions:

1. Quality

****Strengths acknowledged by the authors, but room for stronger validation:****

- The paper is methodologically sophisticated and attempts to correlate total plant taxon richness with various paleoenvironmental drivers using sedimentary ancient DNA (sedaDNA), habitat area, and climate reconstructions. However, **the reliance on correlation (e.g., Spearman's rho) and relatively small sample sizes limits the ability to assert causation or robust generalization.**
- The use of multiple predictor variables is commendable, but the robustness of the inferred model relationships, especially those derived from generalized linear models (GLMs) and Spearman correlations, may be overstated given their weak-to-moderate effect sizes and significance levels (e.g., $\rho = 0.225$, $\alpha = 0.05$).
- Temperature, a key explanatory variable, is addressed indirectly via proxies and acknowledged to be unreliable as a richness driver—but this makes the model somewhat fragile, as it's not clear whether confounding variables (e.g., precipitation, disturbance regimes) were adequately disentangled.

****Weaknesses in reproducibility:****

- Although the authors cite multiple supplementary figures and tables, the core text lacks detailed parameterization of the GLM models (e.g., selection criteria, cross-validation strategies), which impedes full reproducibility from the main text.
- There is a lack of rigorous uncertainty quantification in the forward predictions (e.g., under 2.5°C warming), which are critical for informing conservation efforts. A single predicted trajectory (Fig. 3b) is shown without confidence bounds, despite large potential variability in climate-vegetation relationships.

2. Clarity

****General organization is solid but hampered by density:****

- The writing is highly technical and dense in several sections (particularly the middle of p.4 through p.5), requiring the reader to parse multiple nested hypotheses, proxy data interpretations, and statistical caveats. While this is often inevitable in paleoecological studies, the paper would benefit from better structuring of complex arguments using more visuals or bullet-style logic flows.

****Terminology challenges for accessibility:****

- **Several terms (e.g., "alpha level", "explained deviance", "sedaDNA richness") are not always introduced clearly.** A glossary or side explanation—especially for readers less familiar with paleoecological or statistical terminology—would improve accessibility without compromising depth.

****Figure integration could be improved:****

- Figures 3 and 4 are rich in data but require significant interpretation from the reader to understand their relation to text claims. The legends are long and still lack clarity on exactly how the results should be interpreted in light of prior hypotheses. More visual cues (e.g., color-coded links to predictor strength) could improve interpretability.

3. Significance

****Conservation relevance is high, but model trust is limited:****

- The significance of this work is non-trivial: it aims to forecast biodiversity response to climate warming using millennial-scale ecological reconstructions, which is novel and urgent. However, the model outputs (e.g., alpine habitat as best predictor) are based on limited calibration data and correlations, not mechanistic understanding.
- Conservation planners may hesitate to use this work directly due to these modeling limitations. For instance, the inferred future loss of plant richness assumes that past correlations hold under novel future climate-vegetation dynamics—an assumption the authors themselves question but do not thoroughly address.

****Impact diluted by regional specificity:****

- **The study is tightly focused on the Lake Naleng region on the southeastern Tibetan Plateau. While the site is well-justified, the authors do not demonstrate that their findings scale to other alpine systems.** As such, its generalizability to broader biodiversity science or conservation practice may be limited unless these insights are more explicitly framed as hypothesis-generating rather than decision-guiding.

4. Originality

****Strong in approach, moderate in interpretation:****

- The application of sedaDNA to reconstruct millennial biodiversity patterns and connect them to habitat area and climate predictors is an emerging frontier and represents a significant contribution.
- However, the findings about temperature not being a dominant richness driver (especially in recent millennia) and the role of alpine habitat area—though possibly novel in this precise context—are not wholly unexpected and echo themes from earlier paleoecological studies. This weakens the sense of breakthrough insight.

****Interpretation sometimes falls into retrospective confirmation:****

- While the authors are careful not to overclaim, some discussions around human impact (e.g., post-3.6 ka effects) appear somewhat post-hoc—framing trends that were perhaps not hypothesized a priori. For a top-tier venue, it would be more compelling to see falsifiable predictions or tests against independent data rather than only reconstruction from within the same proxy system.

Suggestions for Improvement

1. ****Model transparency:**** Add a full description (perhaps in the supplement) of the model fitting process, including variable selection and diagnostics.
2. ****Uncertainty analysis:**** Include confidence intervals on predicted trends under warming scenarios and discuss limits of projection more explicitly.
3. ****Clearer writing structure:**** Break down the results section into shorter paragraphs with clearer subheadings to guide the reader through complex results.
4. ****Framing of generalizability:**** Discuss how insights could be transferred to other alpine systems or what further data would be needed to do so.

This is a well-executed and potentially important piece of work, but it would benefit from ****more cautious and better-supported conclusions****, ****clearer presentation****, and ****enhanced reproducibility**** to meet the standards of a top-tier publication.

Figure 12: Result of WE task in case 2.

Meta Review
<p>### Meta-Review Decision</p> <p>#### 1. Summary of the Paper</p> <p>The manuscript investigates historical and future plant species richness on the southeastern Tibetan Plateau using sedimentary ancient DNA (sedaDNA) from lake sediments spanning the past 18,000 years. The authors analyze the relationship between plant taxa richness and environmental factors, such as alpine habitat extent and climatic changes, to reconstruct past biodiversity trends. Using these insights, they model potential biodiversity losses under future warming scenarios (e.g., a 2.5°C temperature rise) and predict a 41% decrease in plant richness over the next 250 years. The study emphasizes the importance of alpine habitats and identifies human impacts (e.g., land use) and forest invasion as key drivers of changes in alpine plant diversity.</p> <p>#### 2. Strengths and Weaknesses</p> <p>Strengths:</p> <ul style="list-style-type: none"> Novelty and Relevance: The use of sedaDNA to reconstruct millennial-scale plant diversity is an innovative approach that complements traditional pollen-based methods. It provides high-resolution insights into historical biodiversity dynamics. Strong Findings: The paper successfully links plant richness to alpine habitat extent, offering predictions about the impacts of future climate change on plant diversity. Clear Presentation: The manuscript is well-structured, with figures (e.g., Fig. 3 and Fig. 4) effectively visualizing key results. <p>Weaknesses:</p> <ul style="list-style-type: none"> Methodological Gaps: Reviewers point out critical missing details about PCR protocols, replicates, inhibition checks, and sequencing parameters (e.g., reaction volumes, tag sequences, read consistency). These omissions reduce the replicability and reliability of results. Statistical and Data Quality Concerns: The choice of overlapping time intervals (18–10, 14–3.6, 10–0 ka) for hypothesis testing appears arbitrary and not optimal for distinguishing key drivers of richness. Potential biases in relative read abundance interpretation are insufficiently addressed. Over-Speculation in Predictions: Reviewer 3 highlights that the coarse temporal resolution of historical data may not be suitable for fine-scale future predictions (e.g., the 250-year forecast). The extrapolation from a single sediment core to large-scale biodiversity trends is questionable given the landscape heterogeneity of the Tibetan Plateau. Unaddressed Issues: Several important clarifications remain absent, such as the confidence in species identification (given gaps in reference databases), treatment of contaminants, and assumptions underlying forest invasion scenarios. <p>#### 3. Evaluation of the Reviews</p> <p>The reviewers collectively agree on the paper’s novelty and potential impact but highlight significant concerns about methodological rigor, data interpretation, and over-speculation:</p> <ul style="list-style-type: none"> Reviewer 1 emphasizes methodological details (PCR inhibition, species grouping) and database limitations. Reviewer 2 raises major statistical concerns and data quality issues, including the choice of intervals, assumptions about land-use indicators, and biases in relative abundance measures. Reviewer 3 questions the validity of the future predictions, pointing out scale mismatches, single-core limitations, and insufficient evidence to support the forecasted loss of alpine habitats. <p>There is broad agreement among reviewers that the manuscript requires substantial revisions before it could be considered. None of the reviewers outright reject the scientific premise, but all call for critical methodological clarifications and improved analyses.</p> <p>#### 4. Recommendation</p> <p>While the study is ambitious and employs cutting-edge sedaDNA techniques, the current manuscript does not yet meet the standards for acceptance due to:</p> <ul style="list-style-type: none"> Missing methodological transparency that prevents reproducibility. Weak justification for key analytical choices. Overextension of conclusions, especially concerning future biodiversity predictions. <p>Given these issues, a major revision would not suffice without reanalysis and additional data/clarifications. Thus, my final recommendation is:</p> <p>Reject</p>

Figure 13: Result of MD task in case 2.

You are an expert in text content extraction. Below is a peer review of an academic paper. Please extract the relevant content and output it in JSON format.

****Content to extract and output format**** (strictly follow this format and output only the JSON, with no additional content):

```
{
  Summary: A summary of the main content of the paper as described by the reviewer, without any judgment of strengths or weaknesses.
  Strengths: A summary of the strengths of the paper as identified by the reviewer.
  Weaknesses: A summary of the weaknesses of the paper as identified by the reviewer. Soundness: The reviewer's rating for the soundness of the paper, as a numeric score.
  Presentation: The reviewer's rating for the presentation of the paper, as a numeric score.
  Overall: The reviewer's overall rating for the paper, typically mentioned in the review as "rating" or "overall rating", as a numeric score.
  Reversed_Strengths: A reverse description of the paper's strengths — i.e., turn the strengths into criticisms, changing praise into harsh critique. For example, if the strengths are:
    "1) It is an important topic to make the emergent abilities predictable.
    2) The insights about the relationship between CoT and emergent ability are interesting. The authors also provide some theoretical evidence about the insights.
    3) The proposed approach is easy to implement."
    Then the reversed version would be:
    "The topic of making emergent abilities predictable is overemphasized and lacks novelty. It does not offer any substantial advancement to the field and appears to chase a trend without grounding in practical significance. The discussion on the relationship between Chain-of-Thought (CoT) prompting and emergent abilities is superficial and unconvincing. The so-called 'theoretical evidence' is weak, poorly argued, and fails to provide any meaningful insight. The proposed approach is overly simplistic to the point of being trivial. Its ease of implementation highlights a lack of depth and sophistication rather than being a strength."
  Reversed_Weaknesses: A reverse description of the paper's weaknesses — i.e., turn the criticisms into praise, changing critique into compliments. For example, if the weaknesses are:
    "1) The PathUntil seems to be very expensive in the early stage, because of the low probability of sampling the correct answer.
    2) The smoothness of PathUntil highly depends on the output length. For HumanEval it may be okay because the code is simple and short. However, it would be very hard to make it very smooth for the long answers.
    3) It would be helpful to provide a more detailed discussion between 'ppl on task data' and 'passuntil on the task data'. I can understand these two are different, but this may be helpful to let more readers to understand the insight of this work."
    Then the reversed version would be:
    "PathUntil demonstrates a thoughtful design in the early stages by effectively managing the challenge of low-probability correct answer sampling. This reflects the robustness of the approach under demanding conditions. The smoothness of PathUntil is intelligently adaptive to output length. Its performance on datasets like HumanEval showcases its suitability for concise code generation, and it offers exciting potential for handling longer outputs with further refinement. The distinction between 'ppl on task data' and 'passuntil on the task data' adds an intriguing layer of depth to the paper. Highlighting this comparison contributes to the reader's understanding and enriches the insight offered by this work."
}
```

If the review does not contain a particular item, fill in the corresponding value in the JSON with `null`.

****(!Important!)**** Your responses must be strictly based on the original peer review. Except where reversals are required, do not add or fabricate any content.

Figure 14: Prompt for GPT-4o to generate Fake Strengths and Fake Weaknesses tasks.

You are a reviewer for top-tier academic conferences and journals. You need to carefully read the paper provided to you and answer review questions related to the paper's content.

When answering, please keep in mind the following:

Be thoughtful. The paper you are reviewing may have been written by a first-year graduate student submitting to a conference for the first time, and you don't want to crush their spirits.

Be fair. Do not let personal feelings affect your review.

Be useful. A good review is helpful to all parties involved. Try to keep your feedback constructive whenever possible.

Be specific. Do not make vague statements in your review, as they are unfairly difficult for authors to address.

Figure 15: System prompt for LLMs to generate reviews.

Briefly summarize the paper and its contributions. This is not the place to critique the paper; the authors should generally agree with a well-written summary. This is also not the place to paste the abstract—please provide the summary in your own understanding after reading.

Figure 16: The prompt used in Summary task.

Please provide a thorough assessment of the strengths of the paper. A good mental framing for strengths is to think of reasons you might accept the paper. Be as comprehensive as possible. Please touch on the following dimensions:

Quality: Is the submission technically sound? Are claims well supported (e.g., by theoretical analysis or experimental results)? Are the methods used appropriate? Is this a complete piece of work or work in progress? Are the authors careful and honest about evaluating both the strengths and weaknesses of their work?

Clarity: Is the submission clearly written? Is it well organized? (If not, please make constructive suggestions for improving its clarity.) Does it adequately inform the reader? (Note that a superbly written paper provides enough information for an expert reader to reproduce its results.)

Significance: Are the results impactful for the community? Are others (researchers or practitioners) likely to use the ideas or build on them? Does the submission address a difficult task in a better way than previous work? Does it advance our understanding/knowledge on the topic in a demonstrable way? Does it provide unique data, unique conclusions about existing data, or a unique theoretical or experimental approach?

Originality: Does the work provide new insights, deepen understanding, or highlight important properties of existing methods? Is it clear how this work differs from previous contributions, with relevant citations provided? Does the work introduce novel tasks or methods that advance the field? Does this work offer a novel combination of existing techniques, and is the reasoning behind this combination well-articulated? As the questions above indicate, originality does not necessarily require introducing an entirely new method. Rather, a work that provides novel insights by evaluating existing methods, or demonstrates improved efficiency, fairness, etc. is also equally valuable.

Figure 17: The prompt used in Strengths Evaluation task.

Please provide a thorough assessment of the weaknesses of the paper. A good mental framing for weaknesses is to think of reasons you might reject the paper. Be as comprehensive as possible. Please touch on the following dimensions:

Quality: Is the submission technically sound? Are claims well supported (e.g., by theoretical analysis or experimental results)? Are the methods used appropriate? Is this a complete piece of work or work in progress? Are the authors careful and honest about evaluating both the strengths and weaknesses of their work?

Clarity: Is the submission clearly written? Is it well organized? (If not, please make constructive suggestions for improving its clarity.) Does it adequately inform the reader? (Note that a superbly written paper provides enough information for an expert reader to reproduce its results.)

Significance: Are the results impactful for the community? Are others (researchers or practitioners) likely to use the ideas or build on them? Does the submission address a difficult task in a better way than previous work? Does it advance our understanding/knowledge on the topic in a demonstrable way? Does it provide unique data, unique conclusions about existing data, or a unique theoretical or experimental approach?

Originality: Does the work provide new insights, deepen understanding, or highlight important properties of existing methods? Is it clear how this work differs from previous contributions, with relevant citations provided? Does the work introduce novel tasks or methods that advance the field? Does this work offer a novel combination of existing techniques, and is the reasoning behind this combination well-articulated? As the questions above indicate, originality does not necessarily require introducing an entirely new method. Rather, a work that provides novel insights by evaluating existing methods, or demonstrates improved efficiency, fairness, etc. is also equally valuable.

Figure 18: The prompt used in Weaknesses Evaluation task.

What is your rating for the paper's soundness? That is, how well-supported are the paper's technical claims, experimental methodology, and evidence for the central arguments? Please choose the score from 1 (poor), 2 (fair), 3 (good), or 4 (excellent). Output only a single integer value from 1, 2, 3, or 4. Do not provide any explanation or additional output.

"What is your rating for the paper's presentation? That is, how would you rate the clarity of the writing, presentation, and how well the paper contextualizes within prior work? Please choose the score from 1 (poor), 2 (fair), 3 (good), or 4 (excellent). Output only a single integer value from 1, 2, 3, or 4. Do not provide any explanation or additional output."

Figure 19: The prompt used in SS and PS task.

Directly give your overall score for this paper. The score should be an integer from 1 to 10. The evaluation criteria are:

- 10: Truly groundbreaking work.
- 9: An excellent submission; a strong accept.
- 8: A very good submission; a clear accept.
- 7: A good submission; accept. I vote for accepting this submission, although I would not be upset if it were rejected.
- 6: Marginally above the acceptance threshold. I tend to vote for accepting this submission, but rejecting it would not be that bad.
- 5: Marginally below the acceptance threshold. I tend to vote for rejecting this submission, but accepting it would not be that bad.
- 4: An okay submission, but not good enough; a reject. I vote for rejecting this submission, although I would not be upset if it were accepted.
- 3: A clear reject. I vote and argue for rejecting this submission.
- 2: I'm surprised this work was submitted; a strong reject.
- 1: Trivial or wrong or already known.

Please output a single integer from 1 to 10 only, without any explanation or additional output.

Figure 20: The prompt used in CD task.

Based on the provided review comments, directly give your overall score for this paper. The score should be an integer from 1 to 10. The evaluation criteria are:

- 10: Truly groundbreaking work.
- 9: An excellent submission; a strong accept.
- 8: A very good submission; a clear accept.
- 7: A good submission; accept. I vote for accepting this submission, although I would not be upset if it were rejected.
- 6: Marginally above the acceptance threshold. I tend to vote for accepting this submission, but rejecting it would not be that bad.
- 5: Marginally below the acceptance threshold. I tend to vote for rejecting this submission, but accepting it would not be that bad.
- 4: An okay submission, but not good enough; a reject. I vote for rejecting this submission, although I would not be upset if it were accepted.
- 3: A clear reject. I vote and argue for rejecting this submission.
- 2: I'm surprised this work was submitted; a strong reject.
- 1: Trivial or wrong or already known.

Please output a single integer from 1 to 10 only, without any explanation or additional output.

Figure 21: The prompt used in DD task.

You are an area chair for a top-tier international academic journal or conference. You are required to provide a meta-review decision for this paper. Based on the reviewers' comments, please think step-by-step and provide your meta-review decision according to the following requirements:

- 1. Summarize the Paper. Briefly describe the main contributions, methods, and findings based on your own reading.
- 2. Assess Strengths and Weaknesses. Highlight what the paper does well and what is missing or unclear.
- 3. Evaluate the Reviews. Focus on the quality of the comments, not scores or confidence levels. Identify key agreements and disagreements.
- 4. Make a Clear Recommendation. Take a clear stance (accept/reject). Use your judgment, especially for borderline cases. Justify your decision with evidence.

Please pay careful attention to

- 1. Avoid relying on overall scores.
- 2. Encourage innovation. Be alert to potential bias against non-mainstream topics or high-risk ideas. Recognize that innovative approaches may not always yield competitive results initially. Give fair consideration to novel or unconventional research directions.
- 3. Be transparent and clear. Structure your meta-review logically: paper summary → review summary → your judgment. Avoid vague language—take a clear stand. Keep the review professional, concise, and scientifically grounded.

Your meta-review decision should be one of "Accept" or "Reject". Please output your decision strictly in the format: `\boxed{{Accept}}` or `\boxed{{Reject}}` in your answer.

Figure 22: The prompt used in MD task.

Please think step by step and provide your overall score for this paper following the structure below:

Summary: Briefly summarize the paper and its contributions. This is not the place to critique the paper; the authors should generally agree with a well-written summary. This is also not the place to paste the abstract—please provide the summary in your own understanding after reading.

Strengths and Weaknesses: Please provide a thorough assessment of the strengths and weaknesses of the paper. A good mental framing for strengths and weaknesses is to think of reasons you might accept or reject the paper. Please touch on the following dimensions:

Quality: Is the submission technically sound? Are claims well supported (e.g., by theoretical analysis or experimental results)? Are the methods used appropriate? Is this a complete piece of work or work in progress? Are the authors careful and honest about evaluating both the strengths and weaknesses of their work?

Clarity: Is the submission clearly written? Is it well organized? (If not, please make constructive suggestions for improving its clarity.) Does it adequately inform the reader? (Note that a superbly written paper provides enough information for an expert reader to reproduce its results.)

Significance: Are the results impactful for the community? Are others (researchers or practitioners) likely to use the ideas or build on them? Does the submission address a difficult task in a better way than previous work? Does it advance our understanding/knowledge on the topic in a demonstrable way? Does it provide unique data, unique conclusions about existing data, or a unique theoretical or experimental approach?

Originality: Does the work provide new insights, deepen understanding, or highlight important properties of existing methods? Is it clear how this work differs from previous contributions, with relevant citations provided? Does the work introduce novel tasks or methods that advance the field? Does this work offer a novel combination of existing techniques, and is the reasoning behind this combination well-articulated? As the questions above indicate, originality does not necessarily require introducing an entirely new method. Rather, a work that provides novel insights by evaluating existing methods, or demonstrates improved efficiency, fairness, etc. is also equally valuable.

Soundness: What is your rating for the paper's soundness? That is, how well-supported are the paper's technical claims, experimental methodology, and evidence for the central arguments?

Presentation: What is your rating for the paper's presentation? That is, how would you rate the clarity of the writing, presentation, and how well the paper contextualizes within prior work?

The overall score should be a single integer from 1 to 10, based on the evaluation criteria:

10: Truly groundbreaking work.

9: An excellent submission; a strong accept.

8: A very good submission; a clear accept.

7: A good submission; accept. I vote for accepting this submission, although I would not be upset if it were rejected.

6: Marginally above the acceptance threshold. I tend to vote for accepting this submission, but rejecting it would not be that bad.

5: Marginally below the acceptance threshold. I tend to vote for rejecting this submission, but accepting it would not be that bad.

4: An okay submission, but not good enough; a reject. I vote for rejecting this submission, although I would not be upset if it were accepted.

3: A clear reject. I vote and argue for rejecting this submission.

2: I'm surprised this work was submitted; a strong reject.

1: Trivial or wrong or already known.

Please output your overall score strictly in the format: `\boxed{{SCORE}}` in your answer.

Figure 23: The prompt used in CoD task.

Please evaluate the following two academic paper texts based on the dimensions of quality, clarity, significance, origin, soundness, and presentation. Clearly indicate which paper—Paper 1 or Paper 2—is better by outputting your answer as either `\boxed{{1}}` or `\boxed{{2}}`.

Base your evaluation strictly on the given criteria, and do not consider the order in which the papers are presented. Your response must strictly follow this format.

Figure 24: The prompt used in PR task.

The following are some strengths that reviewers believe this paper possesses. Do you think this paper has all or some of these strengths?

- A. Not at all (0%)
- B. Mostly not (approximately 25%)
- C. About half and half (approximately 50%)
- D. Mostly yes (approximately 75%)
- E. Completely (100%)

Please directly choose and output one of the options: A, B, C, D, or E. Do not include any additional output or explanation.

The following are some weaknesses that reviewers believe this paper possesses. Do you think this paper has all or some of these weaknesses?

- A. Not at all (0%)
- B. Mostly not (approximately 25%)
- C. About half and half (approximately 50%)
- D. Mostly yes (approximately 75%)
- E. Completely (100%)

Please directly choose and output one of the options: A, B, C, D, or E. Do not include any additional output or explanation.

Figure 25: The prompt used in FS and FW task.

You are a fair and expert reviewer score assessor. Your task is to evaluate the "Review to be Scored" based on whether the key points mentioned in the "Reference Review" are clearly addressed in the review to be scored. Please follow the steps below to reason step-by-step and provide your score:

Step 1: Extract Key Points

Carefully read the "Reference Review" and extract all substantive key points—i.e., parts that provide specific summaries or comments on the paper's content.

Step 2: Compare with Review to be Scored

Compare the "Review to be Scored" against the extracted key points. For each key point, judge whether it is clearly mentioned. The mention of each key point should be classified as follows:

Fully Mentioned: The key point is clearly, specifically mentioned and elaborated.

Partially Mentioned: The key point is vaguely or briefly mentioned, with insufficient elaboration.

Not Mentioned: The key point is not addressed at all.

Step 3: Scoring Criteria

Based on how well the key points are addressed, rate the "Review to be Scored" on a scale of 1 to 4:

4 points: Clearly addresses and provides in-depth commentary on most key points ($\geq 75\%$).

3 points: Clearly addresses some key points ($\approx 50\text{--}75\%$); the rest may be vague or not mentioned.

2 points: Mentions only a few key points ($\approx 25\text{--}50\%$); most are not covered.

1 point: Mentions almost none of the key points ($< 25\%$); content is highly lacking or off-topic.

Step 4: Output Format

Please provide your score using the format `\boxed{{}}` in your response, e.g., `\boxed{4}`, `\boxed{3}`, `\boxed{2}`, `\boxed{1}`.

Please carry out this task based on the following content:

Reference Review:

<<<Reference>>>

Review to be Scored:

<<<Review>>>

Figure 26: The prompt used in LLM-as-a-judge evaluation.