

Evaluating the Impact of Reviewer Guideline Design on LLM-Based Automated Peer Review

Anonymous ACL submission

Abstract

Peer review is an essential process in scientific research, yet the growing workload has made its automation increasingly necessary. In this study, we analyze how different types of *reviewer guidelines*, such as official conference guidelines and reviewer-imitating ones generated from high-quality human reviews using LLMs, affect automated peer review. Our experiments show that official conference guidelines produce review results most consistent with human judgments, suggesting that evaluation criteria refined through conference practice serve as effective guidance for automated reviewing as well. In contrast, reviewer-imitating guidelines were generally less effective than official conference guidelines. Furthermore, enforcing strict rubric-style scoring consistently degraded performance, highlighting the importance of allowing subjective and holistic scoring.

1 Introduction

In recent years, the explosive growth in submissions to major machine learning conferences has placed an enormous burden on peer reviewers (Stelmakh et al., 2020; Zhang et al., 2022; Wei et al., 2025; Kim et al., 2025). The current human-based peer review system faces serious challenges including reviewer overload, inconsistent judgments, superficial evaluations, and delayed feedback. At the same time, as research automation advances, the importance of accurate evaluation has only increased, making the automation of the peer review process an urgent necessity.

One clear illustration comes from the ARR guidelines, which ask reviewers to “*Check for common review issues*”¹. These items do not simply caution reviewers against mistakes, but establish concrete criteria for judging the scientific contribution of a paper: for example, a result should not be

undervalued merely because it appears “unsurprising,” nor should a contribution be dismissed solely because it contradicts prior beliefs. By formalizing such principles, ARR effectively provides operational indicators that define what counts as sound, novel, and impactful research. Such directions illustrate how, as noted by Seeber (2020), reviewer guidelines shape the nature and quality of peer review rather than serving as administrative notes.

While several studies have explored automated peer review using Large language models (LLMs) (Zhuang et al., 2025; Shin et al., 2025; Chen et al., 2025; Lin et al., 2025; Li et al., 2025; Weng et al., 2025), the impact of guidelines and rubrics that govern review policies on automated review performance has not been sufficiently investigated. Guidelines direct reviewers’ attention to particular aspects and perspectives when evaluating a paper, and thus are crucial determinants of review quality. However, it remains unclear what types of guidelines are effective, and what types are ineffective, for automated peer review.

Broadly, two approaches can be considered for guideline design in automated peer review. One is to directly provide LLMs with existing, human-written guidelines established for each conference (e.g., official reviewer guideline from ARR) and have them conduct reviews. The other is a reviewer-imitating approach, where high-quality human-written reviews are used as input to a guideline generator (such as an LLM), which then extracts recurring evaluation patterns and key perspectives to generate new guidelines.

In this study, we systematically investigate how review quality is affected when LLMs conduct reviews under different guideline conditions. Specifically, we address the following research questions:

RQ1: Does providing guidelines improve the performance of automated peer review? **A:** We found that providing guidelines improves the

¹<https://aclrollingreview.org/reviewerguidelines>

performance of automated peer review. In particular, using official conference guidelines, such as those from NeurIPS and ARR, proved to be more effective. This suggests that evaluation criteria refined over time by the research community serve as a strong foundation for automated reviewing as well (§3).

RQ2: Are reviewer-imitating guidelines effective? In particular, is it better for models to determine scores subjectively, or to compute them strictly using rubric-style criteria? **A:** We found that reviewer-imitating guidelines were generally less effective than official conference guidelines. In particular, rubric-style formulations further reduced alignment with human judgments. This suggests that allowing LLMs to determine scores subjectively may offer better alignment with human judgment in specific contexts, rather than enforcing rigid, additive scoring logic (§4).

2 Problem Formulation

Types of Reviewer Guidelines This study analyzes two approaches for automated peer review: (1) using human-created guidelines (e.g., official ICLR reviewer guidelines) directly provided to LLMs for reviewing, and (2) generating guidelines from high-quality human reviews (filtered based on criteria in Appendix A) using LLMs and then using these generated guidelines for review.

Guideline-Based Automated Review Both approaches follow a common input–output structure for automated review, as shown in the prompt template in Appendix (Listing 6). The **input** consists of: (i) a paper in text format², and (ii) review guidelines (either human-created or LLM-generated). The **output** is a numerical review score (typically on a 1-10 scale) along with textual justification.

Evaluation To evaluate the quality of automated reviews, we compute the root mean squared error (RMSE) between LLM-generated scores and human review scores for the same papers. A lower RMSE indicates that the model’s ratings are numerically closer to human ratings. This metric allows us to assess how effectively an automated system approximates human peer review (Wei et al., 2025).

²Converted from PDF with references and appendices removed

Model	No-Guideline	ARR	NeurIPS	ICLR
<i>Official conference guidelines</i>				
Qwen3-30B-A3B	3.06	3.18	2.07	<u>2.91</u>
DeepSeek-R1-32B	3.86	3.27	2.72	<u>2.91</u>
Qwen3-32B	3.58	3.32	2.51	<u>2.89</u>
<i>Guidelines rewritten into instruction-style prompts</i>				
Qwen3-30B-A3B	3.06	2.97	3.08	<u>2.99</u>
DeepSeek-R1-32B	3.86	<u>3.06</u>	3.18	3.05
Qwen3-32B	3.58	<u>3.16</u>	3.32	3.09

Table 1: RMSE of automated review scores under No-Guideline and official conference guidelines. Upper block shows results using *official reviewer guidelines*, while lower block uses the same content rewritten into *instruction-style prompts* (converted with ChatGPT). Best results are in **bold**, second-best are underlined. **Lower values are better**, indicating closer alignment with human scores.

3 Effectiveness of Conference Reviewer Guidelines for Automated Review

This section evaluates how well reviewer guidelines improve automated peer review. We address whether providing guidelines improves accuracy and which venue’s guidelines are most effective.

Setup We randomly sampled 1,500 submissions (papers and reviews) from the ICLR 2024 peer-review dataset³. For each paper, the human reference label is the average of its reviewers’ overall rating. We evaluate three models: Qwen3-30B-A3B (MoE model)⁴(Yang et al., 2025), DeepSeek-R1-Distill-Qwen3-32B⁵(DeepSeek-AI et al., 2025), and Qwen3-32B⁶(Yang et al., 2025). “No-Guideline” denotes prompting the model with the paper only (no guideline) to output an integer overall rating, as shown in Appendix B. We also evaluate using official reviewer guidelines from ICLR⁷/NeurIPS⁸/ARR⁹. Because official documents often include non-evaluative administrative notes and explanatory prose, we additionally test

³We collected data from ICLR 2024 submissions on OpenReview, including 7,262 paper PDFs and 28,028 open-access reviews from <https://openreview.net/group?id=ICLR.cc/2024/Conference>

⁴<https://huggingface.co/Qwen/Qwen3-30B-A3B>

⁵<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen3-32B>

⁶<https://huggingface.co/Qwen/Qwen3-32B>

⁷<https://iclr.cc/Conferences/2024/ReviewerGuide>

⁸<https://neurips.cc/Conferences/2024/ReviewerGuidelines>

⁹<https://aclrollingreview.org/reviewerguidelines>

an instruction-style variant where the same evaluative content is rewritten as imperative, bullet-pointed instructions with explicit output requirements converted using ChatGPT 5. The prompts used for all guideline-based conditions is shown in [Appendix C](#).

Q1: Do guidelines improve automated reviewing? We compare No-Guideline to runs with conference guidelines. As shown in [Table 1](#), RMSE consistently decreases across all models when guidelines are provided. For example, DeepSeek-R1-32B demonstrates a substantial reduction in error across all three conference venues compared to the No-guideline baseline. Similarly, the instruction-style prompts outperform the baseline for all tested models, indicating that explicit evaluative criteria improve agreement with human judgements.

Q2: Which conference guidelines are more effective? We next compare ICLR/NeurIPS/ARR against each other. With the official guideline text, all three models achieve the lowest RMSE under NeurIPS, followed by ICLR, then ARR (e.g., for Qwen3-30B-A3B: $2.07 < 2.91 < 3.18$). When the same content is converted into instruction-style prompts, the best venue varies by model: Qwen3-30B-A3B prefers ARR (2.97), while DeepSeek-R1-32B and Qwen3-32B prefer ICLR (3.05 and 3.09, respectively). In other words, beyond the guideline content, the presentation (verbatim policy text vs. instruction-style prompt) can alter model behavior and flip the relative ranking among venues.

4 Effectiveness of Reviewer-imitating Guidelines in Automated Peer Review

In addition to using human-written conference reviewer guidelines, another approach is to *construct* guidelines with an LLM. In this section, we examine whether *reviewer-imitating guidelines*, generated from human-written reviews, can improve automated peer review. We evaluate review quality using several indicators (such as review length) and extract common perspectives from high-quality reviews, which are then converted into guidelines by an LLM. Through this filtering process, we divide reviews into three groups: *Good Review*, *Middle*, and *Bad Review*. We generate guidelines from each group and compare automated review performance under these different settings.

Overview of Guideline Generation We generate reviewer-imitating guidelines and apply them to

automated peer review through the following steps:

Listing 1: Example of reviewer-imitating guideline extracted from high-quality (“Good”) reviews.

```
Provide the following based on the Peer Review Behavior Checklist:

peer_review_behavior_checklist:
  good_paper:
    contribution_and_impact:
      - question: "Does the paper tackle a significant and timely problem?"
        description: "Reviewers positively note papers that address important and relevant research topics, seeing it as a sign of potential impact."
      - question: ...
    clarity_and_presentation: ...
  bad_paper:
    novelty_and_scholarship:
      - question: "Is the contribution incremental or a re-implementation of existing ideas?"
        description: "This is a frequent and critical flaw. Reviewers are highly knowledgeable and will reject papers they deem to be minor variations of prior work, often citing the specific papers."
      - question: ...
    clarity_and_soundness: ...
```

- **STEP1: Review Data Filtering:** To ensure guidelines are distilled from thorough reviews, we filter data using programmatic proxies: Length, Submission Date, Citations, and Consistency. While manual expert annotation is a gold standard, these structural metrics provide a scalable and reproducible foundation for assessing reviewer engagement and thoroughness. Reviews are grouped into *Good*, *Middle*, and *Bad* based on a composite score (3-17 points).

- **STEP2: Guideline Generation:** We provide the filtered reviews from each group (30 reviews) to Gemini 2.5 Pro ([Google Gemini Team, 2025](#)) to generate reviewer-imitating guidelines. The resulting guidelines are structured as a checklist, with separate items for positive aspects (good papers) and negative aspects (bad papers). They capture frequently mentioned criticisms, recurring evaluation criteria, and characteristic perspectives. For example, [Listing 1](#) shows an actual generated guideline (excerpted for brevity). Further details are described in [§A.3](#).

- **STEP3: Automated Review Using Guidelines:** This step follows the same procedure

Model	Setting	Base.	Good	Mid.	Bad
Qwen3-30B-A3B	w/o Rubric	3.06	3.22	3.24	<u>3.14</u>
	w/ Rubric	N/A	4.42	4.10	<u>4.28</u>
DeepSeek-R1-32B	w/o Rubric	3.86	3.13	<u>3.12</u>	3.06
	w/ Rubric	N/A	<u>4.26</u>	3.92	<u>4.26</u>
Qwen3-32B	w/o Rubric	3.58	3.48	<u>3.51</u>	3.69
	w/ Rubric	N/A	<u>4.57</u>	4.31	4.62

Table 2: RMSE (lower values are better) under reviewer-imitating guidelines, comparing settings with and without rubric. “Good” and “Bad” denote guidelines distilled from *high-quality* and *low-quality* human reviews, respectively.

as in the conference guideline experiments. The generated reviewer-imitating guidelines are given to the LLMs together with the paper text. The model then outputs an overall rating on a 1–10 scale and gives a brief justification.

Q1: Do reviewer-imitating guidelines outperform No-Guideline? We compare No-Guideline (the model uses only the paper) against reviewer-imitating guidelines from the Good/Bad groups. The impact of reviewer-imitating guidelines is model-dependent (see Table 2 for the comparison with the No-Guideline baseline). While DeepSeek-R1-32B shows a notable improvement in alignment, other models such as Qwen3-30B-A3B experience performance degradation when using guidelines distilled from human reviews. This suggests that reviewer imitation does not yield the same uniform gains observed with official conference guidelines.

Q2: How do reviewer-imitating guidelines compare to conference guidelines? We compare reviewer-imitating guidelines (Good/Bad) to official conference reviewer guidelines (ICLR, NeurIPS, ARR) under the same evaluation setup. Cross-referencing Table 1, conference guidelines consistently achieve lower RMSE than reviewer imitation. For example, with Qwen3-30B-A3B, Good = 3.22 and Bad = 3.14, while NeurIPS = 2.07 and ICLR = 2.91. These results suggest that the structural clarity and explicit requirements in official guidelines calibrate models more effectively than behavior-only signals distilled from reviews.

Q3: Does rubricizing reviewer-imitating guidelines help? We compare free-form (non-rubric) instructions to rubric-style scoring for reviewer-imitating guidelines. In the rubric setting, the LLM

judges whether each checklist item (Listing 1) in the “good_paper” section is satisfied and assigns +1 point if so, while items in the “bad_paper” section contribute -1 point. The total score is then normalized to a 1–10 scale and used as the overall rating. As shown in Table 2, the non-rubric setting consistently outperforms rubric-style formulation across all models and quality groups. In all cases, enforcing rigid point-allocation system significantly increased RMSE compared to allowing the model to provide free-form justifications.

5 Related Work

Automatic peer review research can be broadly categorized into two approaches: learning-based and prompt-based. The learning-based approach trains models to generate reviews or predict scores using task-specific supervision and intermediate reasoning. Early studies generated explainable reviews from structured evidence built through knowledge graphs (Gao et al., 2020). More recent work employs multi-stage pipelines that emulate expert analytical processes to train end-to-end review models (Zhu et al., 2025). In addition, recent surveys have organized datasets and tasks, helping to systematize what is learned and how models are evaluated (Yuan et al., 2022; Dycke et al., 2023). The prompt-based approach guides large language models during inference using rubrics or instructions. Studies on prompt-based automated reviewing have shown that well-designed evaluation criteria improve the fidelity and granularity of judgments (Kim et al., 2023). Benchmarking frameworks further standardize prompt formats and evaluation procedures (Zheng et al., 2023). These studies collectively highlight the importance of explicit instructions for LLMs. Building on this insight, we analyze how using reviewer guidelines, which serve as instructions to human reviewers, as prompts affects the quality of automated peer review. A related study (Kirtani et al., 2025) uses existing conference guidelines to generate and evaluate LLM-generated reviews. In contrast, we not only compare official guidelines from multiple venues but also propose generating reviewer-imitating guidelines from human reviews using LLMs, and systematically evaluate how these different guideline types affect scoring alignment with human judgments

335 Limitations

336 This study is limited in several aspects. First, we
337 only analyze review data written in English, as the
338 ICLR 2024 dataset contains exclusively English
339 reviews. Consequently, our findings may not gen-
340 eralize to non-English scientific contexts where
341 linguistic structures and cultural factors influence
342 evaluation criteria.

343 Furthermore, our evaluation relies primarily on
344 Root Mean Squared Error (RMSE) between LLM-
345 generated scores and human averages. While
346 RMSE effectively measures numerical alignment,
347 it serves as a narrow metric that does not fully
348 capture the quality, helpfulness, or logical depth
349 of the textual justifications. A significant limita-
350 tion of this work is the absence of systematic tex-
351 tual analysis, such as measuring Review Length
352 or Vocabulary Richness (type-token ratio), which
353 would provide a more multidimensional view of
354 how guidelines influence the granularity and lin-
355 guistic diversity of automated reviews. Addition-
356 ally, our "No-Guideline" baseline exhibits a ten-
357 dency toward score centralization, which could
358 potentially mask performance nuances in extreme
359 scoring ranges; future work should examine the full
360 distribution of predictions to mitigate this effect.

361 Finally, while our reviewer-imitating guidelines
362 are derived from high-quality human reviews, their
363 effectiveness depends on the specific logic and po-
364 tential biases of the generator model (Gemini 2.5
365 Pro) and the structural metrics used for filtering.
366 Future research should explore cross-model guide-
367 line generation and human-in-the-loop validation
368 to ensure that automated criteria remain aligned
369 with evolving scientific standards.

370 References

371 Shiping Chen, Duncan Brumby, and Anna Cox. 2025.
372 [Envisioning the future of peer review: Investigating](#)
373 [llm-assisted reviewing using chatgpt as a case study.](#)
374 *In Proceedings of the 4th Annual Symposium on*
375 *Human-Computer Interaction for Work, CHIWORK*
376 *'25*, New York, NY, USA. Association for Computing
377 Machinery.

378 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,
379 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
380 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,
381 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-
382 hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.
383 2025. [Deepseek-r1: Incentivizing reasoning capa-](#)
384 [bility in llms via reinforcement learning.](#) *Preprint*,
385 arXiv:2501.12948.

Nils Dycke, Iliia Kuznetsov, and Iryna Gurevych. 2023. [NLPeer: A unified resource for the computational study of peer review.](#) *In Proceedings of ACL 2023 (Long Papers)*. 386
387
388
389

Xue Gao and 1 others. 2020. [Reviewrobot: Explainable paper review generation based on knowledge graphs.](#) *In Proceedings of INLG 2020*. 390
391
392

Google Gemini Team. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities.](#) *arXiv preprint arXiv:2507.06261*. Technical Report / preprint. 393
394
395
396
397

Jaeho Kim, Yunseok Lee, and Seulki Lee. 2025. [Position: The ai conference peer review crisis demands author feedback and reviewer rewards.](#) *Preprint*, arXiv:2505.04966. 398
399
400
401

Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. [Prometheus: Inducing fine-grained evaluation capability in language models.](#) *Preprint*, arXiv:2310.08491. 402
403
404
405
406

Chhavi Kirtani, Madhav Krishan Garg, Tejash Prasad, Tanmay Singhal, Murari Mandal, and Dhruv Kumar. 2025. [Revieweval: An evaluation framework for ai-generated reviews.](#) *CoRR*, abs/2502.11736. 408
409
410
411

Ruochi Li, Haoxuan Zhang, Edward Gehringer, Ting Xiao, Junhua Ding, and Haihua Chen. 2025. [Unveiling the merits and defects of llms in automatic review generation for scientific papers.](#) *Preprint*, arXiv:2509.19326. 412
413
414
415
416

Tzu-Ling Lin, Wei-Chih Chen, Teng-Fang Hsiao, Hou-I Liu, Ya-Hsin Yeh, Yu Kai Chan, Wen-Sheng Lien, Po-Yen Kuo, Philip S. Yu, and Hong-Han Shuai. 2025. [Breaking the reviewer: Assessing the vulnerability of large language models in automated peer review under textual adversarial attacks.](#) *Preprint*, arXiv:2506.11113. 417
418
419
420
421
422
423

Marco Seeber. 2020. How do journals of different rank instruct peer reviewers? reviewer guidelines in the field of management. *Scientometrics*, 122(3):1387–1405. 424
425
426
427

Hyungyu Shin, Jingyu Tang, Yoonjoo Lee, Nayoung Kim, Hyunseung Lim, Ji Yong Cho, Hwajung Hong, Moontae Lee, and Juho Kim. 2025. [Mind the blind spots: A focus-level evaluation framework for llm reviews.](#) *Preprint*, arXiv:2502.17086. 428
429
430
431
432

Ivan Stelmakh, Nihar B. Shah, Aarti Singh, and Hal Daumé III. 2020. [A novice-reviewer experiment to address scarcity of qualified reviewers in large conferences.](#) *Preprint*, arXiv:2011.15050. 433
434
435
436

Qiyao Wei, Samuel Holt, Jing Yang, Markus Wulfmeier, and Mihaela van der Schaar. 2025. [The ai imperative: Scaling high-quality peer review in machine learning.](#) *Preprint*, arXiv:2506.08134. 437
438
439
440

441 Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo
442 Zhang, Jindong Wang, Yue Zhang, and Linyi
443 Yang. 2025. [Cycleresearcher: Improving auto-](#)
444 [mated research via automated review](#). *Preprint*,
445 arXiv:2411.00816.

446 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
447 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
448 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-
449 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao
450 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40
451 others. 2025. [Qwen3 technical report](#). *CoRR*,
452 abs/2505.09388.

453 Weizhe Yuan, Pengfei Liu, and Graham Neubig. 2022.
454 [Can we automate scientific reviewing?](#) *Journal of*
455 *Artificial Intelligence Research*, 75:171–212.

456 Yichi Zhang, Fang-Yi Yu, Grant Schoenebeck, and
457 David Kempe. 2022. [A system-level analysis of con-](#)
458 [ference peer review](#). In *Proceedings of the 23rd ACM*
459 *Conference on Economics and Computation*, EC ’22,
460 page 1041–1080, New York, NY, USA. Association
461 for Computing Machinery.

462 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
463 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
464 Zhuohan Li, Dacheng Li, Eric P. Xing, and 1 others.
465 2023. [Judging LLM-as-a-judge with MT-bench and](#)
466 [chatbot arena](#). *Preprint*, arXiv:2306.05685.

467 Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang.
468 2025. [Deepreview: Improving llm-based paper re-](#)
469 [view with human-like deep thinking process](#). In *Pro-*
470 *ceedings of ACL 2025 (Long Papers)*.

471 Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu,
472 Yuwen Jiang, and Jialiang Lin. 2025. [Large lan-](#)
473 [guage models for automated scholarly paper review:](#)
474 [A survey](#). *Preprint*, arXiv:2501.10326.

475 A Detailed Description of the 476 Reviewer-imitating Guideline

477 A.1 Review Data Collection and Filtering

478 We collected data from ICLR 2024 submissions
479 on OpenReview, including 7,262 paper PDFs and
480 28,028 open-access reviews as mentioned in § 3.
481 We selected ICLR due to the abundance of open-
482 access data and its status as a representative confer-
483 ence in the machine learning field. Table 3 shows
484 an example of review data used in our study. Each
485 review includes the following elements: review
486 summary, item-specific scores (soundness, presen-
487 tation, contribution), strengths, weaknesses, ques-
488 tions, overall score, and confidence level. However,
489 not all reviews are of high quality. Some reviews
490 are too short, have unclear evidence, or have incon-
491 sistent scores and content. To generate guidelines
492 from high-quality reviews, it is necessary to filter
493 the review data as preprocessing. Filtering uses

the following five criteria, based on our intuitions
about what constitutes high-quality reviews.

1. **Submission Date:** Measures deviation from
the average review submission date. We hy-
pothesize that reviewers who submit reviews
promptly are more conscientious and likely to
produce higher quality reviews.
2. **Review Length:** Evaluation based on char-
acter count, with longer reviews assumed to
indicate deeper engagement with the paper.
We calculate scores based on multiples of the
median length.
3. **Reference Quality:** Evaluates mention of one
or more references, pointing out insufficient
related work, and critical discussion of prior
work (0-3 points). We expect that reviewers
with higher domain expertise, evidenced by
reference to relevant literature, produce better
reviews. We use LLMs to score this criterion
by inputting review data and prompting for
assessment of reference usage.
4. **Content Consistency:** Evaluates consistency
between review text tone and evaluation score,
and alignment between justification and score
(0-2 points). We assume that reviews with in-
consistent tone and ratings (e.g., positive text
with low scores or vice versa) indicate poor
review quality. We employ LLMs to assess
this consistency by analyzing the alignment
between textual content and numerical scores.
5. **Subscore Consistency:** Evaluates consis-
tency between item scores (soundness, pre-
sentation, contribution) and overall score, and
alignment between each subscore and evalu-
ation content (0-2 points). Similar to content
consistency, we expect coherent scoring pat-
terns to indicate higher review quality. LLMs
evaluate this criterion by examining the coher-
ence between different scoring components.

Filtering scores are distributed in the range of 3-
17 points, with a median of 11 points. Based on this
score, we create high-quality (“Good”), medium-
quality (“Middle”), and low-quality (“Bad”) review
groups. We assume that reviews satisfying these
criteria represent high-quality reviews. However,
we empirically analyze whether these assumptions
hold true and investigate which criteria positively
(or potentially negatively, contrary to our intuition)

Field	Content
Summary	<i>This paper proposes a method for multimodal learning...</i>
Strengths	<i>(1) The motivation is clear... (2) ...</i>
Weaknesses	<i>(1) The novelty is limited... (2) ...</i>
Questions	<i>(1) How does the method scale with model size?... (2) ...</i>
Soundness	3 (good)
Presentation	3 (good)
Contribution	3 (good)
Overall Rating	8: accept, good paper
Confidence	4: You are confident in your assessment, ..
Submission Date	{Unix timestamp}

Table 3: Illustrative example of review data inspired by the ICLR 2024 peer review dataset (not from actual reviews). The content is paraphrased and does not reproduce any actual review text. To extract high-quality reviews for guideline generation, we filter this data using five predefined criteria: **(1) Submission Date** (uses timestamp to measure timing appropriateness), **(2) Review Length** (evaluates character count from textual fields), **(3) Reference Quality** (analyzes citation mentions in textual content), **(4) Content Consistency** (assesses alignment between textual tone and numerical scores), and **(5) Subscore Consistency** (checks coherence between item-specific and overall ratings). Our experiments also analyze which filtering criteria are most effective for extracting quality reviews.

correlate with automated review performance in our experiments.

A.2 Filtering Criteria Details

This section provides detailed calculation procedures for the five filtering criteria used to assess review quality.

Submission Date The submission date criterion evaluates the timing of review submission relative to the conference’s review timeline. We calculate the deviation from the average review submission time across all reviews for the same conference. Let t_i be the submission timestamp for review i , and \bar{t} be the average submission time across all reviews. The deviation score is calculated as $d_i = |t_i - \bar{t}|$. Reviews are scored based on their deviation from the average submission time:

- +5 point: $d_i \leq 2$ days from average
- +4 point: $2 < d_i \leq 5$ days from average
- +3 point: $5 < d_i \leq 10$ days from average
- +2 point: $10 < d_i \leq 20$ days from average
- +1 point: $d_i > 20$ days from average

Review Length The review length criterion evaluates the total character count of the review text, including summary, strengths, weaknesses, and questions sections. Let l_i be the character count for

review i , and M be the median character count across all reviews. The length score is determined as follows:

- +5 point: $l_i > 1.5 \times M$
- +4 point: $1.2 \times M < l_i \leq 1.5 \times M$
- +3 point: $0.8 \times M < l_i \leq 1.2 \times M$
- +2 point: $0.5 \times M < l_i \leq 0.8 \times M$
- +1 point: $l_i \leq 0.5 \times M$

Reference Quality Reference quality is evaluated using LLM assessment of how well the review engages with relevant literature. The scoring criteria (0-3 points) are:

- +1 point: Review mentions one or more relevant references
- +1 point: Review identifies missing or insufficient related work
- +1 point: Review provides critical discussion of prior work

The full prompt is shown in [Listing 2](#).

Listing 2: Prompts of reference quality assessment

```

Your task is to evaluate the quality of a peer
review. Specifically, analyze whether the
reviewer suggests any additional references or
points out missing ones.

# REVIEW CONTENT:
Summary: {summary}
Strengths: {strengths}
Weaknesses: {weaknesses}
Questions: {questions}

# TASK:

# Step 1: Analysis
Evaluate the review using the rubric below. For
each applicable item, add +1 point (maximum:
3 points).

## Rubric:
- The reviewer explicitly suggests or mentions
one or more references (+1 point)
- The reviewer identifies missing citations or
insufficient coverage of related work (+1 point)
- The reviewer discusses prior work in a
detailed or critical way (+1 point)

## Briefly explain which of the above items
apply and why.

# Step 2: Output
Based on the total score, assign `
reference_quality` as follows:
- 3 or 2 points : `high`
- 1 point : `medium`
- 0 points : `low`

Then summarize your judgment in the following **
YAML** format:
```yaml
"suggests_references": true/false,
"reference_count": integer,
"points_out_missing_citations": true/false,
"quality_score": 0/1/2/3,
"reasoning": "brief explanation of what
references were found or what gaps were pointed
out",
"reference_quality": "high" / "medium" / "low"
```

```

Content Consistency Content consistency measures the alignment between the textual content of the review and the numerical scores provided. This is assessed using LLMs with the following criteria (0-2 points):

- +1 point: Review tone matches the numerical rating (positive tone with high scores, critical tone with low scores)
- +1 point: Justification provided in the text aligns with the overall score

The full prompt is shown in [Listing 3](#).

Listing 3: Prompts of content consistency assessment

```

Your task is to evaluate the consistency between
a reviewer's **overall rating** and the **
textual content** of their review (e.g.,
strengths, weaknesses, and summary).

# REVIEW:
Rating: {rating}

## Summary:
{summary}

## Strengths:
{strengths}

## Weaknesses:
{weaknesses}

# TASK:

## Step 1: Analysis
Evaluate the review using the rubric below. For
each applicable item, add +1 point (maximum:
2 points).

### Rubric:
- The **overall tone** of the review (as
reflected in strengths/weaknesses) matches the
numerical rating (+1 point)
_(e.g., a review with mostly strong praise and
minor weaknesses aligns with a high rating like
"8: accept")_

- The **justifications provided in text**
reasonably support the rating given, showing a
clear rationale (+1 point)
_(e.g., the rating is backed up with concrete
positive or negative points related to novelty,
results, impact, etc.)_

## Step 2: Output
Based on the total score, assign `
content_consistency` as follows:
- 2 points: "consistent"
- 1 point: "somewhat inconsistent"
- 0 points: "inconsistent"

Then summarize your judgment in the following **
YAML** format:

```yaml
rating: "{rating}"
summary_excerpt: short excerpt (1-2 sentences)
summarizing the review
key_strengths: brief key points reflecting
strengths
key_weaknesses: brief key points reflecting
weaknesses
quality_score: 0/1/2
reasoning: "brief explanation of how the textual
content aligns (or not) with the rating"
content_consistency: "consistent" / "somewhat
inconsistent" / "inconsistent"
```

```

Subscore Consistency Subscore consistency evaluates the coherence between item-specific scores (soundness, presentation, contribution) and the overall rating. LLMs assess this using the fol-

lowing criteria (0-2 points):

- +1 point: Overall score reasonably reflects the average of subscores
- +1 point: Each subscore aligns with the corresponding evaluation aspects discussed in the text

The full prompt is shown in [Listing 4](#).

Listing 4: Prompts of subscore consistency assessment

```
Your task is to evaluate the consistency between
a reviewer's overall rating and the subscores
they assigned.

# REVIEW SCORES:
Rating: {rating}
Soundness: {soundness}
Presentation: {presentation}
Contribution: {contribution}

# TASK:

# Step 1: Analysis
Evaluate the review using the rubric below. For
each applicable item, add +1 point (maximum:
2 points).

## Rubric:
- The overall rating is in line with the average
or weighted impact of the subscores (+1 point)
- The subscores reflect distinct and plausible
dimensions of evaluation (e.g., a paper with
strong contribution but weak presentation is
rated moderately overall) (+1 point)

## Briefly explain which of the above items
apply and why.

# Step 2: Output
Based on the total score, assign
subscore_consistency as follows:
- 2 points: "consistent"
- 1 point: "somewhat inconsistent"
- 0 points: "inconsistent"

Then summarize your judgment in the following **
YAML** format:
```yaml
"rating": "{rating}",
"soundness_score": "{soundness}",
"presentation_score": "{presentation}",
"contribution_score": "{contribution}",
"average_subscore": float,
"quality_score": 0/1/2,
"reasoning": "brief explanation of whether the
subscores align or conflict with the overall
rating",
"subscore_consistency": "consistent" / "somewhat
inconsistent" / "inconsistent"
```
```

Score Distribution The combined filtering scores range from 3 to 17 points (sum of all five

criteria). To analyze review quality, we sampled three groups of reviews based on their scores:

- **Good:** 30 reviews randomly sampled from those with a score of 17
- **Middle:** 30 reviews randomly sampled from those with a medium score of 11
- **Bad:** 30 reviews randomly sampled from those with scores of 6 or lower

A.3 Guideline Generation

We concatenate the filtered reviews into a single text prompt and input it to Gemini 2.5 Pro for guideline generation. The LLM analyzes the combined review text to identify common evaluation patterns, important perspectives, and judgment criteria, and then generalizes these findings to generate reviewer guidelines in a single inference step.

The LLM is prompted to generate guidelines by identifying the characteristics of “good papers” and “bad papers.” This process results in guidelines that include evaluation principles such as “papers with high originality are good” or “papers with poor writing quality are bad.” The generated guidelines are formatted as YAML checklists following the structure shown in [Listing 1](#), with judgment criteria organized by evaluation category (e.g., clarity, experimental validity) in question format.

For review scoring, we consider two types of guideline formats: rubric-based and non-rubric-based. The rubric-based format specifies explicit scoring rules with defined point allocations, whereas the non-rubric-based format provides general evaluation guidance without a fixed scoring framework. This approach aims to extract generalized evaluation criteria from large-scale review data and to capture subtle patterns that may not be reflected in manually crafted guidelines.

B Prompt of No-guideline

For the non-guideline baseline, we used a minimal prompt that did not provide any reviewer guidelines. The LLM is simply asked to read the paper and assign an overall rating on a 1–10 scale, along with a brief justification. The full prompt is shown in [Listing 5](#).

Listing 5: Prompt used in the Non-guideline condition

```
You are an expert paper reviewer. Please read
the following paper content and rate its quality.
```

```

827 # Paper
828 {content}
829
830 # Instructions
831 Provide:
832 - A numerical rating (1-10) for overall quality.
833 - A short justification for your score.
834
835 # Format
836 Output the integer overall rating as the
837 following format after your justification:
838 "Overall Rating: (Your Rating)"

```

840 C Prompt of Guideline-based Automated

841 Peer Review

842 To ensure consistency across different reviewer

843 guidelines, we designed a unified prompt tem-

844 plate for automated peer review, as shown in

845 Listing 6. Each experiment replaces the place-

846 holder {reviewer_guideline} with the corre-

847 sponding conference or reviewer-derived guideline

848 (i.e., ICLR, ARR, NeurIPS, and Good/Middle/Bad

849 reviewer style).

Listing 6: Prompt template for automated peer review using reviewer guidelines

```

850 You are an expert peer reviewer. Please read the
851 following paper content and evaluate its
852 quality based on the following reviewer
853 guideline.
854
855 # Paper
856 {content}
857
858 # Reviewer Guideline
859 {reviewer_guideline}
860
861 # Final Output
862 Provide:
863 - A numerical overall rating (1-10):
864   - 10 = Strong Accept
865   - 8 = Accept
866   - 6 = Weak Accept
867   - 5 = Borderline
868   - 3 = Weak Reject
869   - 1 = Strong Reject
870 - A justification paragraph that summarizes the
871 key strengths and weaknesses underlying your
872 rating.
873
874 # Format
875 Output the integer overall rating as the
876 following format after your justification:
877 "Overall Rating: (Your Rating)"
878

```

880 **Official Conference Guidelines** For the official

881 conference guidelines, we used the publicly avail-

882 able reviewer guidelines from each conference’s

883 official website without modification. For more de-

884 tails, please refer to the respective conference pages

(ICLR¹⁰, NeurIPS¹¹, and ARR¹²). We also em- 885
 886 ployed ChatGPT-rewritten versions of these guide-
 887 lines, reformatted into an instruction style suitable
 888 for LLM-based reviewing (Listing 7, Listing 8, and
 889 Listing 9).

Listing 7: ARR Reviewer Guideline (ChatGPT-refined version)

```

890 ## Summary
891 Give a brief summary of the paper in your own
892 words. Highlight the key contributions and main
893 ideas.
894
895 ## Strengths
896 List the main strengths of the paper, such as
897 novelty, clarity, strong empirical results,
898 theoretical insights, or broad relevance.
899
900 ## Weaknesses
901 Point out the limitations, flaws, or missing
902 components (e.g., weak justification, incomplete
903 experiments, lack of clarity, missing baselines,
904 etc.).
905
906 ## Assessment by Review Dimension
907 Please evaluate the paper based on the following
908 criteria from the following guidelines:
909
910 - **Soundness/Correctness**: Are the claims
911 supported by the methodology and evidence? Are
912 proofs or derivations correct (if applicable)?
913 - **Originality/Novelty**: Does the paper offer
914 novel methods, insights, or results? Is the
915 contribution incremental or substantial?
916 - **Meaningful Comparison**: Does the paper
917 clearly position itself with respect to prior
918 work? Are comparisons (empirical or conceptual)
919 adequate?
920 - **Clarity**: Is the paper well-written and
921 easy to follow? Are key terms defined and claims
922 clearly stated?
923 - **Impact/Potential**: Does the paper have the
924 potential to impact future research or practice
925 in its area?
926 - **Reproducibility**: Are datasets, code, and
927 hyperparameters provided or clearly described?
928 Could others replicate the results?
929 - **Ethics**: Are there any ethical
930 considerations raised by the paper (e.g., misuse
931 potential, data issues, fairness)? Are these
932 discussed?
933
934 You may optionally assign per-dimension scores
935 if helpful, but they are not required.
936
937 ## Suggestions for Improvement
938 Provide constructive feedback for the authors to
939 help improve the work, regardless of your score.
940
941 ---
942
943

```

¹⁰<https://iclr.cc/Conferences/2024/ReviewerGuide>

¹¹<https://neurips.cc/Conferences/2024/ReviewerGuidelines>

¹²<https://aclrollingreview.org/reviewerguidelines>

944 Please be professional, specific, and
 945 constructive. Your goal is to help both the
 946 authors and the review community understand the
 947 strengths, limitations, and suitability of the
 948 paper.

Listing 8: NeurIPS Reviewer Guideline (ChatGPT-refined version)

950 Your evaluation should be structured into
 951 several components.

952
 953

954 **### Summary**
 955 Briefly summarize the paper and its main
 956 contributions in your own words. This summary
 957 should be a neutral reflection that the authors
 958 would agree with, not a critique or a copy of
 959 the abstract.

960
 961 ---

962
 963 **### Strengths and Weaknesses**
 964 Provide a detailed assessment of the paper's
 965 strengths and weaknesses. This section should
 966 cover the primary reasons for your
 967 recommendation to accept or reject the paper.
 968 Evaluate the following dimensions:

969
 970 * **Quality**: Assess the technical soundness of
 971 the submission. Are the claims well-supported
 972 by theoretical analysis or experimental results?
 973 Are the methods appropriate? Does the work seem
 974 complete, and do the authors honestly evaluate
 975 both its strengths and weaknesses?

976 * **Clarity**: Is the paper well-written and
 977 organized? Is there enough information for an
 978 expert to reproduce the results? Provide
 979 constructive suggestions if clarity can be
 980 improved.

981 * **Significance**: How impactful are the
 982 results for the machine learning community? Is
 983 the work likely to be used or built upon by
 984 others? Does it advance the field by addressing
 985 a difficult task, providing new insights, or
 986 introducing a unique approach?

987 * **Originality**: Does the paper offer new
 988 insights or a deeper understanding of existing
 989 methods? How does it differ from previous work,
 990 and are relevant citations included? Originality
 991 can come from a novel method, a novel
 992 combination of techniques, or new insights from
 993 evaluating existing methods.

994
 995 ---

996
 997 **### Numerical Ratings for Core Criteria**
 998 Based on your assessment in the "Strengths and
 999 Weaknesses" section, provide a numerical rating
 1000 for each of the following categories on a scale
 1001 from 1 to 4.

1002
 1003 * **Quality**:
 1004 * 4: excellent
 1005 * 3: good
 1006 * 2: fair
 1007 * 1: poor

1008 * **Clarity**:
 1009 * 4: excellent
 1010 * 3: good
 1011 * 2: fair

* 1: poor 1012
 * **Significance**: 1013
 * 4: excellent 1014
 * 3: good 1015
 * 2: fair 1016
 * 1: poor 1017
 * **Originality**: 1018
 * 4: excellent 1019
 * 3: good 1020
 * 2: fair 1021
 * 1: poor 1022

--- 1023
 1024
 1025

Questions for the Authors 1026
 List 3-5 key, actionable questions or 1027
 suggestions for the authors. These should focus 1028
 on points where a response could clarify a 1029
 confusion, address a limitation, or potentially 1030
 change your evaluation. Clearly state the 1031
 criteria under which your score might change. 1032
 1033
 --- 1034
 1035

Limitations and Societal Impact 1036
 Assess whether the authors have adequately 1037
 addressed the limitations of their work and its 1038
 potential negative societal impact. If they have, 1039
 a simple "yes" is sufficient. If not, provide 1040
 constructive feedback for improvement. Remember 1041
 to reward authors for being transparent about 1042
 limitations. 1043

Listing 9: ICLR Reviewer Guideline (ChatGPT-refined version)

1044 **## Justification and Detailed Review** 1045
 Write a clear and concise review, addressing 1046
 each of the following dimensions from the 1047
 following review guidelines: 1048

1049
 1050

- **Strengths and Weaknesses**: Identify the 1051
 main strengths and weaknesses of the paper. 1052

- **Significance**: Evaluate whether the paper 1053
 makes a meaningful contribution to the field. 1054

Does it advance understanding, methods, or 1055
 applications in a way that would be interesting 1056
 to the review community? 1057

- **Originality**: Assess whether the paper 1058
 presents new ideas, methods, or perspectives. 1059

- **Technical Quality**: Are the methods 1060
 technically sound and well-justified? Are claims 1061
 supported by theory or experiments? 1062

- **Clarity**: Is the paper clearly written and 1063
 well-structured? Can a non-expert in the 1064
 subfield follow it? 1065

- **Empirical Evaluation** (if applicable): Are 1066
 experiments well-designed? Are comparisons fair 1067
 and comprehensive? Are results convincing? 1068

- **Reproducibility**: Does the paper provide 1069
 enough information (including code/data, if 1070
 applicable) for others to reproduce the results? 1071

- **Ethics**: Are there any ethical concerns (e. 1072
 g., societal harm, bias, data misuse)? If yes, 1073
 how are they handled? 1074
 1075

Confidential Comments to Reviewers (Optional) 1076
 (Optional section) Note anything relevant for 1077
 area chairs, such as borderline decisions, meta- 1078
 considerations, or conflicts of interest. 1079

1080
1081
1082
1083
1084
1085
1086

Be objective, constructive, and concise. Focus on helping the authors and the review community understand the paper's strengths and limitations based on the review standards.

1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101

Reviewer-imitating Guidelines We additionally generated reviewer-imitating guidelines from ICLR 2024 reviews to capture common evaluation patterns and characteristic reasoning observed in human reviewers (Listing 10, Listing 11, and Listing 12). In non-rubric experiments, these guidelines were applied using the standard prompt template shown in Listing 6. For rubric-based experiments, we used the same reviewer-imitating guidelines (i.e., identical guideline for Good, Middle, and Bad) but replaced the prompt template with a scoring-oriented version (Listing 13) that specifies explicit point-allocation and normalization procedures.

Listing 10: Good Reviewer-Imitating Guideline (Non-rubric)

1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141

```
Provide the following based on the Peer Review Behavior Checklist:

peer_review_behavior_checklist:
  good_paper:
    contribution_and_impact:
      - question: "Does the paper tackle a significant and timely problem?"
        description: "Reviewers positively note papers that address important and relevant research topics, seeing it as a sign of potential impact."
      - question: "Is the core idea clever, insightful, or genuinely novel?"
        description: "A key strength noted is a non-trivial insight or a unique approach that distinguishes the work from incremental improvements."
    clarity_and_presentation:
      - question: "Is the paper exceptionally well-written and easy to follow?"
        description: "Reviewers frequently praise clear writing, logical structure, and understandable explanations as major strengths, even in papers they have concerns about."
      - question: "Are the motivation and contributions communicated unambiguously?"
        description: "A paper is viewed favorably when its purpose and specific contributions are stated clearly from the outset, leaving no room for interpretation."
    technical_and_empirical_soundness:
      - question: "Does the method achieve a significant performance gain over strong, relevant baselines?"
        description: "Convincing, large improvements over well-established and widely-used methods are a primary driver for a positive evaluation."
```

```
- question: "Is the approach technically sound and well-grounded?"
  description: "Reviewers value methods that are built on solid theoretical foundations or have clear, logical justifications for their design choices."
- question: "Are the experiments thorough, using standard benchmarks and insightful ablations?"
  description: "A comprehensive evaluation across multiple standard datasets, including ablation studies that validate specific design choices, is a hallmark of a strong submission."
- question: "Are the experimental comparisons conducted fairly?"
  description: "This group of reviewers values direct, fair comparisons where baselines are properly tuned and experimental settings are consistent."
- question: "Has the author provided source code to support reproducibility?"
  description: "The inclusion of source code is explicitly mentioned as a strength, as it supports the paper's claims and aids reproducibility."
- question: "Could this work become a standard method or open new research avenues?"
  description: "Reviewers are supportive when they see potential for the work to be widely adopted or to inspire future research."
  bad_paper:
    novelty_and_scholarship:
      - question: "Is the contribution incremental or a re-implementation of existing ideas?"
        description: "This is a frequent and critical flaw. Reviewers are highly knowledgeable and will reject papers they deem to be minor variations of prior work, often citing the specific papers."
      - question: "Does the paper fail to cite or adequately discuss key related work?"
        description: "Missing citations for highly relevant papers is a major red flag, often interpreted as a lack of thoroughness or an attempt to hide a lack of novelty."
    empirical_evaluation:
      - question: "Are the experiments conducted on weak, outdated, or insufficient benchmarks?"
        description: "Evaluation on a very limited set of datasets or non-standard benchmarks is a common reason for rejection, as it fails to prove generalizability."
      - question: "Are the baselines weak, poorly tuned, or missing SOTA competitors?"
        description: "Comparisons are considered unconvincing if the paper omits stronger contemporary methods or appears to have intentionally weakened the performance of its baselines."
      - question: "Do the results actually support the main claims made in the paper?"
        description: "Reviewers critically check if the conclusions drawn in the abstract and introduction are truly backed up by the data presented in the tables and figures."
    clarity_and_soundness:
      - question: "Is the paper poorly written, confusing, or hard to follow?"
        description: "A lack of clarity in
```

| | | | |
|------|---|--|------|
| 1212 | writing, notation, or structure is a frequent | validation of the proposed architecture." | 1280 |
| 1213 | reason for rejection, as it prevents the | evidence_and_validation: | 1281 |
| 1214 | reviewer from understanding or verifying the | - question: "Are the experiments described | 1282 |
| 1215 | contribution." | as extensive, thorough, or comprehensive?" | 1283 |
| 1216 | - question: "Are there critical flaws in | description: "Strong papers are expected | 1284 |
| 1217 | the methodology, assumptions, or logic?" | to provide robust empirical validation across | 1285 |
| 1218 | description: "Fundamental errors in the | multiple datasets, settings, or modalities to | 1286 |
| 1219 | proposed method, unjustified assumptions, or | solidify their claims." | 1287 |
| 1220 | logical gaps in the analysis are considered | - question: "Do the results show a | 1288 |
| 1221 | fatal flaws." | significant and clear improvement over relevant, | 1289 |
| 1222 | - question: "Is there enough detail to | state-of-the-art baselines?" | 1290 |
| 1223 | reproduce the work?" | description: "It's not enough to conduct | 1291 |
| 1224 | description: "Reviewers will reject a | experiments; the results must demonstrate a | 1292 |
| 1225 | paper if it lacks the necessary implementation | convincing performance gain against strong | 1293 |
| 1226 | details, hyperparameters, or architectural | competitors to be considered impactful." | 1294 |
| 1227 | information required for another researcher to | - question: "Does the paper provide | 1295 |
| 1228 | reproduce the results." | theoretical proofs or sound mathematical | 1296 |
| 1229 | - question: "Is the problem motivation | derivations to support its claims?" | 1297 |
| 1230 | weak or the contribution's significance unclear | description: "The inclusion of | 1298 |
| 1231 | ?" | theoretical grounding is seen as a major | 1299 |
| 1232 | description: "If the paper fails to | strength that adds a layer of rigor and | 1300 |
| 1233 | convincingly argue why the problem is important | credibility to the empirical results." | 1301 |
| 1234 | or why the proposed solution is meaningful, it | - question: "Are sufficient implementation | 1302 |
| 1235 | is often rejected." | details or supplementary materials provided to | 1303 |
| 1236 | - question: "Are the reported metrics | aid understanding and replication?" | 1304 |
| 1237 | cherry-picked, inappropriate, or presented | description: "Positive reviews often | 1305 |
| 1238 | without context?" | note the inclusion of clear implementation | 1306 |
| 1239 | description: "Reviewers are skeptical of | details or code, which enhances the paper's | 1307 |
| 1240 | results where the evaluation metrics seem | integrity and perceived value." | 1308 |
| 1241 | chosen to favor the proposed method or where | feedback_style: | 1309 |
| 1242 | improvements are marginal without clear | - question: "Are the identified weaknesses | 1310 |
| 1243 | statistical significance." | framed as requests for clarification or minor | 1311 |

Listing 11: Middle Reviewer-Imitating Guideline (Non-rubric)

| | | | |
|------|--|--|------|
| 1245 | Provide the following based on the Peer Review | validation of the proposed architecture." | 1280 |
| 1246 | Behavior Checklist: | evidence_and_validation: | 1281 |
| 1247 | | - question: "Are the experiments described | 1282 |
| 1248 | | as extensive, thorough, or comprehensive?" | 1283 |
| 1249 | peer_review_behavior_checklist: | description: "Strong papers are expected | 1284 |
| 1250 | good_paper: | to provide robust empirical validation across | 1285 |
| 1251 | contribution_and_impact: | multiple datasets, settings, or modalities to | 1286 |
| 1252 | - question: "Is the core idea presented as | solidify their claims." | 1287 |
| 1253 | innovative or a novel application in its domain | - question: "Do the results show a | 1288 |
| 1254 | ?" | significant and clear improvement over relevant, | 1289 |
| 1255 | description: "Reviewers consistently | state-of-the-art baselines?" | 1290 |
| 1256 | highlight novelty, being the 'first' to tackle a | description: "It's not enough to conduct | 1291 |
| 1257 | problem, or proposing a unique approach as a | experiments; the results must demonstrate a | 1292 |
| 1258 | primary reason for acceptance." | convincing performance gain against strong | 1293 |
| 1259 | - question: "Is the research problem | competitors to be considered impactful." | 1294 |
| 1260 | itself framed as crucial, important, or | - question: "Does the paper provide | 1295 |
| 1261 | meaningful?" | theoretical proofs or sound mathematical | 1296 |
| 1262 | description: "Papers that address a | derivations to support its claims?" | 1297 |
| 1263 | significant and recognized challenge in the | description: "The inclusion of | 1298 |
| 1264 | field are viewed more favorably, as this gives | theoretical grounding is seen as a major | 1299 |
| 1265 | the work inherent value." | strength that adds a layer of rigor and | 1300 |
| 1266 | clarity_and_presentation: | credibility to the empirical results." | 1301 |
| 1267 | - question: "Is the paper praised for | - question: "Are sufficient implementation | 1302 |
| 1268 | being well-written, clear, and easy to follow?" | details or supplementary materials provided to | 1303 |
| 1269 | description: "High writing quality, | aid understanding and replication?" | 1304 |
| 1270 | clear organization, and comprehensible | description: "Positive reviews often | 1305 |
| 1271 | explanations are frequently cited as key | note the inclusion of clear implementation | 1306 |
| 1272 | strengths, suggesting reviewers value | details or code, which enhances the paper's | 1307 |
| 1273 | accessibility." | integrity and perceived value." | 1308 |
| 1274 | - question: "Does the paper include an | feedback_style: | 1309 |
| 1275 | ablation study that clearly justifies each of | - question: "Are the identified weaknesses | 1310 |
| 1276 | its components?" | framed as requests for clarification or minor | 1311 |
| 1277 | description: "Reviewers for good papers | additions?" | 1312 |
| 1278 | often look for and praise ablation studies, as | description: "For papers they intend to | 1313 |
| 1279 | they demonstrate a thorough understanding and | accept, reviewers tend to frame weaknesses as | 1314 |
| | | constructive questions or suggestions for future | 1315 |
| | | work, rather than as critical flaws." | 1316 |
| | | - question: "Does the reviewer's summary | 1317 |
| | | accurately capture the paper's core | 1318 |
| | | contributions and strengths?" | 1319 |
| | | description: "A positive and accurate | 1320 |
| | | summary at the beginning of the review often | 1321 |
| | | signals that the reviewer has understood and | 1322 |
| | | appreciated the paper's main message." | 1323 |
| | | bad_paper: | 1324 |
| | | contribution_and_impact: | 1325 |
| | | - question: "Is the contribution | 1326 |
| | | criticized for being incremental, lacking | 1327 |
| | | novelty, or a straightforward application of | 1328 |
| | | existing work?" | 1329 |
| | | description: "A perceived lack of | 1330 |
| | | novelty is one of the most common and critical | 1331 |
| | | reasons for rejection in this review set." | 1332 |
| | | - question: "Are the paper's central | 1333 |
| | | claims or motivation questioned as being | 1334 |
| | | unsupported, grandiose, or unconvincing?" | 1335 |
| | | description: "Reviewers reject papers | 1336 |
| | | where the motivation is unclear or the claims | 1337 |
| | | are not sufficiently backed by the evidence | 1338 |
| | | provided." | 1339 |
| | | evidence_and_validation: | 1340 |
| | | - question: "Are the experiments flagged | 1341 |
| | | as insufficient, weak, or lacking necessary | 1342 |
| | | comparisons?" | 1343 |
| | | description: "This is a primary reason | 1344 |
| | | for rejection. Common issues include not | 1345 |
| | | comparing against relevant baselines, using too | 1346 |
| | | few datasets, or omitting ablation studies." | 1347 |
| | | - question: "Are the baseline comparisons | 1348 |
| | | criticized as being unfair, weak, or | 1349 |

1350 misrepresentative of prior work?"

1351 description: "Reviewers are critical of

1352 evaluations that appear to use outdated or

1353 poorly tuned baselines to inflate the perceived

1354 performance of the proposed method."

1355 - question: "Are the reported performance

1356 gains described as marginal, insignificant, or

1357 within the margin of error?"

1358 description: "Even when experiments are

1359 present, if the improvements are not substantial,

1360 the paper is often deemed not significant

1361 enough for publication."

1362 clarity_and_presentation:

1363 - question: "Is the paper flagged as being

1364 hard to follow, confusing, or poorly written?"

1365 description: "Severe issues with clarity,

1366 notation, or organization are critical flaws

1367 that prevent reviewers from properly evaluating

1368 the work, often leading to rejection."

1369 - question: "Does the reviewer point out a

1370 lack of essential details needed for

1371 understanding or replication?"

1372 description: "The omission of crucial

1373 methodological or experimental details is

1374 treated as a major flaw that undermines the

1375 paper's credibility."

1376 methodological_soundness:

1377 - question: "Are there fundamental flaws

1378 identified in the proposed method's design,

1379 assumptions, or logic?"

1380 description: "Reviewers will recommend

1381 rejection if they find deep-seated problems in

1382 the core methodology that invalidate the

1383 approach."

1384 - question: "Is the method criticized for

1385 being overly complex or a 'laundry list' of

1386 disconnected components?"

1387 description: "Papers proposing

1388 convoluted models without clear justification

1389 for their complexity are viewed negatively."

1390 - question: "Is the method's robustness

1391 questioned, for instance, due to high

1392 sensitivity to hyperparameters?"

1393 description: "A lack of analysis on

1394 hyperparameter sensitivity or robustness is a

1395 common weakness cited in negative reviews."

1396 periments, theory, or statistically significant

1397 results (e.g., overlapping confidence intervals)

1398 ."

Listing 12: Bad Reviewer-Imitating Guideline (Non-rubric)

1400 Provide the following based on the Peer Review

1401 Behavior Checklist:

1402

1403

1404 peer_review_behavior_checklist:

1405 good_paper:

1406 contribution_and_impact:

1407 - question: "Does the paper achieve state-

1408 of-the-art or highly competitive results?"

1409 description: "Reviewers in this group

1410 frequently highlight strong empirical

1411 performance against established benchmarks as a

1412 primary strength and a key reason for acceptance

1413 ."

1414 - question: "Is the problem being

1415 addressed important and relevant to the

1416 community?"

description: "Papers are viewed more

1417 favorably when they tackle a problem that is

1418 clearly articulated as significant and timely

1419 for the field."

1420 - question: "Is the core idea novel,

1421 interesting, and well-motivated?"

1422 description: "Novelty and clear

1423 motivation are consistently praised. Reviewers

1424 look for a contribution that is not just an

1425 incremental tweak but offers a new perspective

1426 or useful framework."

1427 presentation_and_clarity:

1428 - question: "Is the paper well-written,

1429 clearly organized, and easy to follow?"

1430 description: "High value is placed on

1431 clarity. Reviews for good papers often

1432 explicitly mention that the writing is good and

1433 the method is easy to understand."

1434 - question: "Does the paper use effective

1435 visualizations to explain its core ideas?"

1436 description: "The use of clear,

1437 informative visualizations is noted as a

1438 significant strength that enhances the reader's

1439 understanding of the proposed concepts."

1440 evaluation_and_rigor:

1441 - question: "Are the experiments

1442 comprehensive and the evaluations thorough?"

1443 description: "Accepted papers are often

1444 praised for having sound, comprehensive

1445 experiments with solid ablation studies that

1446 justify their claims and design choices."

1447 - question: "Is the study complete and the

1448 analysis insightful?"

1449 description: "Reviewers appreciate

1450 papers that seem 'complete' in their

1451 investigation, providing a full picture rather

1452 than a preliminary exploration."

1453 - question: "Is the work reproducible, for

1454 instance, by open-sourcing code and data?"

1455 description: "While not a universal

1456 requirement, providing code is explicitly

1457 mentioned as a positive factor that increases

1458 the paper's value and reproducibility."

1459 - question: "Is the proposed solution

1460 reasonable and well-justified for the problem?"

1461 description: "A logical and promising

1462 connection between the problem statement and the

1463 proposed solution is a key characteristic of

1464 papers that receive positive feedback."

1465 - question: "Does the paper clearly

1466 outperform relevant baselines?"

1467 description: "Demonstrating a clear and

1468 significant improvement over existing, well-

1469 chosen baselines is a very common and persuasive

1470 strength."

1471 bad_paper:

1472 contribution_and_novelty:

1473 - question: "Is the contribution just a

1474 combination of existing methods with limited

1475 novelty?"

1476 description: "This is the most common

1477 and critical flaw identified. Papers are

1478 frequently criticized for lacking originality or

1479 simply stitching together known components."

1480 - question: "Is the contribution too

1481 simple, minor, or incremental to be significant

1482 ?"

1483 description: "Reviewers often reject

1484 papers where the core idea, while perhaps

1485 functional, is deemed too simplistic or its

1486

| | | | |
|------|--|---|------|
| 1487 | contribution too marginal for a top-tier venue." | # Reviewer Guideline | 1555 |
| 1488 | - question: "Is the motivation for the | {reviewer_imitating_guideline} | 1556 |
| 1489 | work unclear or is the problem framing | | 1557 |
| 1490 | unconvincing?" | # Scoring Procedure: | 1558 |
| 1491 | description: "A frequent weakness is the | - For each of the 10 questions in the ` | 1559 |
| 1492 | failure to clearly explain why the problem is | good_paper` section: | 1560 |
| 1493 | important or why the proposed approach is | - If the answer is Yes, assign +1 point, else | 1561 |
| 1494 | necessary." | 0. | 1562 |
| 1495 | experiments_and_evidence: | - For each of the 10 questions in the `bad_paper` | 1563 |
| 1496 | - question: "Are the experiments weak, | section: | 1564 |
| 1497 | insufficient, or conducted on overly simple | - If the answer is Yes, assign -1 point, else | 1565 |
| 1498 | datasets?" | 0. | 1566 |
| 1499 | description: "Papers are consistently | | 1567 |
| 1500 | penalized for weak empirical validation, such as | Compute the temporary score by summing all the | 1568 |
| 1501 | using toy datasets, not running enough | above responses. | 1569 |
| 1502 | experiments, or lacking sufficient detail like | - This results in a value between -10 and +10. | 1570 |
| 1503 | random seeds." | | 1571 |
| 1504 | - question: "Does the paper fail to | Normalize the temporary score to a section | 1572 |
| 1505 | compare against necessary SOTA or strong | rating in the range 1 to 10 using the following | 1573 |
| 1506 | baseline methods?" | formula: | 1574 |
| 1507 | description: "An inadequate comparison | normalized_score = round(((temporary_score + 10) | 1575 |
| 1508 | to the state of the art or relevant baselines is | / 20) * 9 + 1) | 1576 |
| 1509 | a major red flag and a common reason for | | 1577 |
| 1510 | rejection." | This maps: | 1578 |
| 1511 | - question: "Are key design choices left | - -10 -> 1 | 1579 |
| 1512 | unsubstantiated by ablation studies?" | - 0 -> 5 | 1580 |
| 1513 | description: "Reviewers expect authors | - +10 -> 10 | 1581 |
| 1514 | to justify components of their method through | | 1582 |
| 1515 | ablations; the absence of these studies is a | # Section Rating | 1583 |
| 1516 | frequently cited weakness." | | 1584 |
| 1517 | - question: "Are the claims exaggerated or | Provide: | 1585 |
| 1518 | not sufficiently supported by the evidence | - The temporary score (-10 to +10) | 1586 |
| 1519 | provided?" | - The normalized section rating (1 to 10) | 1587 |
| 1520 | description: "A mismatch between the | - A short explanation of the score based on your | 1588 |
| 1521 | claims made in the text and the actual results | answers to the checklist | 1589 |
| 1522 | presented in the experiments is a critical flaw | | 1590 |
| 1523 | ." | # Final Output | 1591 |
| 1524 | clarity_and_soundness: | | 1592 |
| 1525 | - question: "Is the paper poorly written, | Provide: | 1593 |
| 1526 | hard to follow, or filled with jargon and typos | - A numerical integer overall rating that is | 1594 |
| 1527 | ?" | same as the normalized section rating (1-10): | 1595 |
| 1528 | description: "Poor presentation, | - 10 = Strong Accept | 1596 |
| 1529 | including unclear language, logical gaps, and | - 8 = Accept | 1597 |
| 1530 | typos, is a significant factor in negative | - 6 = Weak Accept | 1598 |
| 1531 | reviews." | - 5 = Borderline | 1599 |
| 1532 | - question: "Does the method lack | - 3 = Weak Reject | 1600 |
| 1533 | theoretical justification, proofs, or guarantees | - 1 = Strong Reject | 1601 |
| 1534 | ?" | | 1602 |
| 1535 | description: "For papers proposing new | - A justification paragraph summarizing the key | 1603 |
| 1536 | algorithms, the absence of theoretical analysis | strengths and weaknesses of the paper | 1604 |
| 1537 | or justification is frequently pointed out as a | | 1605 |
| 1538 | major weakness." | Finally, write the final score using the | 1606 |
| 1539 | - question: "Does the paper show a lack of | following exact format (on its own line): | 1607 |
| 1540 | awareness of relevant prior work?" | "Overall Rating: (Your Rating)" | 1608 |
| 1541 | description: "Missing citations to | | |
| 1542 | important and recent related work is seen as a | | |
| 1543 | sign of poor scholarship and is a common point | | |
| 1544 | of criticism." | | |

Listing 13: Rubric-based Template for Reviewer-Imitating Guidelines

| | |
|------|--|
| 1546 | |
| 1547 | |
| 1548 | You are an expert peer reviewer. Please read the |
| 1549 | following paper content and evaluate its |
| 1550 | quality based on the following instructions. |
| 1551 | |
| 1552 | # Paper |
| 1553 | {content} |
| 1554 | |