

ARR-Dv2: New Data from ACL Rolling Review and its Potential to Improve Reviewing in the ACL Community

Anonymous ACL submission

Abstract

Peer review is a central part of the scientific publication cycle, and the peer review service has long been an honor and obligation that researchers gladly fulfill. With the considerable increase in submissions, this aspect of voluntary contribution to the community can become a burden due to the increased workload. We present a new dataset to *support* AI-assisted and human studies of the peer reviewing process, based on the ARR Open Review Platform. Unlike previous similar data collections, ours (ARR-Dv2) not only contains the **reviews** but also author-reviewer **discussions** (aka rebuttal phase) and **meta-reviews**. We present statistical analyses of the data with respect to the overall score and the correlations among the overall score, confidence, and soundness scores. Additionally, we extend the task of *Review Score Prediction* also to include rebuttal data and analyze its effect on score prediction. Finally, we introduce a new task: *Meta Review Score Prediction*, which is based on a set of up to three reviews rather than a single review. Our initial results in a zero-shot setting indicate that the rebuttal data adds valuable information to the score prediction and enables reliable predictions.

1 Introduction

Peer review is the key mechanism for quality control of scientific manuscripts across academia (Birukou et al., 2011), playing a vital role in ensuring research integrity and reliable progress (Bornmann, 2011). Commonly, a group of researchers (the *reviewers*) scrutinizes the manuscript and each composes a *review report* that serves as the basis for a *meta-reviewer* to create a condensed overview of the individual reviews, known as *meta-review report*. Despite its importance to the scientific community with the growing number of publications in science (Landhuis, 2016) and Artificial Intelligence (AI) research in particular (Sculley et al.,

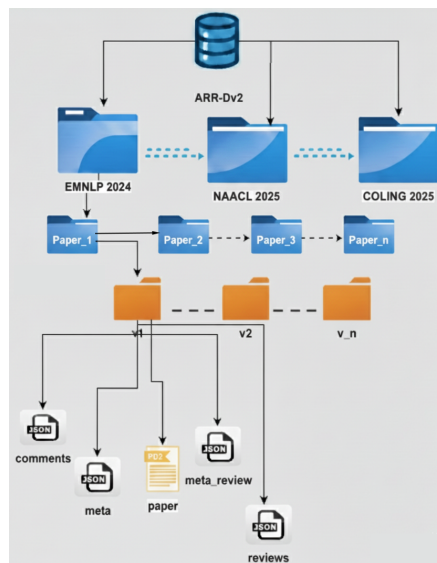


Figure 1: The structure of ARR-Dv2 dataset

2018), the reviewing load places significant burden on researchers. This amplifies existing reviewer biases (Lee et al., 2013), encourages outsourcing to junior researchers (Kuznetsov et al., 2024b), or – with the rise of generative AI (GenAI) tools – to Large Language Models (LLMs) such as ChatGPT (Achiam et al., 2023) and the like (Liang et al., 2024). This, in turn, introduces new biases and quality issues (Du et al., 2024). Putting the reliability of peer review at risk conceivably has severe consequences for the scientific discovery process as a whole (Ebadi et al., 2025; Hanafi et al., 2025). Human-AI collaboration and assistance for peer reviewers is a promising tool to mitigate these issues by speeding up the review process without sacrificing quality, and while maintaining full control for human reviewers (Kuznetsov et al., 2024a). Training and evaluation data are essential for the development of such systems. While there is a plethora of peer reviewing corpora (Yuan et al., 2022; Zhang et al., 2025), most of them come from AI research and open reviewing systems. However, reviewing behavior varies between domains

and reviewing systems, hereby limiting the generalization of findings to other domains (Kuznetsov et al., 2024b). Additionally, using publicly available reviewing data for evaluation poses risks of contamination (Balloccu et al., 2024; Jiang et al., 2024; Ravaut et al., 2024). In this work, we introduce ARR-Dv2, a new peer review corpus collected from ACL Rolling Review (ARR) – a closed reviewing system that ensures its data has not been seen during training of current LLMs – to support the development of peer review assistance systems tailored to reviewing in natural language processing (NLP) research. Figure 1 shows the structure of ARR-Dv2.

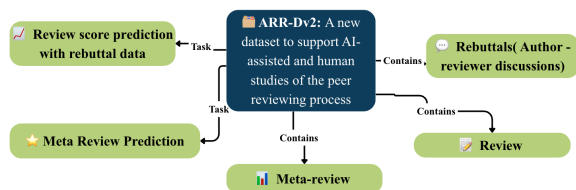


Figure 2: ARR-Dv2 Dataset

This new corpus complements previous work by Dycke et al. (2023) covering NLP reviewing data prior to 2023. However, NLPEER includes only review reports without author rebuttals and reviewer discussion or meta-reviews. ARR-Dv2 includes these important reviewing artifacts (see Figure 2) and is the largest and most recent peer review dataset from the ACL community so far, including the conference of Empirical Methods in Natural Language Processing (EMNLP) 2024, Nations of the Americas Chapter of the ACL (NAACL) 2025, and the International Conference on Computational Linguistics (COLING) 2025. We make the data easily accessible for further research. The dataset enables studying the submitted and accepted version of each paper, their reviews, and the reviewer-author discussion with meta-reviews with open licenses and ethical clearance.

Based on the newly collected ARR-Dv2, we quantify the reviewers’ behavior across time and venues. We find that reviewing in NLP and AI across time does show a focus on similar aspects of papers, such as novelty and experimentation, but with varying importance. Additionally, we turn to a practical assistance task during peer review, *review score prediction*, aiming to support reviewers in choosing adequate scores for their textual report; this is particularly challenging for junior reviewers and reviewers from other fields, and may ease their adaptation to the NLP domain. We use the new reviewer-author

discussion (rebuttals) and meta-review data from our ARR-Dv2, to investigate these artifacts’ effect on score prediction performance. In summary, we contribute

- an extension to the ethically and legally aware collection protocol of Dycke et al. (2022),
- a new corpus of papers, reviews, meta-reviews and author-reviewer discussions from the NLP community from 2024-2025,
- a new perspective on rebuttal and meta-review data during review score prediction lying the groundwork for the systematic investigation of rebuttals’ effect on reviewers’ scoring, and
- the exploration of the task of meta-review score prediction in the ACL community.

2 Related Work

Peer Reviewing Corpora have been published over the past years. Most datasets originate from machine learning (ML) conferences such as ICLR and NEURIPS (Staudinger et al., 2024) crawled from open venues hosted on OpenReview¹ (Yuan et al., 2022; Kennard et al., 2022, i.a.) with the most recent and complete version presented by Zhang et al. (2025). Peer review data from other domains includes closed corpora, such as the Elsevier corpus (Willems, 2022), and open reviewing platforms, such as F1000Research², which are included in NLPEER (Dycke et al., 2023). Unlike most other datasets, ARR-Dv2 includes peer reviews from a closed-reviewing system, which is currently under-represented despite its dominance in science (Dycke et al., 2023) and ensures contamination-free evaluation, as the reviews cannot have been part of the pretraining data of current LLMs yet. Only the PeerRead corpus (Kang et al., 2018) and the ARR and COLING subsets of NLPEER (Dycke et al., 2023) cover peer reviewing data from the NLP community. Both corpora do not include meta-reviews or author responses, and they cover substantially fewer papers and reviews compared to ARR-Dv2. Table 1 summarizes the key differences between NLPeer and ARR-Dv2.

Review Score Prediction aims to infer the numerical review scores (e.g., Soundness, Excitement) typically accompanying review reports in the NLP community. These scores are critical to informing meta-reviewers and acceptance decisions. This task has gained increasing attention in recent

¹<https://openreview.net>

²<https://f1000research.com/>

years. Kang et al. (2018) are the first to study review score prediction, training a simple regression model. Ribeiro et al. (2021) investigate the feasibility of review score prediction based on sentiment labeling of review texts and SVMs. Bharti et al. (2021) are the first to propose a neural model that combines both the paper and review text to predict review scores. Dycke et al. (2023) evaluate BERT-based models for review score prediction across domains. Fernandes and Vaz-de Melo (2024) conduct a large-scale evaluation of review score prediction models and identify specific inputs that challenge current models. Finally, review score prediction accuracy also has become an important metric for assessing the quality of automatic review generation systems (Yuan and Liu, 2022; Lu et al., 2024; Zhou et al., 2024; Idahl and Ahmadi, 2025; Zhu et al., 2025; Weng et al., 2025). With the newly collected ARR-Dv2 dataset, it is now possible to extend this task to meta-review score prediction using data from the ACL community, whereas prior work has primarily relied on review data from other communities, most notably the machine learning community (Bharti et al., 2021; Lu et al., 2024; Zhu et al., 2025).

3 ARR-Dv2

The following section presents details on the data collection process, the resulting dataset, and initial analyses.

3.1 Data Collection

Peer reviewing in the ACL community is a highly confidential process, where great care is taken to keep the information about authors and reviewers secret. Dycke et al. (2022) presents a donation-based data collection protocol — called the Yes-Yes-Yes- or 3Y-workflow — designed for collecting data from sensitive domains with multiple stakeholders and applying it to peer reviews. Their protocol, developed in consultation with the ACL ethics committee, offers an approved framework for collecting and handling peer review data. Our work extends the existing protocol to include author-reviewer discussions and meta-reviews in the licensing workflow, thereby increasing coverage of the donated artifacts. All modifications have been approved and developed in coordination with the ACL Ethics Committee. Figure 3, summarizes the data collection procedure.



Figure 3: The procedure and checks we performed to ensure that the data collection was as ethically sound as possible.

Extended Protocol The 3Y-Workflow is a default opt-out three-step workflow that first asks reviewers to donate their reports, then checks for acceptance of papers, and finally reaches out to authors of accepted papers for their approval of data release. The coverage of donated papers is low, which is arguably attributed to the long time period between submitting a paper and its acceptance, disincentivising authors to donate their review data and drafts. We modify the procedure by asking authors for a donation decision during submission with an option to revise this decision throughout the review process.

As a second major change, we extend the license transfer agreements of the 3Y-workflow to include the author-reviewer-discussion and rebuttals asking authors and reviewers for license transfer for these artifacts during donation. Analogous to reviewers, we ask meta-reviewers to donate their meta-reviews during the reviewing period if the authors approve. All steps involve explicit disclaimers for all stakeholders. The donation of data is voluntary and can be revoked at any time. See Appendix A for more details on the collection workflow.

Collection We collected ARR-Dv2 in the time period from 2024 to 2025. In accordance with the ACL Rolling Review board, we retrieved the reviewing data using the OpenReview API based on the codebase of Dycke et al. (2022). Compared to the earlier NLPeer dataset, which covers 476 ARR submissions and 684 reviews, ARR-Dv2 offers a substantial expansion: it includes 2838 papers and 1923 reviews, with 52.5% of papers accompanied by rebuttals.

Coverage Voluntary data donation results in sub-selection, which limits coverage relative to the full set of review data. We quantify coverage as the proportion of donated papers among all accepted papers at each venue. Higher coverage is desirable and helps reduce self-selection bias in the dataset (Dycke et al., 2022). Our dataset covers 23.8% of accepted papers for EMNLP 2024, 19.4% for

	ARR@NLPeer	ARR-Dv2
time	2021-2022	2024-2025
#papers	476	2838
% accepted	100	100
#reviews	684	1923
#papers w/ rebuttals	-	1403(52.5 %)
Avg. sent per rebuttal	-	15.0 ± 13.4
#reviews per paper	1.43 ± 1.1	0.68 ± 0.75
#sent per paper	220 ± 62.9	538.9 ± 199.6
#sent per review	31 ± 27.0	22.9 ± 12.3
total #tok in reviews	266k	822k
total #tok in rebuttals	-	1.8M

Table 1: Comparison of NLPeer (Dycke et al., 2023) and ARR-Dv2 in terms of papers, reviews, rebuttals, etc. We report the mean and standard deviation of per-paper and per-review statistics.

247 NAACL 2025, and 22% for COLING 2025. This
248 means that ARR-Dv2 covers, on average, roughly
249 20% of the accepted submissions, a rate nearly dou-
250 ble that of the previous collection (approx. 12%).

251 3.2 Meta-data Analysis

252 We analyze the metadata of ARR-Dv2 collected
253 from EMNLP 2024, COLING 2025, and NAACL
254 2025 to determine differences from the previously
255 collected data from the ACL community, also with
256 respect to self-selection bias.

257 **Comparison to NLPEER** As shown in Table 1,
258 ARR-Dv2 covers more data compared to the ARR
259 subset of NLPEER. The number of papers in-
260 creased from 476 to 2838, and reviews grew nearly
261 threefold from 684 to 1923. ARR-Dv2 includes
262 rebuttals for more than half of the papers (52.5%),
263 with an average of 15 sentences per rebuttal. No-
264 tably, the papers from 2024/25 are longer, while
265 their reviews are shorter, compared to the data up
266 to 2022. This may be because reviewers had more
267 papers to review as the number of submissions in-
268 creased, and also because the papers themselves
269 often included longer appendices and reference
270 lists.

271 **Review Score Distribution** Review scores are a
272 useful proxy for reviewer sentiment. We estimate
273 their distribution in the dataset to quantify the diver-
274 sity of papers with respect to quality and reviews
275 in terms of sentiment in ARR-Dv2.

276 Figure 4 shows the distribution of reviewers over-
277 all scores (1-5) for each venue in ARR-Dv2. We
278 observe that the distribution of the scores is more
279 shifted toward lower values for COLING 2025,
280 whereas NAACL 2025 and EMNLP 2024 concen-
281 trate around 3 to 4, and the highest score (5) is

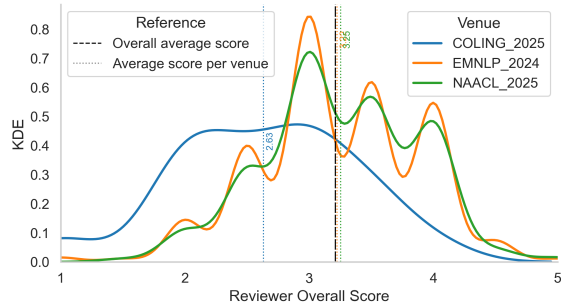


Figure 4: Distribution of review overall scores by venue. Kernel density estimates (KDE) curves show the score distribution for each venue. Colored dotted lines show the average review score per venue (COLING \approx 2.63, EMNLP \approx 3.23, NAACL \approx 3.25), while the black dashed line indicates the overall average scores across all reviews.

282 rarely given to papers.

283 **Meta-review Score Distribution** A meta-review,
284 usually written by the Area Chair (AC), summa-
285 rizes the reviews and discussion for a submitted
286 paper (Wu et al., 2022). Figure 5 shows their distri-
287 bution in ARR-Dv2: the typical score, how spread
288 out the scores are, and whether different venues
use the scoring scale differently. The meta-review

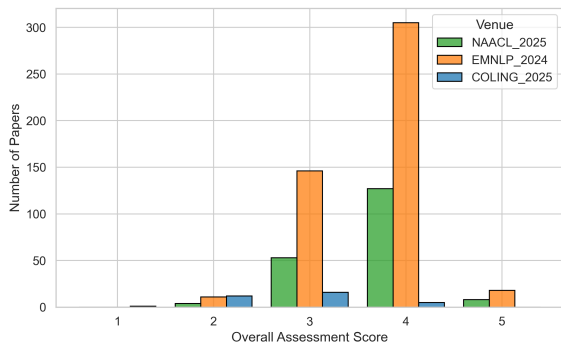


Figure 5: Number of papers by meta-review overall assessment score distribution for each venue.

289 scores for EMNLP 2024 and NAACL 2025 are pri-
290 marily at 4, with some at 3, while COLING 2025
291 has more papers scoring between 2 and 3 and fewer
292 at 4. We observe that some papers in the dataset
293 have notably low meta-review overall assessment
294 scores and were nevertheless accepted. The number
295 of low-scoring papers accepted by each conference
296 are shown in Table 2.

297 **Scoring Behavior** We are interested in scoring
298 behavior to explore if evaluations across venues are
299 comparable and to understand whether confidence
300 and soundness relate to overall scores. We want to
301

venue	lowest score	#main	#findings
EMNLP 2024	2	2	9
COLING 2025	1	1	-
	2	12	-
NAACL 2025	2	1	3

Table 2: The lowest meta-review overall assessment scores in each venue (**lowest score**), the number of main conference acceptances (**#main**) and findings acceptances (**#findings**). COLING 2025 does not have findings proceedings.

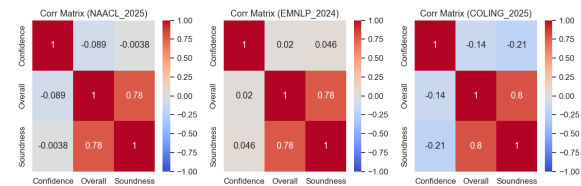


Figure 6: Review-level Pearson correlations between reviewer confidence, overall score, and soundness by venue.

answer the following research question by exploring scoring behavior. **RQ**: How do soundness and reviewer confidence relate to the overall score, and does this relationship change by venue?

Figure 6 shows the correlation between the various scores. In all three venues, there is a strong positive correlation between overall score and soundness (NAACL (0.78), EMNLP (0.78), COLING (0.80)), confirming that reviewers’ overall assessments are closely tied to their evaluation of a paper’s correctness. However, the relationship between reviewer confidence and scoring varies more across venues. At NAACL and EMNLP, the correlations between confidence and both overall and soundness scores are near zero (ranging from -0.09 to 0.046), indicating a weak relationship. In contrast, COLING shows a moderate negative correlation between confidence and soundness (-0.21) and overall score (-0.14), suggesting that more confident reviewers may tend to be slightly more critical. This variation at COLING could point to either differences in reviewer behavior or a possible mismatch between self-assessed confidence and actual scoring patterns in that venue.

3.3 Reviewer Focus Analysis

Review aspects are characteristics of a paper that a reviewer makes a judgment on when evaluating the paper (Lu et al., 2025). The distribution of review aspects measures the importance of different paper dimensions to the reviewers, useful for quantifying

differences in review behavior across domains and time. In this section, we analyze reviewer focus using the NLPeer dataset, the ICLR dataset, and the newly collected ARR-Dv2 dataset. We present a cross-temporal and cross-venue comparison to show how reviewer focus differs.

Aspect Tagging We perform an aspect-based analysis of the newly collected data using the scheme and tagger by Lu et al. (2025). Their scheme includes two sets of labels based on granularity: the **COARSE** and **FINE** label sets. Table 7 in Appendix C gives examples of aspects for both.

Table 3 shows the average number of aspects identified per review for each dataset. The number of aspects per review is generally similar for both the **COARSE** and **FINE** label sets.

dataset	COARSE	FINE
NLPeer ARR22	6.28	9.74
NLPeer COLING2020	7.62	13.42
ARR-Dv2 EMNLP 2024	6.35	10.71
ARR-Dv2 COLING 2025	6.31	11.12

Table 3: Average number of aspects identified per review for each dataset, using the **COARSE** and **FINE** label sets.

Figure 7 shows the five most frequent **COARSE** and **FINE** aspects identified in the reviews for each venue. We also include ICLR 2025 data³ for comparison, which is not part of the ARR-Dv2 data. For **COARSE** aspects, older data (NLPeer ARR-22 and COLING 2020) show almost identical distributions to the new data. This suggests that the reviewer’s focus remains largely stable over time. One notable exception is NLPeer COLING 2020, which includes the aspect **DDDDEI** (Definition, Description, Detail, Discussion, Explanation, Interpretation), suggesting an emphasis on explainability and interpretability. In all venues, the most frequent **COARSE** aspects are Methodology, Result, Data/Task, and Presentation. There is no major difference between ACL venues and ICLR 2025, indicating similar reviewer focus between the NLP and ML communities at a **COARSE** level.

FINE aspects show similar results to the **COARSE** aspect: the reviewer’s focus remains stable across time and venues. Common aspects such as Model, Data, Experiment, and Method are consistently mentioned across both earlier and more recent venues. Note that ICLR 2025 has a more balanced aspect distribution. This may suggest that ICLR

³Publicly available through the [OpenReview API](#).

reviewers place more uniform emphasis across different aspects.

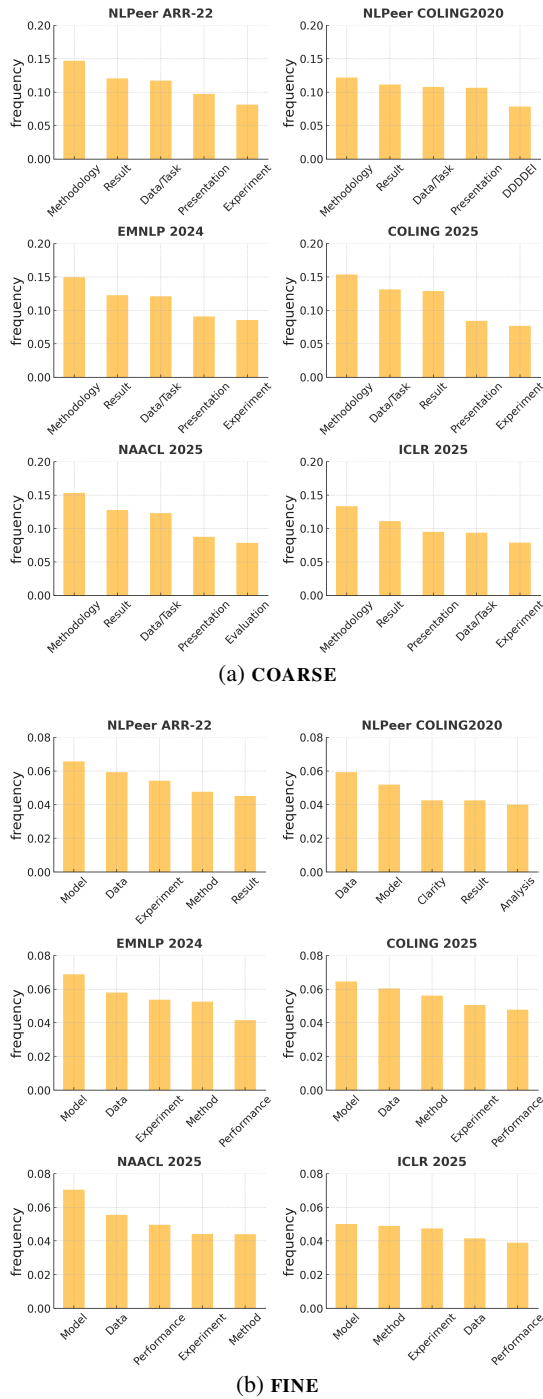


Figure 7: The 5 most frequent aspects in the reviews for each dataset using the COARSE and FINE label sets.

Figure 8 shows the five most frequent FINE aspects related to Methodology in the reviews. Reviewers consistently focus on Model and Method, especially in EMNLP 2024 and NAACL 2025. Compared to ACL venues, ICLR 2025 shows a slightly lower emphasis on Model and Method.

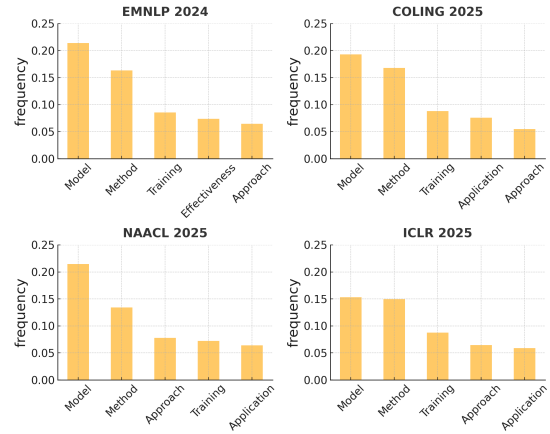


Figure 8: The 5 most frequent FINE aspects of Methodology in the reviews.

4 Review Score Prediction

In review score prediction, the goal is to predict the numerical review score for a given review report plus potential additional artifacts, such as the review template and scoring rubrics. Solving this task requires modeling the relationship between the review text and the overall reviewer opinion which may depend on the reviewing venue and reviewing habits. Here, review score prediction models can serve as assistant tools for reviewers to make more consistent and accurate scoring decisions.

ARR-Dv2 allows us to extend previous work in two important ways: 1) Review score prediction incorporates discussion data (comments by authors and reviewers regarding the initial review and clarification of those comments). This allows us to explore whether incorporating rebuttal data improves the accuracy of review score prediction, which provides insights into how rebuttals influence reviewers' scoring. 2) Meta-review score prediction gives insight to the relationship of scores and text for meta-reviews, too.

4.1 Experimental Settings

For all experiments, we provide the review text together with venue-specific review rubrics in the prompt. In addition to predicting scores, we instruct the model to generate justifications for its score predictions to encourage consistency (see Tables 8 and 9 in the Appendix for the prompts). All experiments are conducted in a zero-shot setting using Deepseek-R1-Distill-Qwen-14B (?), Qwen3-(14B & 32B) (Yang et al., 2025), LLaMA3-(3B & 8B) (Grattafiori et al., 2024). We select these models as representative recent open-source models with strong reported performance.

Individual Review Score Prediction. We consider two input settings: the individual review alone (**review**) and the individual review plus its corresponding rebuttal (**review+rebuttal**).

Meta-Review Score Prediction. This setting is the major extension of this work, where we do not have any comparison data. We use the same setup as in individual review score prediction and use it as an initial test.

Evaluation. We evaluate model performance using exact match accuracy (**acc@0.0**) and the percentage of score predictions that differ from the human-assigned score by 0.5 points (**acc@0.5**). We also report **sum**, which denotes the sum of **acc@0.0** and **acc@0.5**. Together, these metrics offer an interpretable view of model performance: **acc@0.0** is the most strict metric, while **acc@0.5** allows for a small, reasonable margin of error, and **sum** provides a more balanced view. In addition, we report the rooted mean squared error (**RMSE**) and macro F1 score. See Appendix D.1 for more details on the experimental settings.

4.2 Results

Table 4 shows the review score prediction results.

Individual Review Score Prediction. Overall, the models achieve reasonably strong agreement with human-assigned scores. A large number of predictions fall within 0.5 points of human-assigned scores (i.e., **sum** is above 0.8). Model size generally helps: larger models tend to yield better performance, and Qwen3-32B is the best performing model in most cases. The comparison between R1-Qwen-14B and Qwen3-14B shows that the R1-distilled or reasoning-enhanced objective does not improve performance on this task.

Interestingly, incorporating the rebuttal into the input (**review+rebuttal**) does not seem to improve performance. One possible reason is that rebuttals tend to be positive, which may bias the model toward higher score predictions. To further investigate this, we conduct a more detailed analysis of the results, distinguishing between reviews that include reviewer responses (where the reviewer replies, discusses with the author, or acknowledges the author’s responses) and those without reviewer responses (among the 1923 reviews, 878 include reviewer responses). Table 5 shows results on the subset of reviews with reviewer responses. When reviewer responses are present, using **re-**

input field	model	acc@0.0	acc@0.5	sum	RMSE	f1	
M-R	OA	majority	0.64	-	-	0.65	0.20
		random	0.21	-	-	1.70	0.14
		R1-Qwen-14B	<u>0.69</u>	-	-	0.57	0.32
		Qwen3-14B	0.71	-	-	<u>0.54</u>	0.38
		Qwen3-32B	0.71	-	-	0.53	<u>0.35</u>
		LLaMA3-3B	0.63	-	-	0.66	0.17
	LLaMA3-8B	0.67	-	-	0.58	0.31	
R	S	majority	0.31	0.33	0.64	0.72	0.05
		random	0.10	0.20	0.30	1.75	0.07
		R1-Qwen-14B	0.30	<u>0.46</u>	0.76	0.65	0.10
		Qwen3-14B	0.38	<u>0.46</u>	0.85	0.53	0.17
		Qwen3-32B	<u>0.35</u>	<u>0.46</u>	0.81	<u>0.57</u>	<u>0.12</u>
		LLaMA3-3B	0.24	<u>0.46</u>	0.70	0.73	0.10
	LLaMA3-8B	0.34	0.47	<u>0.82</u>	<u>0.57</u>	<u>0.12</u>	
R+R	OA	majority	0.32	0.38	0.70	0.68	0.05
		random	0.12	0.18	0.30	1.72	0.08
		R1-Qwen-14B	0.24	<u>0.48</u>	0.73	0.71	0.09
		Qwen3-14B	0.36	0.45	0.81	0.57	0.15
		Qwen3-32B	0.36	0.44	<u>0.80</u>	<u>0.58</u>	<u>0.14</u>
		LLaMA3-3B	<u>0.33</u>	0.40	0.74	0.65	0.10
	LLaMA3-8B	0.30	0.49	0.79	0.61	0.11	
R	S	majority	0.31	0.33	0.64	0.72	0.05
		random	0.10	0.20	0.30	1.75	0.07
		R1-Qwen-14B	0.28	0.37	0.65	0.75	0.08
		Qwen3-14B	<u>0.32</u>	<u>0.41</u>	<u>0.74</u>	<u>0.65</u>	<u>0.11</u>
		Qwen3-32B	0.33	0.42	0.75	0.64	<u>0.11</u>
		LLaMA3-3B	0.15	0.38	0.53	0.93	0.04
	LLaMA3-8B	0.33	0.34	0.68	0.70	0.13	
R+R	OA	majority	0.32	<u>0.38</u>	<u>0.70</u>	0.68	0.05
		random	0.12	0.18	0.30	1.72	0.08
		R1-Qwen-14B	0.19	<u>0.38</u>	0.58	0.91	0.07
		Qwen3-14B	0.36	<u>0.38</u>	0.75	<u>0.65</u>	<u>0.14</u>
		Qwen3-32B	<u>0.35</u>	0.40	0.75	0.64	0.15
		LLaMA3-3B	0.27	0.32	0.60	0.83	0.08
	LLaMA3-8B	0.32	0.35	0.68	0.74	0.13	

Table 4: Results of review score prediction under different input settings: meta-review (M-R), review (R), review+rebuttal (R+R). **sum** denotes the sum of **acc@0.0** and **acc@0.5**. The meta-review rubrics do not have 0.5-point scoring, so **acc@0.5** and **sum** are marked as “-.” **Bold** indicates the best performance and underline indicates the second best. “OA” is Overall Assessment and “S” is Soundness.

view+rebuttal achieves higher **acc@0.0** on Overall Assessment score prediction than using the review alone. This suggests that reviewer responses help mitigate the positivity bias in the rebuttals.

Meta-Review Score Prediction. Overall, the models achieve good performance, with **acc@0.0** around 0.7 for most models. Larger models generally perform better, while reasoning-enhanced objective does not lead to improvements. Note that meta-review scores are integer values and do not have decimal-based scoring.

Positivity Bias. Results are shown in Table 6, which reports the percentage of cases where the model’s predicted score is higher (+%) or lower

field	model	acc@0.0	
		r	r+r
Soundness	R1-Qwen-14B	0.31	0.31
	Qwen3-14B	0.40	0.35
	Qwen3-32B	0.36	0.36
	LLaMA3-3B	0.25	0.16
	LLaMA3-8B	0.34	0.32
Overall Assessment	R1-Qwen-14B	0.26	0.25
	Qwen3-14B	0.38	0.42
	Qwen3-32B	0.38	0.41
	LLaMA3-3B	0.34	0.30
	LLaMA3-8B	0.31	0.37

Table 5: Results of review score prediction on reviews with reviewer responses.

(-%) than the human-assigned score. Positivity bias is much lower for meta-reviews than for other inputs. When rebuttal is included, models show a stronger positivity bias compared to using the review alone. This further supports our hypothesis that rebuttals introduce a positivity bias in model predictions.

input	field	model	+%	-%
M-R	OA	R1-Qwen-14B	21.56	8.83
		Qwen3-14B	14.16	14.31
		Qwen3-32B	17.14	11.05
		LLaMA3-3B	30.73	5.70
		LLaMA3-8B	25.70	6.70
R	S	R1-Qwen-14B	49.48	20.42
		Qwen3-14B	33.66	28.04
		Qwen3-32B	47.68	16.66
		LLaMA3-3B	63.02	12.57
		LLaMA3-8B	43.73	21.29
R+R	OA	R1-Qwen-14B	54.32	21.09
		Qwen3-14B	34.13	29.44
		Qwen3-32B	40.29	23.27
		LLaMA3-3B	46.56	19.75
		LLaMA3-8B	50.57	19.39
R+R	S	R1-Qwen-14B	61.98	9.79
		Qwen3-14B	55.47	11.87
		Qwen3-32B	56.87	9.57
		LLaMA3-3B	80.10	4.59
		LLaMA3-8B	59.30	7.02
R+R	OA	R1-Qwen-14B	71.56	8.75
		Qwen3-14B	47.50	16.04
		Qwen3-32B	52.91	11.34
		LLaMA3-3B	65.97	6.26
		LLaMA3-8B	61.40	5.79

Table 6: Proportions of outputs where the LLM predicted score is higher (+%) or lower (-%) than the human-assigned score.

We further observe that for COLING 2025, models show a strong positivity bias across all tasks and input settings (see Table 10 in the Appendix). One possible reason is incomplete data coverage⁴. An-

⁴COLING 2025 includes both direct and ARR submissions,

other possible explanation is that the review scores in COLING 2025 are relatively low compared to other venues (see Figures 4 and 5 in Section 3.2). As a result, models tend to overestimate scores, leading to a more noticeable positivity bias.

Our analysis and interpretation assume that human-assigned scores are the gold standard. However, in some cases, human-assigned scores do not fully align with the corresponding review texts. From this perspective, the observed positivity bias in model predictions may also reflect mismatches between review texts and human-assigned scores. To further investigate this possibility, we conduct a small-scale qualitative analysis of cases with large discrepancies between human-assigned and LLM-predicted scores (see Table 11 in the Appendix). We find that while some discrepancies indeed stem from inconsistencies between the review text and the human-assigned scores, in most cases the model predictions tend to be overly positive.

5 Conclusion and Future Work

We introduced ARR-Dv2, a peer review dataset that considerably extends previous data collections in terms of quantity and artifacts, including data from the rebuttal phase and meta-reviews. We analyzed ARR-Dv2 on differences between venues, for example, regarding the link between the overall, confidence and soundness scores. Additionally, we compared the dataset to previous data collections regarding the aspect distribution in reviews. Further, we conducted experiments on review score prediction, observing that the inclusion of rebuttal information aids in score prediction. Finally, we extended review score prediction to the task of meta-review score prediction, which shows promising initial results on which we can build in the future.

As we have not done any fine-tuning in our experiments so far, this would be an obvious next step. Future work also includes the collection of unpublished data, as the current data collections show a strong bias towards accepted papers. Additionally, we are able to tackle various research questions and peer review tasks outlined by Kuznetsov et al. (2024b), such as (meta-) review summarization and providing feedback to reviewers on their scores based on their written reviews.

while our data only covers ARR submissions (we only have access to ARR submissions).

537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587

Limitations

The data analysis and review score prediction experiments may be influenced by potential bias. The current dataset is not comprehensive, as a number of reviews are missing due to constraints on reviewer consent, and our dataset only includes reviews for accepted papers (see Section 3.1 for coverage statistics). As a result, the statistics may over-represent characteristics specific to accepted submissions, and the experimental results may not accurately reflect the models’ true predictive capability across a broader range of reviews. Our study focuses on data analysis and score prediction tasks which enable us to measure representativeness of the data and relevance of its constituent the artifacts. We do not consider generation tasks, such as review or rebuttal generation, due to the lack of objective and standardized evaluation metrics for these scenarios. We leave the study of these tasks for future work.

Additionally, our main goal and focus in this work was to introduce and analyse the ARR-Dv2 dataset, rather than to optimise model performance. Our experiments are intentionally lightweight and use zero-shot inference as a baseline to illustrate how the new data, especially rebuttals and meta-reviews, can be used. Fine-tuning and architecture-specific optimization would shift the paper’s scope from dataset and analysis to model development, which is outside this initial release. We also acknowledge, that our limitation to open source models does not allow for conclusions to closed-source state of the art models. Nevertheless, open source models enables local deployment and full control, which are crucial for ensuring confidentiality and data protection in the peer-review context. We do not claim that open-weight models are inherently ethically or legally compliant; rather, open access makes it possible to meet compliance obligations by avoiding transmission of unpublished manuscripts or reviews to third-party servers. In this study, we used only accepted and consented data, but in real-world deployment, AI-assisted reviewer tools process unpublished, confidential submissions. In such settings, using closed-source, cloud-hosted APIs even if private would still need sharing sensitive materials externally, contrary to the confidentiality principles outlined in the ACL Publication Ethics. Open-weight models, by contrast, enable on-premise, auditable, and reproducible deployment that aligns with these ethical

requirements. 588

Reproducing the data collection requires approval by the ACL Ethics Committee as well as support from the ACL Exec. Both can be readily obtained. In order to reproduce the experimental, technical aspect of our work we refer to the Appendix, where the technical details (models, versions, parameters and hardware specifications) are given. 589
590
592
593
594
595
596

Ethical Considerations 597

This work is associated with several ethical issues. First, the data collection itself and its publication come with their own ethical and legal challenges. We addressed those in the paper. Second, we are not promoting the automation of the peer review process. On the contrary, science requires feedback from peers, and the *peer* review process is at the core of this feedback mechanism. What we are suggesting is using technical *support* to reduce the workload on the reviewers while keeping the human reviewers in the loop. 598
599
600
601
602
603
604
605
606
607
608

References 609

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*. 610
611
612
613
614

Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dusek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source llms](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian’s, Malta, March 17-22, 2024*, pages 67–93. Association for Computational Linguistics. 615
616
617
618
619
620
621
622

Prabhat Kumar Bharti, Shashi Ranjan, Tirthankar Ghosal, Mayank Agrawal, and Asif Ekbal. 2021. Peerassist: leveraging on paper-review interactions to predict peer review decisions. In *Towards Open and Trustworthy Digital Societies: 23rd International Conference on Asia-Pacific Digital Libraries, ICADL 2021, Virtual Event, December 1–3, 2021, Proceedings 23*, pages 421–435. Springer. 623
624
625
626
627
628
629
630

Aliaksandr Birukou, Joseph R. Wakeling, Claudio Bartolini, Fabio Casati, Maurizio Marchese, Katsiaryna Mirylenka, Nardine Osman, Azzurra Ragone, Carles Sierra, and Aalam Wassef. 2011. [Alternatives to peer review: Novel approaches for research evaluation](#). *Frontiers in Computational Neuroscience*, volume 5 - 2011. 631
632
633
634
635
636
637

638	Lutz Bornmann. 2011. Scientific peer review . <i>Annual Review of Information Science and Technology</i> , 45(1):197–245.	
639		
640		
641	Jiangshu Du, Yibo Wang, Wenting Zhao, Zhongfen	Filip Radenovic, Francisco Guzmán, Frank Zhang,
642	Deng, Shuaiqi Liu, Renze Lou, Henry Peng Zou,	Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-
643	Pranav Narayanan Venkit, Nan Zhang, Mukund Srin-	derson, Govind Thattai, Graeme Nail, Gregoire Mil-
644	nath, Haoran Ranran Zhang, Vipul Gupta, Yinghui	alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
645	Li, Tao Li, Fei Wang, Qin Liu, Tianlin Liu, Pengzhi	Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
646	Gao, Congying Xia, Chen Xing, Cheng Jiayang,	Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-
647	Zhaowei Wang, Ying Su, Raj Sanjay Shah, Ruohao	han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
648	Guo, Jing Gu, Haoran Li, Kangda Wei, Zihao Wang,	Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
649	Lu Cheng, Surangika Ranathunga, Meng Fang, Jie	Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
650	Fu, Fei Liu, Ruihong Huang, Eduardo Blanco, Yixin	Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,
651	Cao, Rui Zhang, Philip S. Yu, and Wenpeng Yin.	Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
652	2024. LLMs assist NLP researchers: Critique paper	Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-
653	(meta)-reviewing . In <i>Proceedings of the 2024 Con-</i>	teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,
654	<i>ference on Empirical Methods in Natural Language</i>	Kartikaya Upasani, Kate Plawiak, Ke Li, Kenneth
655	<i>Processing</i> , pages 5081–5099, Miami, Florida, USA.	Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
656	Association for Computational Linguistics.	Kshitiz Malik, Kuenley Chiu, Kunal Balla, Kushal
657	Nils Dycke, Iliia Kuznetsov, and Iryna Gurevych. 2022.	Lakhotia, Lauren Rantala-Yearly, Laurens van der
658	Yes-yes-yes: Proactive data collection for ACL	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,
659	rolling review and beyond . In <i>Findings of the Associ-</i>	Louis Martin, Lovish Madaan, Lubo Malo, Lukas
660	<i>ation for Computational Linguistics: EMNLP 2022</i> ,	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline
661	pages 300–318, Abu Dhabi, United Arab Emirates.	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar
662	Association for Computational Linguistics.	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew
663	Nils Dycke, Iliia Kuznetsov, and Iryna Gurevych. 2023.	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-
664	NLPeer: A unified resource for the computational	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,
665	study of peer review . In <i>Proceedings of the 61st An-</i>	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-
666	<i>annual Meeting of the Association for Computational</i>	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,
667	<i>Linguistics (Volume 1: Long Papers)</i> , pages 5049–	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick
668	5073, Toronto, Canada. Association for Computa-	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vas-
669	tional Linguistics.	sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,
670	Saman Ebadi, Hassan Nejadghanbar, Ahmed Rawdhan	Praveen Krishnan, Punit Singh Koura, Puxin Xu,
671	Salman, and Hassan Khosravi. 2025. Exploring the	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj
672	impact of generative ai on peer review: Insights from	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,
673	journal reviewers. <i>Journal of Academic Ethics</i> , pages	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,
674	1–15.	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-
675	Gustavo Lúcius Fernandes and Pedro OS Vaz-de Melo.	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
676	2024. Enhancing the examination of obstacles in an	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-
677	automated peer review system. <i>International Journal</i>	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-
678	<i>on Digital Libraries</i> , 25(2):341–364.	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-
679	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	ran Narang, Sharath Raparthy, Sheng Shen, Shengye
680	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Wan, Shruti Bhosale, Shun Zhang, Simon Vanden-
681	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	hende, Soumya Batra, Spencer Whitman, Sten
682	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	Sootla, Stephane Collot, Suchin Gururangan, Syd-
683	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	dney Borodinsky, Tamar Herman, Tara Fowler, Tarek
684	tra, Archie Sravankumar, Artem Korenev, Arthur	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias
685	Hinsvark, Arun Rao, Aston Zhang, Aurelien Rod-	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal
686	riguez, Austen Gregerson, Ava Spataru, Baptiste	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh
687	Roziere, Bethany Biron, Binh Tang, Bobbie Chern,	Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-
688	Charlotte Caucheteux, Chaya Nayak, Chloe Bi,	ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-
689	Chris Marra, Chris McConnell, Christian Keller,	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-
690	Christophe Touret, Chunyang Wu, Corinne Wong,	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-
691	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinf-
692	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-
693	Danny Wyatt, David Esiobu, Dhruv Choudhary,	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,
694	Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,
695	Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing
696	Elina Lobanova, Emily Dinan, Eric Michael Smith,	Chen, Zoe Papanikos, Aaditya Singh, Aayushi Sri-
		vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,
		Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
		Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
		Baevski, Allie Feinstein, Amanda Kallet, Amit San-
		gani, Amos Teo, Anam Yunus, Andrei Lupu, An-
		dres Alvarado, Andrew Caples, Andrew Gu, Andrew
		Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-

761	dani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	825
762	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	826
763	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	827
764	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	828
765	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	829
766	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	830
767	Brian Gamido, Britt Montalvo, Carl Parker, Carly	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	831
768	Burton, Catalina Mejia, Ce Liu, Changhan Wang,	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	832
769	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	833
770	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	834
771	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	Subramanian, Sy Choudhury, Sydney Goldman, Tal	835
772	Daniel Kreymer, Daniel Li, David Adkins, David	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	836
773	Xu, Davide Testuggine, Delia David, Devi Parikh,	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	837
774	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	Matthews, Timothy Chou, Tzook Shaked, Varun	838
775	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	839
776	Elaine Montgomery, Eleonora Presani, Emily Hahn,	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	840
777	Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	841
778	Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	842
779	Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	843
780	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	844
781	Seide, Gabriela Medina Florez, Gabriella Schwarz,	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	845
782	Gada Badeer, Georgia Swee, Gil Halpern, Grant	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	846
783	Herman, Grigory Sizov, Guangyi, Zhang, Guna	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	847
784	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	848
785	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	849
786	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	Zhiwei Zhao, and Zhiyu Ma. 2024. The Llama 3	850
787	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	Herd of Models . <i>arXiv preprint</i> . ArXiv:2407.21783	851
788	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	[cs].	852
789	Irina-Elena Veliche, Itai Gat, Jake Weissman, James		853
790	Geboski, James Kohli, Janice Lam, Japhet Asher,	Ahmed M Hanafi, Mohamed M Al-mansi, Omar A Al-	854
791	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	Sharif, et al. 2025. Generative ai in academia: A	855
792	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	comprehensive review of applications and implica-	856
793	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	tions for the research process. <i>International Journal</i>	857
794	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	<i>of Engineering and Applied Sciences-October 6 Uni-</i>	858
795	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	<i>versity</i> , 2(1):91–110.	859
796	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	Maximilian Idahl and Zahra Ahmadi. 2025. OpenRe-	860
797	delwal, Katayoun Zand, Kathy Matosich, Kaushik	viewer: A specialized large language model for gener-	861
798	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	ating critical scientific paper reviews . In <i>Proceedings</i>	862
799	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	<i>of the 2025 Conference of the Nations of the Amer-</i>	863
800	Huang, Lailin Chen, Lakshya Garg, Lavender A,	<i>icas Chapter of the Association for Computational</i>	864
801	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	<i>Linguistics: Human Language Technologies (System</i>	865
802	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	<i>Demonstrations)</i> , pages 550–562, Albuquerque, New	866
803	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	Mexico. Association for Computational Linguistics.	867
804	Martynas Mankus, Matan Hasson, Matthew Lennie,	Minhao Jiang, Ken Liu, Ming Zhong, Rylan Schaeffer,	868
805	Matthias Reso, Maxim Groshev, Maxim Naumov,	Siru Ouyang, Jiawei Han, and Sanmi Koyejo. 2024.	869
806	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	Does data contamination make a difference? insights	870
807	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	from intentionally contaminating pre-training data	871
808	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	for language models . In <i>ICLR 2024 Workshop on</i>	872
809	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	<i>Navigating and Addressing Data Problems for Founda-</i>	873
810	Mo Metanat, Mohammad Rastegari, Munish Bansal,	<i>tion Models</i> .	874
811	Nandhini Santhanam, Natascha Parks, Natasha	Dongyeop Kang, Waleed Ammar, Bhavana Dalvi,	875
812	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	Madeleine van Zuylen, Sebastian Kohlmeier, Eduard	876
813	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	Hovy, and Roy Schwartz. 2018. A dataset of peer	877
814	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	reviews (PeerRead): Collection, insights and NLP	878
815	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	applications . In <i>Proceedings of the 2018 Conference</i>	879
816	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	<i>of the North American Chapter of the Association for</i>	880
817	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	<i>Computational Linguistics: Human Language Tech-</i>	881
818	Dollar, Polina Zvyagina, Prashant Ratanchandani,	<i>nologies, Volume 1 (Long Papers)</i> , pages 1647–1661,	882
819	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	New Orleans, Louisiana. Association for Computa-	883
820	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	tional Linguistics.	884
821	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,		
822	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky		
823	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,		
824			

885	Neha Kennard, Tim O’Gorman, Rajarshi Das, Akshay Sharma, Chhandak Bagchi, Matthew Clinton, Pranay Kumar Yelugam, Hamed Zamani, and Andrew McCallum. 2022. DISAPERE: A dataset for discourse structure in peer review discussions . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1234–1249, Seattle, United States. Association for Computational Linguistics.	943
886		944
887		945
888		
889		946
890		947
891		948
892		949
893		950
894		951
895	Ilia Kuznetsov, Osama Mohammed Afzal, Koen Dercksen, Nils Dycke, Alexander Goldberg, Tom Hope, Dirk Hovy, Jonathan K. Kummerfeld, Anne Lauscher, Kevin Leyton-Brown, Sheng Lu, Mausam, Margot Mieskes, Aurélie Névéol, Danish Pruthi, Lizhen Qu, Roy Schwartz, Noah A. Smith, Thamar Solorio, Jingyan Wang, Xiaodan Zhu, Anna Rogers, Nihar B. Shah, and Iryna Gurevych. 2024a. What can natural language processing do for peer review? <i>Preprint</i> , arXiv:2405.06563.	952
896		953
897		954
898		955
899		956
900		
901		957
902		958
903		
904		
905	Ilia Kuznetsov, Osama Mohammed Afzal, Koen Dercksen, Nils Dycke, Alexander Goldberg, Tom Hope, Dirk Hovy, Jonathan K. Kummerfeld, Anne Lauscher, Kevin Leyton-Brown, et al. 2024b. What can natural language processing do for peer review? <i>arXiv preprint arXiv:2405.06563</i> .	959
906		960
907		961
908		962
909		963
910		964
911	Esther Landhuis. 2016. Scientific literature: Information overload. <i>Nature</i> , 535(7612):457–458.	965
912		966
913	Carole J Lee, Cassidy R Sugimoto, Guo Zhang, and Blaise Cronin. 2013. Bias in peer review. <i>Journal of the American Society for information Science and Technology</i> , 64(1):2–17.	967
914		968
915		969
916		970
917	Weixin Liang, Zachary Izzo, Yaohui Zhang, Haley Lepp, Hancheng Cao, Xuandong Zhao, Lingjiao Chen, Haotian Ye, Sheng Liu, Zhi Huang, Daniel Mcfarland, and James Y. Zou. 2024. Monitoring AI-modified content at scale: A case study on the impact of ChatGPT on AI conference peer reviews . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 29575–29620. PMLR.	971
918		972
919		973
920		974
921		975
922		976
923		977
924		978
925		979
926	Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob N. Foerster, Jeff Clune, and David Ha. 2024. The AI scientist: Towards fully automated open-ended scientific discovery . <i>CoRR</i> , abs/2408.06292.	980
927		981
928		
929		
930	Sheng Lu, Ilia Kuznetsov, and Iryna Gurevych. 2025. Identifying aspects in peer reviews . <i>arXiv preprint arXiv:2504.06910</i> .	982
931		983
932		984
933	Mathieu Ravaut, Bosheng Ding, Fangkai Jiao, Hailin Chen, Xingxuan Li, Ruochen Zhao, Chengwei Qin, Caiming Xiong, and Shafiq Joty. 2024. How much are llms contaminated? A comprehensive survey and the llmsanitize library . <i>CoRR</i> , abs/2404.00699.	985
934		986
935		
936		
937		
938	Ana Carolina Ribeiro, Amanda Sizo, Henrique Lopes Cardoso, and Luís Paulo Reis. 2021. Acceptance decision prediction in peer-review through sentiment analysis. In <i>EPIA Conference on Artificial Intelligence</i> , pages 766–777. Springer.	987
939		988
940		989
941		990
942		991
	D Sculley, Jasper Snoek, and Alex Wiltschko. 2018. Avoiding a tragedy of the commons in the peer review process . <i>arXiv preprint arXiv:1901.06246</i> .	992
		993
	Moritz Staudinger, Wojciech Kusa, Florina Piroi, and Allan Hanbury. 2024. An analysis of tasks and datasets in peer reviewing . In <i>Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)</i> , pages 257–268, Bangkok, Thailand. Association for Computational Linguistics.	994
		995
	Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. 2025. Cyclereviewer: Improving automated research via automated review . In <i>The Thirteenth International Conference on Learning Representations</i> .	996
		997
	Linda Willems. 2022. New dataset offers unique insights into peer review . Accessed: 2025-06-06.	998
		999
	Po-Cheng Wu, An-Zi Yen, Hen-Hsen Huang, and Hsin-Hsi Chen. 2022. Incorporating peer reviews and rebuttal counter-arguments for meta-review generation. In <i>Proceedings of the 31st ACM International Conference on Information & Knowledge Management</i> , pages 2189–2198.	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	
	Weizhe Yuan and Pengfei Liu. 2022. Kid-review: knowledge-guided scientific review generation with oracle pre-training . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 11639–11647.	
	Weizhe Yuan, Pengfei Liu, and Graham Neubig. 2022. Can we automate scientific reviewing? <i>Journal of Artificial Intelligence Research</i> , 75:171–212.	
	Daoze Zhang, Zhijian Bao, Sihang Du, Zhiyi Zhao, Kuangling Zhang, Dezheng Bao, and Yang Yang. 2025. ReQ: A consistency-ensured dataset for full-stage peer review and multi-turn rebuttal discussions . <i>arXiv preprint arXiv:2505.07920</i> .	
	Ruiyang Zhou, Lu Chen, and Kai Yu. 2024. Is LLM a reliable reviewer? a comprehensive evaluation of LLM on automatic paper reviewing tasks . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources</i>	

and Evaluation (LREC-COLING 2024), pages 9340–9351, Torino, Italia. ELRA and ICCL.

Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. 2025. [Deepreview: Improving llm-based paper review with human-like deep thinking process](#). *CoRR*, abs/2503.08569.

A More on Data Collection

We use the [OpenReview Python package](#) to collect the data. The data collection process adheres to the Yes-Yes-Yes protocol (Dycke et al., 2022), which involves two key license agreements: **Peer Reviewer Content License Agreement** (hereinafter referred to as the *reviewer agreement*) and **Blind Submission License Agreement** (hereinafter referred to as the *author agreement*). The main steps are as follows:

- **Obtain reviewer agreements**

Use `or_client.get_all_notes(invitation=...)`, where `invitation` is the link to the reviewer agreement in OpenReview.

- **Retrieve accepted papers**

Directly from the [ACL Anthology Repository](#).

- **Collect submission details**

Use `or_client.get_all_notes(invitation=...)`, where `invitation` is the link to the submissions in OpenReview, including reviews, rebuttals, and the author agreements.

- **Filter the submissions**

Iterate through all submissions and select those that meet all of the criteria of the Yes-Yes-Yes protocol.

B Data analysis - soundness and overall score by confidence

Figure 9 shows an analysis of our data with respect to confidence and soundness. To maintain readability, the y-axis is skewed in each of these plots. But we can see that we have very few papers with a confidence of 1 or 2. The majority of the reviewers give their confidence as either 3 or 4, and a few reviewers claim a confidence of 5. We can also see that the distribution focuses around a soundness score of 3, indicating that the distribution approximately follows a Gaussian Distribution.

Figure 10 shows a similar distribution for confidence in relation to the overall score. Few reviewers give a low confidence rating of 1 or 2, while the majority note a confidence rating of 3 or 4, with a smaller group giving a confidence rating of 5.

Again, the scores resemble a Gaussian distribution, with the majority of scores being 3.

Both Figures 9 and 10 show that contrary to the previous data collection our dataset also contains papers that would probably be rejected for the conference, as overall scores of 1 and 2 (a total of 80 papers for NAACL, 182 papers for EMNLP, and 31 papers for COLING) would probably be a clear reject, while papers with 3 (217 papers for NAACL, 657 papers for EMNLP and 26 papers for COLING) are borderline and papers with an overall score of 5 (3 papers for NAACL, 6 papers for EMNLP) are very likely to be accepted. Across all three NLP venues, both the soundness and overall score distributions per reviewer confidence show a strong central tendency in which the vast majority of reviews, regardless of venue or confidence level, assign moderate scores, with very few ratings for 1 or 5.

C More on Data Analysis

Table 7 shows examples of aspects derived by Lu et al. (2025).

COARSE	FINE
Contribution	Contribution
DDDDEI	Definition, Description, Detail, Discussion, Explanation, Interpretation
IJMV	Intuition, Justification, Motivation, Validation
Novelty	Innovation/Novelty/Originality
Presentation	Clarity, Figure, Grammar, Presentation, Typo
Related Work	Citation/Literature/Related Work
Significance	Impact, Importance, Significance
Ablation	Ablation
Analysis	Analysis
Comparison	Comparison
Data/Task	Annotation, Benchmark, Data, Task
Evaluation	Evaluation, Metric
Experiment	Experiment
Methodology	Algorithm, Implementation, Method
Theory	Theory
Result	Findings, Improvement, Performance, Result

Table 7: Examples of aspects in the COARSE and FINE label sets (Lu et al., 2025).

D More on Review Score Prediction

D.1 More on Experiment Settings

For all experiments, we set `temperature=0.5`, `max_tokens=1024`, and `seed=2266`. We used 4-bit quantization for model inference. All experiments

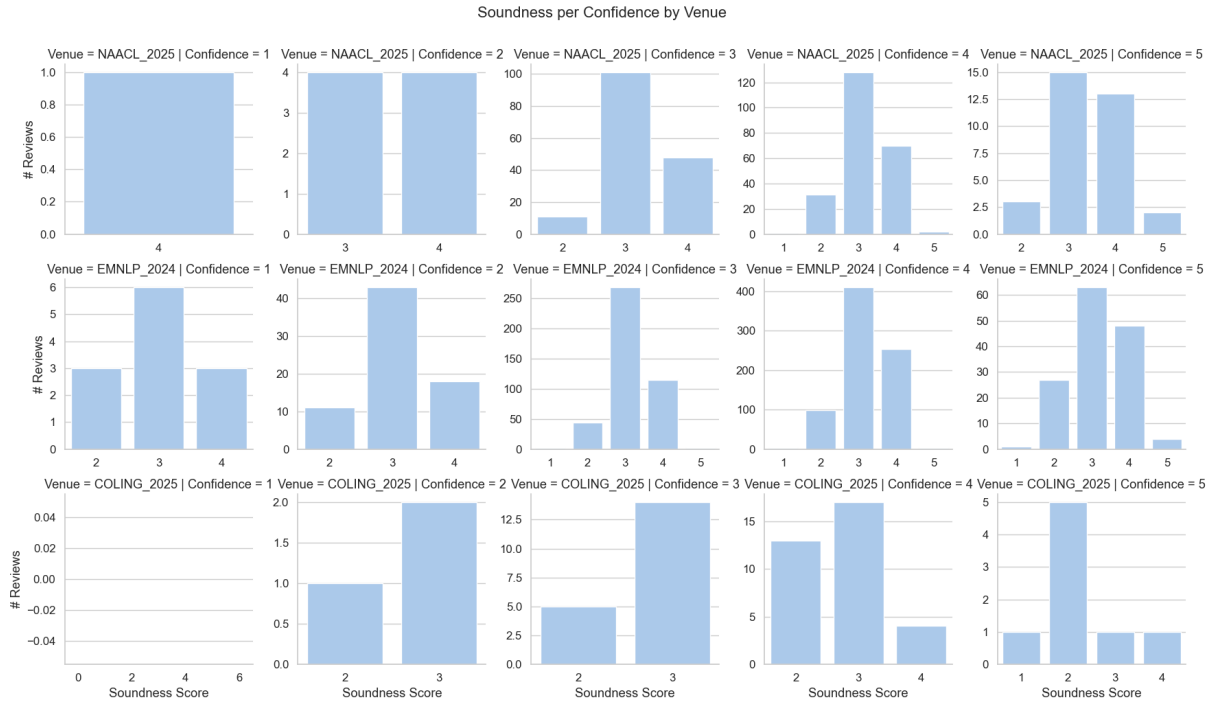


Figure 9: The bars show the average scores for *reviews*, with respect to *confidence* and *soundness* for the papers in our data collection.

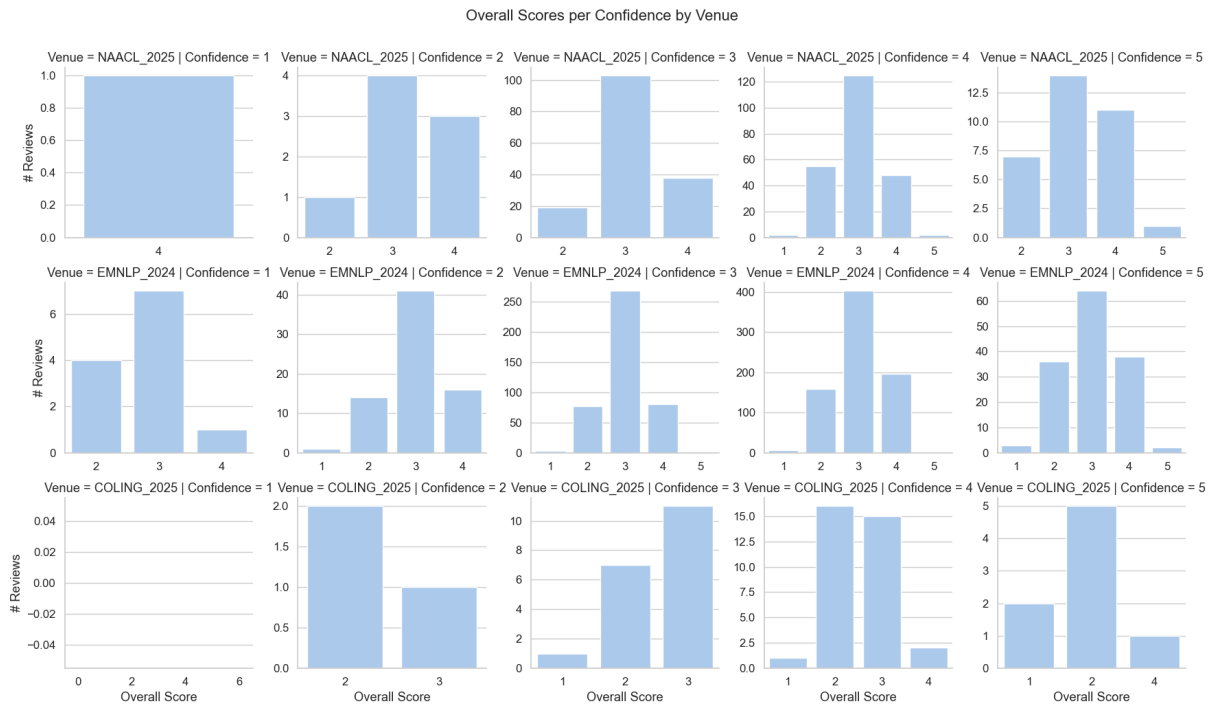


Figure 10: The bars show the average scores for *reviews*, with respect to *confidence* and *overall* score for the papers in our data collection.

1075 were done on either an NVIDIA A100 or H100
 1076 GPU. Tables 8 and 9 show the prompts used in the
 1077 experiments.

D.2 More on Results

Table 10 shows the positivity bias results for COLING 2025. Table 11 shows examples of human-assigned and LLM-predicted review scores with

1078

1079

1080

1081

Given a paper metareview and the corresponding metareview form, recommend appropriate Overall Assessment score that the metareviewer might consider giving based on their review.

Instructions:

1. Read the review text carefully.
2. Analyze the review and recommend an Overall Assessment score.
3. Provide a brief explanation for the recommended score.
4. Output a json dictionary:
`{“overall_assessment_score”: ...,
“overall_assessment_justification”: ...}`
5. Output only the json dictionary and follow the json schema exactly, no extra keys in the output.

Review form:

Overall Assessment

- 5 = The paper is largely complete and there are no clear points of revision
 - 4 = There are minor points that may be revised
 - 3 = There are major points that may be revised
 - 2 = The paper would need significant revisions to reach a publishable state
 - 1 = Even after revisions, the paper is not likely to be publishable at an *ACL venue
-

Table 8: The prompt used for the meta-review score prediction task.

large discrepancies.

Given a paper review and the corresponding review form, recommend appropriate Soundness and Overall Assessment scores that the reviewer might consider giving based on their review.

Instructions:

1. Read the review text carefully.
2. Analyze the review and recommend a Soundness score and an Overall Assessment score.
3. Provide a brief explanation for each recommended score.
4. Output a json dictionary: {"soundness_score": ..., "soundness_justification": ..., "overall_assessment_score": ..., "overall_assessment_justification": ...}
5. Output only the json dictionary and follow the json schema exactly, no extra keys in the output.

Review form:

Soundness

How sound and thorough is this study? Does the paper clearly state scientific claims and provide adequate support for them? For experimental papers: consider the depth and/or breadth of the research questions investigated, technical soundness of experiments, methodological validity of evaluation. For position papers, surveys: consider the current state of the field is adequately represented, and main counter-arguments acknowledged. For resource papers: consider the data collection methodology, resulting data & the difference from existing resources are described in sufficient detail. Please adjust your baseline to account for the length of the paper.

5 = Excellent: This study is one of the most thorough I have seen, given its type.

4.5

4 = Strong: This study provides sufficient support for all of its claims/arguments. Some extra experiments could be nice, but not essential.

3.5

3 = Acceptable: This study provides sufficient support for its major claims/arguments. Some minor points may need extra support or details.

2.5

2 = Poor: Some of the main claims/arguments are not sufficiently supported. There are major technical/methodological problems.

1.5

1 = Major Issues: This study is not yet sufficiently thorough to warrant publication or is not relevant to ACL.

Overall Assessment

Would you personally like to see this paper presented at an *ACL event that invites submissions on this topic? For example, you may feel that a paper should be presented if its contributions would be useful to its target audience, deepen the understanding of a given topic, or help establish cross-disciplinary connections. Note: Even high-scoring papers can be in need of minor changes (e.g. typos, non-core missing refs, etc.).

5 = Top-Notch: This is one of the best papers I read recently, of great interest for the (broad or narrow) sub-communities that might build on it

4.5

4 = This paper represents solid work, and is of significant interest for the (broad or narrow) sub-communities that might build on it

3.5

3 = Good: This paper makes a reasonable contribution, and might be of interest for some (broad or narrow) sub-communities, possibly with minor revisions

2.5

2 = Revisions Needed: This paper has some merit, but also significant flaws, and needs work before it would be of interest to the community

1.5

1 = Major Revisions Needed: This paper has significant flaws, and needs substantial work before it would be of interest to the community

0 = This paper is not relevant to the *ACL community (for example, is in no way related to natural language processing)

Table 9: The prompt used for the individual review score prediction.

input	field	model	+%	-%
meta-review	Overall Assessment	R1-Qwen-14B	61.76	5.88
		Qwen3-14B	47.06	2.94
		Qwen3-32B	50.00	2.94
		LLaMA3-3B	75.00	0.00
		LLaMA3-8B	65.38	3.85
review	Soundness	R1-Qwen-14B	73.44	7.81
		Qwen3-14B	48.44	17.19
		Qwen3-32B	62.50	6.25
		LLaMA3-3B	68.52	7.41
		LLaMA3-8B	60.00	0.00
	Overall Assessment	R1-Qwen-14B	73.44	6.25
		Qwen3-14B	59.38	14.06
		Qwen3-32B	64.06	15.62
		LLaMA3-3B	64.81	12.96
		LLaMA3-8B	80.00	20.00
review+rebuttal	Soundness	R1-Qwen-14B	84.38	0.00
		Qwen3-14B	75.00	3.12
		Qwen3-32B	68.75	3.12
		LLaMA3-3B	90.00	2.00
		LLaMA3-8B	81.25	6.25
	Overall Assessment	R1-Qwen-14B	87.50	0.00
		Qwen3-14B	71.88	14.06
		Qwen3-32B	70.31	10.94
		LLaMA3-3B	84.00	2.00
		LLaMA3-8B	87.50	6.25

Table 10: Proportions of outputs where the LLM predicted score is higher (+%) or lower (-%) than the human-assigned score for COLING 2025.

Human review	Human		LLM		LLM justifications
	S	OA	S	OA	
<p>Strengths: The paper addresses formatting issues of web mathematics datasets and proposes a reasonable pipeline.</p> <p>Weaknesses: Missing details on key components and no validation across domains.</p>	2	1	3	3	<p>Soundness: Reasonable approach but lacks detail and cross-domain validation.</p> <p>Overall: Addresses a specific problem but limited generalizability; minor revisions could improve impact.</p>
<p>Strengths: Well-organized and smoothly written manuscript; dataset focuses on complex scenarios.</p> <p>Weaknesses: None.</p>	2.5	3	5	5	<p>Soundness: Thorough and sound research with clear writing.</p> <p>Overall: High-quality, valuable dataset of broad interest—top-notch contribution.</p>
<p>Strengths: K2Q increases diversity and complexity of VQA questions, improving robustness. Covers multiple domains and allows flexible expansion.</p> <p>Weaknesses: Some unnatural questions; domain coverage limitations.</p>	3	2	4	4	<p>Soundness: Technically sound dataset generation method with clear claims; some refinement needed.</p> <p>Overall: Significant contribution to dataset creation; limitations do not detract from value.</p>
<p>Strengths: Multidimensional dataset for summarization evaluation; detailed annotation; systematic LLM comparison.</p> <p>Weaknesses: Missing automatic metrics and slightly imbalanced data sources.</p>	2	1.5	3.5	3.5	<p>Soundness: Comprehensive description and evaluation, but missing key metrics.</p> <p>Overall: Useful dataset with room for improvement in size and balance.</p>
<p>Strengths: Clearly written; effective experiments.</p> <p>Weaknesses: Limited novelty; assumptions not validated; mechanism unclear.</p>	2	1	3	3	<p>Soundness: Clear problem and solid experiments, but lacks novelty and validation depth.</p> <p>Overall: Reasonable contribution needing minor revisions to strengthen claims.</p>
<p>Strengths: Systematic study on temperature parameter’s effect.</p> <p>Weaknesses: Obvious findings; ignores other parameters (e.g., top-p).</p>	1	1	3	3	<p>Soundness: Supports main claims but lacks depth on related parameters.</p> <p>Overall: Reasonable contribution, though limited in novelty and breadth.</p>
<p>Strengths: Investigates stopword and lemmatization effects.</p> <p>Weaknesses: Lacks dataset details and regex explanation.</p>	1	1	3	3	<p>Soundness: Clear claim but insufficient methodological detail.</p> <p>Overall: Modest contribution; needs clarification for better reproducibility.</p>
<p>Strengths: Includes ablation and clear visualizations.</p> <p>Weaknesses: Lacks comparison with post-2022 work and more challenging datasets.</p>	2.5	1	3.5	3	<p>Soundness: Adequate support for claims, but limited scope and depth.</p> <p>Overall: Reasonable contribution, but relevance could be improved with newer benchmarks.</p>

Table 11: Examples of human-assigned and LLM-predicted review scores with large discrepancies. S is for Soundness score, and OA is for Overall Assessment score.