

Who’s Your Judge? On the Detectability of LLM-Generated Judgments

Anonymous ACL submission

Abstract

Large Language Model (LLM)-based judgments leverage powerful LLMs to efficiently evaluate candidate content and provide judgment scores. However, the inherent biases and vulnerabilities of LLM-generated judgments raise concerns, underscoring the urgent need for distinguishing them in sensitive scenarios like academic peer reviewing. In this work, we propose and formalize the task of judgment detection and systematically investigate the detectability of LLM-generated judgments. Unlike LLM-generated text detection, judgment detection relies solely on judgment scores and candidates, reflecting real-world scenarios where textual feedback is often unavailable in the detection process. Our preliminary analysis shows that existing LLM-generated text detection methods perform poorly given their incapability to capture the interaction between judgment scores and candidate content—an aspect crucial for effective judgment detection. Inspired by this, we introduce *J-Detector*, a lightweight and transparent neural detector augmented with explicitly extracted linguistic and LLM-enhanced features to link LLM judges’ biases with candidates’ properties for accurate detection. Experiments across diverse datasets demonstrate the effectiveness of *J-Detector* and show how its interpretability enables quantifying biases in LLM judges. Finally, we analyze key factors affecting the detectability of LLM-generated judgments and validate the practical utility of judgment detection in real-world scenarios.

1 Introduction

Taking advantage of the powerful Large Language Models (LLMs), the paradigm of LLM-based judgment (Zheng et al., 2023; Li et al., 2024) has been proposed, designed to automate and scale up various annotation and reviewing applications (Lee et al.; Zhu et al., 2025; Chang et al., 2025). By combining powerful LLMs with well-designed

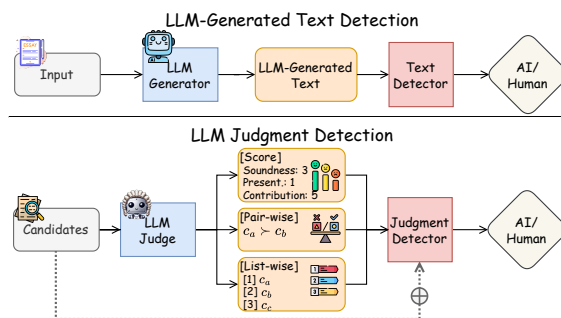


Figure 1: Comparison between LLM-generated judgment detection and text detection.

prompting strategies, LLM-based judgment enables human-like evaluation of long-form and open-ended generation in a more cost-efficient manner. For example, LLM-based judgment has been increasingly used in the peer review of leading AI conferences (Liang et al., 2024).

Despite this remarkable progress, many recent studies point out various biases of LLM-generated judgment toward spurious features, such as length and affinity (Ye et al., 2024; Li et al., 2025a; Zhao et al., 2025a). Besides, the vulnerability of the LLM judgment system has also been revealed, that several maliciously-designed and hard-to-detect tokens or words can fool the LLM judges to give much inconsistent scores despite the candidates’ genuine quality (Shi et al., 2024; Zhao et al., 2025b). Recently, in the scenario of academic peer reviewing, some researchers sneak prompts, which are usually concealed as white text on a white background, into their papers to instruct LLMs to only provide positive feedback and thus trick AI reviewers¹. All these challenges highlight the importance of distinguishing LLM-generated judgments to guarantee the assessment’s fairness and reliability.

To address this concern, we propose the judgment detection task, which aims at examining

¹https://www.theregister.com/2025/07/07/scholars_try_to_fool_llm_reviewers/

071	the detectability of LLM-generated judgments	<i>J-Detector</i> to enable bias quantification in LLM	123
072	across diverse scenarios. Unlike existing machine-	judges. Finally, we analyze key factors affecting	124
073	generated text detection task that focuses on textual	the detectability of LLM-generated judgments and	125
074	content (Mitchell et al., 2023), judgment de-	demonstrate a real-world application that integrates	126
075	tection targets at distinguishing LLM-generated	judgment detection with text-based detection to	127
076	from human-produced judgments solely based on	identify AI-generated reviews in an academic peer	128
077	the <i>candidate content</i> and <i>judgment scores</i> (as il-	reviewing scenario. In summary, our key contribu-	129
078	lustrated in Figure 1). For instance, in academic	tions are:	130
079	paper reviewing, judgment detection will be per-	• We propose, for the first time, the judgment de-	131
080	formed using only the candidate paper and its as-	tection task, which aims at distinguishing human	132
081	signed ratings (e.g., soundness, novelty, overall	and LLM judgments based on judgment scores and	133
082	score), without accessing the full review text. This	candidate content.	134
083	setting is particularly important for real-world sce-	• We design <i>J-Detector</i> , a lightweight and inter-	135
084	narios where textual feedback is often unavailable	pretable detection method, that effectively bridges	136
085	in the detection process. For example, reviewers	candidate and judgment information with linguistic	137
086	who adopt AI-generated reviews may intention-	and LLM-enhanced features.	138
087	ally submit minimal textual content, such as “N/A”	• Through extensive experiments, we demonstrate	139
088	to evade detection. Moreover, in the evaluation	the advantages of <i>J-Detector</i> , identify key factors	140
089	data labeling scenario, annotators are typically re-	driving judgment detectability, and show the utility	141
090	quired to provide only the judgment scores. Score-	of judgment detection in real-world applications.	142
091	based judgment detection is especially critical in		
092	these scenarios to identify the illegal use of LLM-	2 Related Work	143
093	generated judgment and guarantee assessment reli-		
094	ability.	LLM-as-a-judge , first introduced by Zheng et al.	144
095	Developing a good LLM-generated judgment	(2023), leverages powerful LLMs (Zhang et al.,	145
096	detector is not trivial. In our warm-up analy-	2024a,b; Wang et al., 2024a) to automatically eval-	146
097	sis, we identify two key types of information for	uate candidate content and assign scores as judg-	147
098	judgment detection which are not jointly consid-	ment results. This paradigm has been expanded to	148
099	ered in existing related approaches: ❶ Judgment-	diverse applications to judge various types of candi-	149
100	Intrinsic Features , which capture patterns within	dates, including paper quality assessing (Jin et al.,	150
101	the judgment score distribution, and ❷ Judgment-	2024), document relevance measurement (Gao	151
102	Candidate Interaction Features , which capture	et al., 2023; Rahmani et al., 2024), and reasoning	152
103	the interaction between judgment scores and can-	trace correctness verification (Zhang et al.), driving	153
104	didate content. Building on them, we find that	substantial progress in automatic assessment (Li	154
105	existing LLM-generated text detection methods	et al., 2025b; Tan et al., 2025; Beigi et al., 2024;	155
106	fail to capture Judgment-Candidate Interaction Fea-	Hu et al., 2024; Jeong et al., 2024). Despite these	156
107	tures, leading to subpar performance—especially	advances, recent studies highlight notable limita-	157
108	in single-dimension settings, where each judgment	tions. Research has uncovered systematic biases in	158
109	consists of a single score assessing one aspect of	LLM-generated judgments, where evaluations are	159
110	the candidates. To address this, we introduce <i>J-</i>	influenced by spurious features such as response	160
111	<i>Detector</i> , a lightweight and interpretable neural	length or superficial affinity rather than genuine	161
112	detector designed specifically for LLM-generated	content quality (Ye et al., 2024; Li et al., 2025a;	162
113	judgment detection. <i>J-Detector</i> is augmented with	Jiang et al., 2024; Yang et al., 2024). Moreover,	163
114	explicitly extracted linguistic and LLM-enhanced	adversarial work demonstrates that LLM judges	164
115	features to capture systematic correlations between	can be manipulated with a few carefully crafted,	165
116	judgment scores and candidate features that LLM	hard-to-detect tokens or phrases, which induce dis-	166
117	judges are often biased toward, thereby effectively	proportionately high scores misaligned with ac-	167
118	leveraging these biases for more accurate detection.	tual candidate quality (Shi et al., 2024; Zhao et al.,	168
119	Experiments across diverse judgment datasets	2025b). To mitigate these issues, methods such	169
120	demonstrate the effectiveness of <i>J-Detector</i> and	as bias quantification (Ye et al., 2024) and human-	170
121	the two types of augmented features. Besides, we	in-the-loop calibration (Wang et al., 2023a) have	171
122	showcase how to leverage the interpretability of	been proposed. Building on this line of research,	172

we introduce a new task, judgment detection, that aims to distinguish and prevent the misuse of LLM-generated judgments.

AI-generated Text Detection aims to distinguish machine-generated from human-produced text, evolving from early stylometric and perplexity-based methods (Gehrmann et al., 2019; Zellers et al., 2019) to supervised classifiers (Ippolito et al., 2020; Mitchell et al., 2023), and more recently toward robust, generalizable approaches such as zero-shot prompting and watermarking (Sun and Lv, 2025; Mao et al., 2025). Another relevant line of work for us is the detection of LLM-generated peer reviews (Tao et al.; Yu et al.; Rao et al., 2025), where detectors are designed to distinguish machine-written reviews from human-authored ones. However, these approaches rely on textual review content, which is often unavailable in broader judgment settings. In this work, we borrow insights from both fields and propose judgment detection to explore the detectability of LLM-produced judgment, using judgment scores without accessing textual feedback.

3 Task Statement

A *judgment* refers to an assessment made over one or more candidates $c \in \mathcal{C}$, where $|\mathcal{C}|$ denotes the size of the candidate set. A judgment score is denoted by $j = (j_1, \dots, j_d) \in \mathcal{Y}^d$. It can be either *single-dimensional* ($d = 1$), reflecting an assessment toward a single aspect, or *multi-dimensional* ($d > 1$), where each component J_i corresponds to a distinct evaluation aspect (e.g., relevance, fluency, coherence). With these definitions, we formulate the task as follows:

Definition 3.1 (Judgment Detection). LLM-generated judgment detection is defined over *judgment groups*. A judgment group is given by $G = \{(c^i, j^i)\}_{i=1}^k$, where each candidate $c^i \in \mathcal{C}$ is paired with a judgment score $j^i \in \mathcal{J}$. The task is to classify whether a group G originates from a human judge or from an LLM. Formally, the label space is $L = \{0, 1\}$, where $\ell = 0$ denotes human-produced judgments and $\ell = 1$ denotes LLM-generated judgments. The goal is to learn a function $f_\theta : G \rightarrow [0, 1]$, where $f_\theta(G)$ outputs the probability that G was generated by an LLM. The final prediction is obtained as $\hat{y} = \mathbb{I}[f_\theta(G) \geq \tau]$, with threshold $\tau \in [0, 1]$ and indicator function $\mathbb{I}[\cdot]$.

When the group size is 1, *i.e.*, $|G| = 1$, the task is degraded to an i.i.d. (instance-level) detection

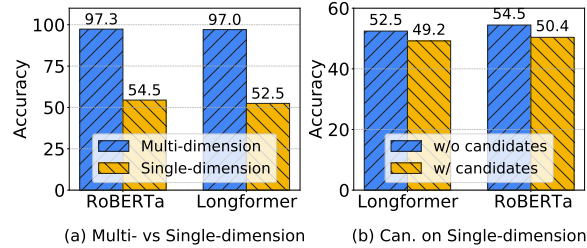


Figure 2: Multi- vs Single-dimension and Candidate Effect.

setting, where each judgment is treated independently. When $|G| > 1$, the group setting better reflects real practice, since judgments are usually produced in batches (e.g., a reviewer scores multiple papers or an annotator evaluates a set of model outputs), and collective patterns across the group can reveal whether the judgments are human-produced or LLM-generated.

4 Warm-up Analysis: What Matters for LLM-generated Judgment Detection?

To understand the key ingredients of a reliable judgment detector, we first conduct a warm-up study by adapting LLM-generated text detection methods to the judgment detection setting. Specifically, we employ small language models (SLM)-based detectors (Wu et al., 2024), *RoBERTa* and *Longformer*, as f_θ and evaluate them on four datasets: *Helpsteer2*, *Helpsteer3*, *NeurIPS*, and *ANTIQUe*. More information about implementation and dataset can be found in Section 6.1.

Multi-dimension vs Single-dimension performance. As shown in Figure 2 (a), both *RoBERTa* and *Longformer* achieve high accuracy in the *multi-dimension* scenarios (*Helpsteer2* and *NeurIPS*) but perform poorly in the *single-dimension* scenarios (*Helpsteer3* and *ANTIQUe*). We assume that this discrepancy arises because, in multi-dimension settings, the detectors can exploit distributional differences in how humans and LLMs assign scores across multiple judgment dimensions, whereas in single-dimensional settings, such distributional cues are almost absent.

Adding candidate information. We further extend the single-dimension setting by providing candidate texts alongside their judgments, exploring whether the detectors can extract and leverage judgment-candidate interaction information. As shown in Figure 2 (b), however, adding candidates does not lead to any performance improvement. This suggests that SLM-based detectors are unable to

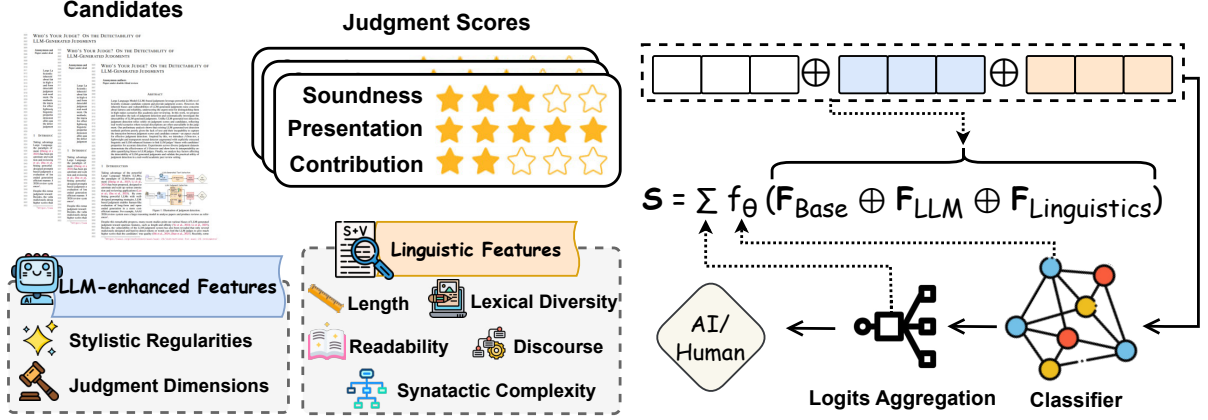


Figure 3: The overview pipeline of our *J-Detector* for LLM-generated judgment detection.

263 directly capture and utilize the interaction between
 264 judgments and candidate content from raw input.
 265 **Takeaway.** From this warm-up study, we identify
 266 two complementary types of information that
 267 a reliable judgment detector should exploit: ❶
 268 **Judgment-Intrinsic Features**, revealed by the
 269 large performance gap between multi-dimension
 270 and single-dimension settings, indicating that dis-
 271 tributional patterns within judgment scores them-
 272 selves can be highly informative; and ❷ **Judgment-**
 273 **Candidate Interaction Features**, which capture
 274 how judgment scores relate to the underlying can-
 275 didate content but remain largely unexplored by
 276 existing methods. These findings highlight that ex-
 277 isting SLM-based text detection methods mainly
 278 leverage judgment-intrinsic patterns but fail to cap-
 279 ture judgment–candidate interactions, which are
 280 especially critical in single-dimension scenarios.

281 5 *J-Detector*: A Lightweight and 282 interpretable Detector

283 To address the limitation of existing text detectors
 284 and design an effective and robust approach for
 285 LLM-generated judgment detection, we first identify
 286 three criteria that a good LLM-generated judg-
 287 ment detector should embody:

- 288 • **(Accurate)** The detection method should be able
 289 to leverage both Judgment-Intrinsic Features and
 290 Judgment-Candidate Interaction Features to deliver
 291 reliable detection results in various scenarios.
- 292 • **(Efficient)** Both the training and inference of the
 293 detector should incur minimal computational over-
 294 head, enabling the method to be deployed in large-
 295 scale judgment detection scenarios.
- 296 • **(Interpretable)** The detection method should be in-
 297 terpretable to support bias analysis in LLM judges.

298 Following these principles, we design *J-*
 299 *Detector*, an accurate, lightweight and interpretable
 300 detector involving the following components. The
 301 overview pipeline is presented in Figure 3.
 302 **Feature Augmentation.** Let \mathbf{F} denote the instance-
 303 level feature vector used by *J-Detector*. We con-
 304 struct it by concatenating three types of features
 305 together:

$$\mathbf{F} = \mathbf{F}_{\text{base}} \oplus \mathbf{F}_{\text{LLM}} \oplus \mathbf{F}_{\text{linguistic}}, \quad (1) \quad 306$$

307 where \mathbf{F}_{base} contains the *given judgment scores*.
 308 \mathbf{F}_{LLM} and $\mathbf{F}_{\text{linguistic}}$ are *LLM-enhanced features*
 309 and *linguistic features* we extract from candidates
 310 content, which act as distilled information of can-
 311 didates and are leveraged to link judgment scores
 312 with candidates’ content.

313 **LLM-enhanced Features.** Borrowing insights
 314 from LLM-based text detection methods (Bao et al.,
 315 2024), we propose LLM-enhanced features to pro-
 316 duce the following types of features:

- 317 • **Stylistic regularities:** scores reflecting surface pol-
 318 ish and presentation patterns of the candidates, in-
 319 cluding *style*, *wording*, and *format*. These aim to
 320 capture the spurious preference LLM judges tend
 321 to have over superficial attributes (Li et al., 2025a).
- 322 • **Judgment-aligned dimensions:** scores aligned to
 323 the same dimensions used in the given judgment
 324 scores. These aim to enhance features by leverag-
 325 ing the similarity of biases across LLM judges.

326 By injecting these high-level, bias-informed sig-
 327 nals, LLM-enhanced Features enable the detector
 328 to better capture subtle judgment patterns that are
 329 difficult to learn from raw candidate content alone.

330 **Linguistic Features.** We further introduce lin-
 331 guistic features $\mathbf{F}_{\text{linguistic}}$ to capture low-level lin-
 332 guistic regularities that often correlate with sys-

333 thematic biases of LLM judges. Specifically, we
334 extract the following aggregated features from the
335 candidate content:

- 336 • *Length*: total token and character counts, as well
337 as average sentence length, to capture the *length*
338 *bias* where LLM judges favor lengthy content and
339 responses (Wei et al.).
- 340 • *Lexical diversity*: unique-token ratio and average
341 word length, which reflect the *surface beauty bias*
342 of LLM-generated judgments compared to human-
343 produced ones (Chen et al., 2024).
- 344 • *Readability*: a composite readability index (e.g.,
345 Coleman–Liau), measuring the *fluency bias* where
346 LLMs tend to favor superficially fluent texts, disre-
347 garding their true quality (Wu and Aji, 2025).
- 348 • *Syntactic complexity*: dependency tree depth and
349 average dependency distance, used to identify the
350 *complexity bias* often observed in LLM judges (Ye
351 et al., 2024).
- 352 • *Discourse/hedging*: the frequency of discourse
353 markers and hedging expressions, capturing the
354 *presentation bias* of LLM, which prefer content
355 with confident tones (Kharchenko et al., 2025).

356 These features provide a compact yet informa-
357 tive summary of linguistic cues, enabling the detec-
358 tor to exploit stable and interpretable signals that
359 are complementary to LLM-enhanced features.

360 **Model Training.** Given labeled instances (\mathbf{F}, y) ,
361 we train a lightweight binary classifier f_θ (e.g.,
362 RandomForest (Breiman, 2001)) to output a *logit*
363 $z \in \mathbb{R}$ indicating the likelihood that the judgment
364 was generated by an LLM ($y = 1$) or by a human
365 ($y = 0$). The classifier is trained using the aug-
366 mented feature \mathbf{F} and serves as the instance-level
367 building block for group-level decisions.

368 **Group-level Aggregation.** To enable the group-
369 level detection setting, we propose a simple ag-
370 gregation method to produce the group-level label
371 give each single prediction. Given a group G con-
372 sisting of k judgments with instance-level logits
373 $\{\hat{z}_1, \dots, \hat{z}_k\}$, we aggregate the evidence using sum
374 aggregation: $\text{score}(G) = \sum_{i=1}^k \hat{z}_i$.

375 In summary, *J-Detector* is designed to satisfy the
376 three criteria identified at the beginning of this sec-
377 tion. First, by incorporating both LLM-enhanced
378 and linguistic features, it is able to capture not
379 only Judgment-Intrinsic Features but also criti-
380 cal Judgment–Candidate Interaction Features, en-
381 abling accurate detection across single-dimensional
382 and multi-dimensional scenarios. Second, it builds
383 on a lightweight binary classifier, making both
384 training and inference highly efficient and thus

385 suitable for large-scale deployment. Third, since
386 the features are semantically clear and the classi-
387 fier itself is simple, the framework offers strong
388 interpretability, which can be leveraged to system-
389 atically quantify and analyze the biases of LLM
390 judges.

391 6 Main Experiment

392 6.1 Experiment Settings

393 **Datasets.** We build a comprehensive LLM-
394 generated judgment detection dataset, *JD-Bench*,
395 which integrates four representative datasets cover-
396 ing three judgment types: pointwise, pairwise and
397 listwise (Li et al., 2024). Among them, *HelpSteer2*
398 provides large-scale pointwise human ratings of
399 LLM responses for helpfulness evaluation, while
400 *HelpSteer3* extends this with pairwise human pref-
401 erence comparisons. The *NeurIPS Review dataset*
402 offers expert peer reviews with multi-dimensional
403 scores such as soundness and novelty, representing
404 high-stakes evaluation. Finally, *ANTIQUÉ* sup-
405 plies listwise human judgments for ranking doc-
406 uments in non-factoid question answering. All
407 four datasets contain human-labeled judgments as
408 reliable references, and we further collect LLM-
409 generated judgments from a diverse pool of mod-
410 els. In total, *JD-Bench* covers a wide spectrum
411 of model families, including *OpenAI*, *Anthropic*,
412 and *Google* for closed-source models, and *LLaMA*,
413 *Qwen*, *Mistral*, and *DeepSeek* for open-source mod-
414 els, ensuring diversity in judgment patterns.

415 **Compared Methods.** In our main experiment, we
416 compare our proposed *J-Detector* against a series
417 of baseline methods, all of which are listed as fol-
418 lows:

- 419 • **SLM-based Detector.** In line with SLM-based
420 text detectors (Yu et al., 2025), this approach feeds
421 either the judgment scores alone or the judgment
422 scores together with the candidate content (w/ can-
423 didates) to train a small language model-based clas-
424 sifier to predict whether the judgment was produced
425 by a human or from an LLM.
- 426 • **LLM-as-a-judge-detector.** Inspired by logits-
427 based detection in AI-generated text detec-
428 tion (Mitchell et al., 2023), where a surrogate LLM
429 is used to compute likelihoods, we adopt a single
430 LLM that first generates judgment scores and then
431 compares them with the judgment scores to be de-
432 tected, making the detection decision based on their
433 similarity.

Method	Helpsteer2		Helpsteer3		NeurIPS		ANTIQUÉ		AVG	
	F1	AUROC	F1	AUROC	F1	AUROC	F1	AUROC	F1	AUROC
<i>SLM-based methods</i>										
RoBERTa	98.1	99.6	50.9	64.5	96.2	99.4	30.0	56.8	68.8	80.1
RoBERTa w/ candidates	98.1	99.6	50.0	63.4	96.3	99.3	27.6	56.6	68.0	79.7
Longformer	98.1	99.7	54.5	65.7	96.2	99.5	30.6	56.6	69.9	80.4
Longformer w/ candidates	98.1	99.7	51.4	64.3	96.2	99.4	21.8	48.8	66.9	78.0
<i>LLM-based methods</i>										
LLM	51.5	50.3	50.3	50.1	43.9	50.2	49.6	49.9	48.8	50.1
LLM w/ Sample-level	49.8	49.7	49.6	50.2	50.5	50.4	50.9	50.3	50.2	50.2
LLM w/ Distribution-level	52.1	50.0	48.8	50.3	49.6	49.8	50.7	50.1	50.3	50.1
LLM w/ Sample-level + Distribution-level	58.7	50.4	49.4	49.6	51.2	50.2	50.2	49.9	52.4	50.0
<i>J-Detector (ours)</i>										
LGBM	99.6	100.0	68.1	73.3	98.7	99.9	85.4	93.3	88.0	91.6
RandomForest	99.5	100.0	74.0	77.0	97.0	99.7	82.6	90.6	88.3	91.8
XGB	99.8	100.0	68.5	73.6	98.4	99.8	84.2	92.3	87.7	91.4

Table 1: Main experimental results on *JD-Bench*. We report F1 and AUROC scores, with the best results highlighted in bold. Each experiment is repeated five times, and average scores are reported.

- Sample-level LLM-based Analysis.** Inspired by recent agent-based frameworks that maintain guideline banks for distinguishing human and AI text (Li et al., 2025c), we let the LLM analyze Human-LLM judgment-candidate pairs to extract concise instance-level features (e.g., length bias in LLM judgments), which are stored in a feature bank to capture regularities useful for detection.
- Distribution-level LLM-based Analysis.** Drawing inspiration from recent work that guides LLMs in structured extraction and analysis of visual summaries (Liu et al., 2025), we provide the model with dataset-level summaries (e.g., per-label histograms and correlations), enabling it to incorporate global and distributional cues into the detection decision.

Implementation Details. We implement our *J-Detector* using three models from the Scikit-learn library (Pedregosa et al., 2011): LGBM (Ke et al., 2017), RandomForest (Breiman, 2001), and XGB (Chen and Guestrin, 2016). We employ *Qwen-3-8B* for both feature augmentation and as the backbone for LLM-based baselines. For SLM-based methods, we use *RoBERTa-base* and *Longformer-4096*. For SLM training, we use a batch size of 8 and fine-tune the SLM for 3 epochs on each dataset. In the main experiments, the group size is fixed to $k=4$. More details, including the *JD-Bench* construction, design of baseline methods, and implementation specifics are provided in Appendix C.

6.2 Main Result

SLM-based Methods Analysis. As we discussed in Section 4, SLM-based methods perform strongly on multi-dimensional datasets like Helpsteer2 (98.1% F1 on RoBERTa) and NeurIPS

(96.2% on RoBERTa), but drop to around 50–55% F1 on single-dimensional datasets like Helpsteer3 and Antique. Even adding candidates barely helps. This shows SLMs rely on inter-dimension patterns and fail to link judgments with candidates when such distributional cues are absent.

LLM-based Methods Analysis. Furthermore, all LLM-based methods hover near 50% F1 score across datasets, indicating almost random guessing. When combining with sample-level comparative analysis and distribution-level chart reasoning, LLM-based detection methods yield some gains in multi-dimensional datasets (e.g., from 51.5% to 58.7% F1 score). While this improvement doesn’t appear in Helpsteer3 and ANTIQUÉ, we conclude that LLM-based detectors also suffer from leveraging judgments-candidates interaction, with either sample- or distribution-level methods.

J-Detector Analysis. Compared with them, *J-Detector* achieves the best detection performance across all 4 datasets and 2 metrics, far surpassing all baselines. Noted that in the single-dimensional judgment scenarios, *J-Detector* yields much better detection performance compared with other baselines. This demonstrates that explicitly modeling the distributional patterns and biases of LLM judgments is crucial for accurate detection, enabling robust performance in both single-dimensional and multi-dimensional judgment detection scenarios.

Ablation Study. Figure 5 shows that both LLM-enhanced and linguistic features consistently improve performance across all group sizes. Removing either feature causes the F1 score to drop at every group size—for example, at $k = 16$, removing linguistic features lowers F1 by 5.3%, and removing both leads to a 12.3% drop. This demonstrates

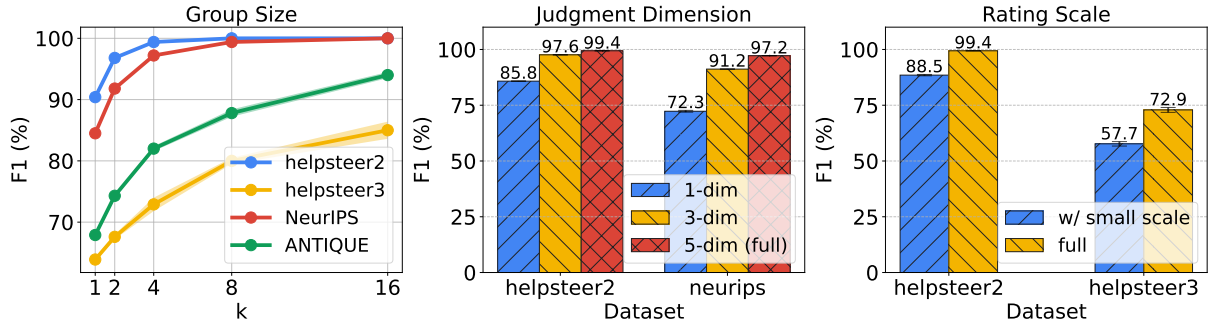


Figure 4: Detectability analysis on group size, judgment dimensions and rating scale.

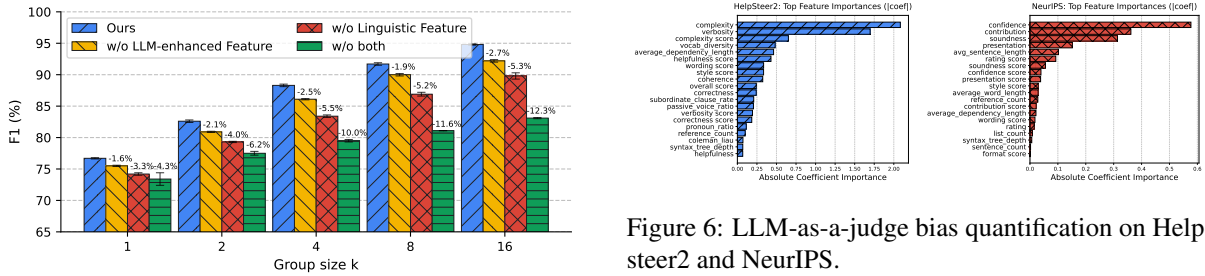


Figure 5: Ablation study on LLM-enhanced and linguistic features.

Figure 6: LLM-as-a-judge bias quantification on Helpsteer2 and NeurIPS.

505 that the two augmented features are complementary
506 and beneficial across all datasets and group-size set-
507 tings.

508 **Bias Quantification with J-Detector.** Addition-
509 ally, we illustrate how the transparency and inter-
510 pretability of *J-Detector* can be leveraged to
511 quantify biases in LLM-as-a-judge by analyzing
512 which features most strongly influence the detec-
513 tor’s decisions. Specifically, we select the top
514 20 most important features ranked by their ab-
515 solute coefficient values, and report the results
516 on the Helpsteer2 and NeurIPS datasets in Fig-
517 ure 6. The analysis reveals that base judgment
518 score features provide strong signals for distin-
519 guishing LLM-generated judgments from human-
520 produced ones, highlighting the critical role of
521 *Judgment-Intrinsic Features*. As shown in the
522 figure, LLM judges exhibit the strongest bias in
523 the *complexity* and *confidence* dimensions for the
524 two datasets, respectively, consistent with prior
525 findings that LLMs tend to favor more complex
526 responses (Ye et al., 2024; Yang et al., 2024)
527 and often display overconfidence (Kadavath et al.,
528 2022). In addition, we observe common cross-
529 dataset biases such as *length bias* (captured by
530 average_dependency_length) and *beauty bias*
531 (reflected in style-related scores), which echo

broader concerns about spurious preference and
532 correlations in LLM-based judgments (Wang et al.,
533 2023b; Shi et al., 2024).
534

7 Further Analysis 535

In this section, we empirically analyze the key fac-
536 tors that influence the detectability of the LLM-
537 generated judgment, as well as present a real-world
538 application to combine LLM-based judgment de-
539 tection with text detection in real-world academic
540 peer reviewing scenarios.
541

7.1 Detectability Analysis 542

**Detectability analysis across group size, judg-
543 ment dimensions, and rating scale.** Figure 4
544 shows that group size is a key factor in the de-
545 tectability of LLM-generated judgments: the F1
546 score consistently improves as the group size in-
547 creases across all four datasets (e.g., F1 score in
548 Helpsteer3 rises from 63.9% at $k = 1$ to 85.0%
549 at $k = 16$). The number of judgment dimen-
550 sions also plays an important role; for instance,
551 when only a single dimension out of the five is used
552 in the NeurIPS dataset, the F1 score drops sub-
553 stantially (from 97.2% to 72.3%). This confirms that
554 multi-dimensional judgments provide richer distri-
555 butional signals as Judgment-Intrinsic Features for
556 detection.
557

In addition, the granularity of the rating scale
558 further impacts detectability: collapsing to a coarse
559

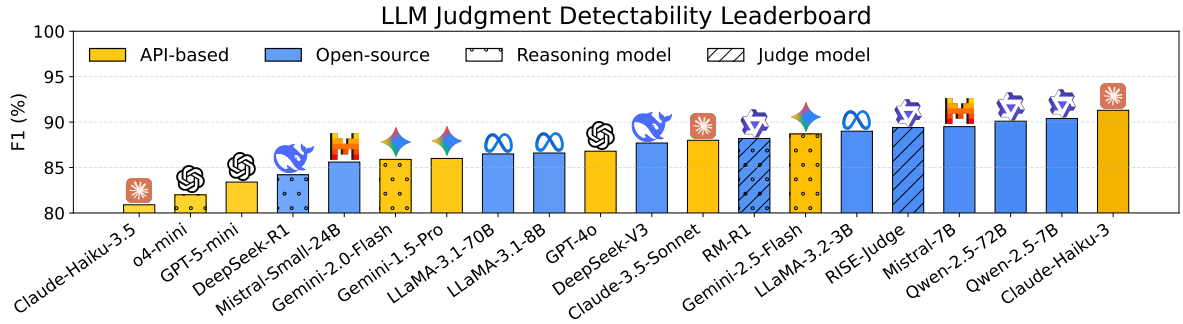


Figure 7: Detectability leaderboard on 20 LLMs. RM-R1 and RISE-Judge are based on Qwen-2.5-7B.

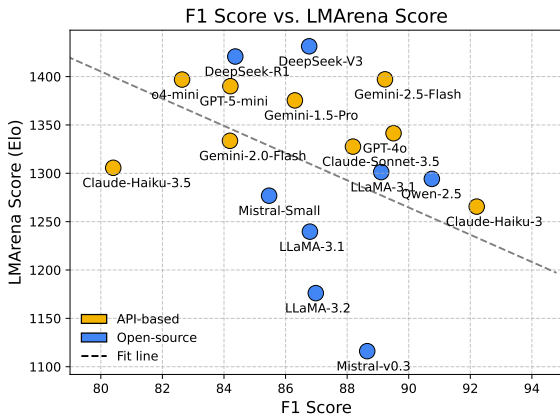


Figure 8: Correlation between judge LLMs' detectability and LMArena score.

scale (e.g., merging $-3/ -2/ -1$ into -1 and $1/2/3$ into 1 in Helpsteer3) leads to degraded performance (e.g., F1 drops from 72.9% to 57.7%). Overall, these results underscore that group size, the number of dimensions, and the rating scale collectively shape how detectable LLM-generated judgments are.

Detectability of Various LLM Judges. Additionally, Figure 7 summarizes the detectability leaderboard across 20 LLMs, averaged over different group sizes. We observe that API-based models (yellow bars) are generally more difficult to detect than open-source models (blue bars), indicating that closed commercial systems such as GPT-5-mini and Claude-Haiku-3 produce judgments that more closely resemble human annotations.

Within the same model families, larger models tend to be less detectable than smaller ones: for instance, among LLaMA-3 and Qwen-2.5 families, larger models consistently achieve lower detectability. Moreover, reasoning models (dotted bars) and specialized judge models (striped bars) consistently achieve higher robustness than standard LLMs, sug-

gesting that models explicitly optimized for reasoning or evaluation align more closely with human judgment distributions and are therefore harder to distinguish from human judges.

As presented in Figure 8, we also study the correlation between the detectability of different LLM judges and their LMArena score (Chiang et al., 2024), which is a proxy of LLMs' alignment degree with human preference and value. We find a clear negative correlation: models with higher alignment scores are systematically less detectable. This observation reinforces our previous findings, supporting the hypothesis that as models become better aligned with human values, the gap between their judgments and human annotations narrows, making their outputs increasingly difficult to distinguish from those of human judges.

For LLM-generated judgment detectability, we also theoretically prove and demonstrate each influence factor's effect and put it in Appendix D.

8 Conclusion

In this work we introduced judgment detection as the task of distinguishing human from LLM-generated judgments and proposed *J-Detector*, a lightweight, interpretable detector enhanced with linguistic and LLM-based features. Experiments on *JD-Bench* show that *J-Detector* consistently outperforms baselines, while our theoretical and empirical analyses reveal that detectability improves with larger group size, richer dimensions, finer rating scales, and greater human-LLM divergence. Using *J-Detector*'s transparency, we further quantified systematic biases in LLM judges, such as complexity, confidence, and length biases, and demonstrated practical value in peer-review authenticity checking. These findings establish LLM-generated judgment detection as a key safeguard for ensuring fairness and accountability in LLM-as-a-judge.

621 Limitations

622 Despite the effectiveness of *J-Detector*, our work
623 has several limitations: First, as the first work to
624 formally propose and define the *judgment detection*
625 task, our primary goal was to establish the task’s
626 feasibility and identify core features (Judgment-
627 Intrinsic and Candidate-Interaction). Consequently,
628 the proposed J-Detector architecture is intention-
629 ally kept lightweight and straightforward. While
630 this ensures transparency and efficiency, we ac-
631 knowledge that more sophisticated neural architec-
632 tures, such as graph-based models or large-scale
633 multi-modal pre-training, might yield even higher
634 detection accuracy. We leave the exploration of
635 these complex methods for future work. Second,
636 our current study primarily focuses on "full" LLM-
637 generated judgments. In real-world scenarios, a
638 "human-in-the-loop" or "AI-assisted" setting is also
639 common, where a human judge might refine or pol-
640 ish an initial score generated by an LLM. Distin-
641 guishing such hybrid judgments from purely hu-
642 man or purely AI ones remains a significant chal-
643 lenge. Future research should investigate the detec-
644 tion of varying degrees of AI involvement in the
645 judgment process.

646 References

647 Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi
648 Yang, and Yue Zhang. 2024. Fast-detectgpt: Efficient
649 zero-shot detection of machine-generated text via
650 conditional probability curvature. In *ICLR*.

651 Alimohammad Beigi, Zhen Tan, Nivedh Mudiam,
652 Canyu Chen, Kai Shu, and Huan Liu. 2024. Model
653 attribution in llm-generated disinformation: A do-
654 main generalization approach with supervised con-
655 trastive learning. In *2024 IEEE 11th International
656 Conference on Data Science and Advanced Analytics
657 (DSAA)*, pages 1–10. IEEE.

658 Leo Breiman. 2001. Random forests. *Machine learning*,
659 45(1):5–32.

660 Yuan Chang, Ziyue Li, Hengyuan Zhang, Yuanbo Kong,
661 Yanru Wu, Zhijiang Guo, and Ngai Wong. 2025.
662 Treereview: A dynamic tree of questions framework
663 for deep and efficient llm-based scientific peer review.
664 *arXiv preprint arXiv:2506.07642*.

665 Guiming Chen, Shunian Chen, Ziche Liu, Feng Jiang,
666 and Benyou Wang. 2024. Humans or llms as the
667 judge? a study on judgement bias. In *Proceedings
668 of the 2024 Conference on Empirical Methods in
669 Natural Language Processing*, pages 8301–8327.

670 Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A
671 scalable tree boosting system. In *Proceedings of*

*the 22nd ACM SIGKDD International Conference
on Knowledge Discovery and Data Mining*, pages
785–794. ACM.

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-
sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,
Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E
Gonzalez, and 1 others. 2024. Chatbot arena: An
open platform for evaluating llms by human prefer-
ence. In *Forty-first International Conference on
Machine Learning*.

Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Ship-
ing Yang, and Xiaojun Wan. 2023. Human-like sum-
marization evaluation with chatgpt. *arXiv preprint
arXiv:2304.02554*.

Sebastian Gehrmann, Hendrik Strobelt, and Alexan-
der M Rush. 2019. **Gltr: Statistical detection and
visualization of generated text**. In *Proceedings of the
57th Annual Meeting of the Association for Compu-
tational Linguistics: System Demonstrations*, pages
111–116. Association for Computational Linguistics.

Helia Hashemi, Mohammad Aliannejadi, Hamed Za-
mani, and W Bruce Croft. 2020. Antique: A non-
factoid question answering benchmark. In *European
Conference on Information Retrieval*, pages 166–173.
Springer.

Lijie Hu, Liang Liu, Shu Yang, Xin Chen, Zhen Tan,
Muhammad Asif Ali, Mengdi Li, and Di Wang.
2024. Understanding reasoning in chain-of-
thought from the hopfieldian view. *arXiv preprint
arXiv:2410.03595*.

Lijie Hu, Chenyang Ren, Zhengyu Hu, Hongbin Lin,
Cheng-Long Wang, Zhen Tan, Weimin Lyu, Jingfeng
Zhang, Hui Xiong, and Di Wang. Editable concept
bottleneck models. In *Forty-second International
Conference on Machine Learning*.

Daphne Ippolito, Daniel Duckworth, Chris Callison-
Burch, and David Eck. 2020. **Automatic detection of
generated text is easiest when humans are fooled**. In
*Proceedings of the 58th Annual Meeting of the Asso-
ciation for Computational Linguistics*, pages 1808–
1822. Association for Computational Linguistics.

Ujun Jeong, Bohan Jiang, Zhen Tan, Russell Bernard,
and Huan Liu. 2024. Bluetempnet: A temporal multi-
network dataset of social interactions in bluesky so-
cial. *IEEE Data Descriptions*.

Bohan Jiang, Dawei Li, Zhen Tan, Xinyi Zhou, Ashwin
Rao, Kristina Lerman, H Russell Bernard, and Huan
Liu. 2024. Assessing the impact of conspiracy the-
ories using large language models. *arXiv preprint
arXiv:2412.07019*.

Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kai-
jie Zhu, Yijia Xiao, and Jindong Wang. 2024. Agen-
t-review: Exploring peer review dynamics with llm
agents. In *EMNLP*.

839	Xuezhi Wang, Jason Wei, Denny Zhou, Ed Chi, Quoc Le, and Dale Schuurmans. 2023b. Adversarial attacks reveal spurious correlations in large language model evaluations. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)</i> .	Hengyuan Zhang, Chenming Shang, Sizhe Wang, Dongdong Zhang, Feng Yao, Renliang Sun, Yiyao Yu, Yujiu Yang, and Furu Wei. 2024a. Shifcon: Enhancing non-dominant language capabilities with a shift-based contrastive framework. <i>arXiv preprint arXiv:2410.19453</i> .	896 897 898 899 900 901
845	Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024b. Helpsteer 2: Open-source dataset for training top-performing reward models. <i>Advances in Neural Information Processing Systems</i> , 37:1474–1501.	Hengyuan Zhang, Yanru Wu, Dawei Li, Sak Yang, Rui Zhao, Yong Jiang, and Fei Tan. 2024b. Balancing speciality and versatility: a coarse to fine framework for supervised fine-tuning large language model. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 7467–7509.	902 903 904 905 906 907
851	Zhilin Wang, Jiaqi Zeng, Olivier Delalleau, Hoo-Chang Shin, Felipe Soares, Alexander Bukharin, Ellie Evans, Yi Dong, and Oleksii Kuchaiev. 2025. Helpsteer3-preference: Open human-annotated preference data across diverse tasks and languages. <i>arXiv preprint arXiv:2505.11475</i> .	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. In <i>The Thirteenth International Conference on Learning Representations</i> .	908 909 910 911 912
857	Hui Wei, Shenghua He, Tian Xia, Fei Liu, Andy Wong, Jingyang Lin, and Mei Han. Systematic evaluation of llm-as-a-judge in llm alignment tasks: Explainable metrics and diverse prompt templates. In <i>ICLR 2025 Workshop on Building Trust in Language Models and Applications</i> .	Chengshuai Zhao, Zhen Tan, Pingchuan Ma, Dawei Li, Bohan Jiang, Yancheng Wang, Yingzhen Yang, and Huan Liu. 2025a. Is chain-of-thought reasoning of llms a mirage? a data distribution lens. <i>arXiv preprint arXiv:2508.01191</i> .	913 914 915 916 917
863	Junchao Wu, Runzhe Zhan, Derek Wong, Shu Yang, Xinyi Yang, Yulin Yuan, and Lidia Chao. 2024. Detectrl: Benchmarking llm-generated text detection in real-world scenarios. <i>Advances in Neural Information Processing Systems</i> , 37:100369–100401.	Yulai Zhao, Haolin Liu, Dian Yu, SY Kung, Haitao Mi, and Dong Yu. 2025b. One token to fool llm-as-a-judge. <i>arXiv preprint arXiv:2507.08794</i> .	918 919 920
868	Minghao Wu and Alham Fikri Aji. 2025. Style over substance: Evaluation biases for large language models. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 297–312.	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in Neural Information Processing Systems</i> , 36:46595–46623.	921 922 923 924 925 926
872	Shiping Yang, Jie Wu, Wenbiao Ding, Ning Wu, Shining Liang, Ming Gong, Hengyuan Zhang, and Dongmei Zhang. 2024. Quantifying the robustness of retrieval-augmented language models against spurious features in grounding data. <i>arXiv preprint arXiv:2503.05587</i> .	Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. 2025. Deepreview: Improving llm-based paper review with human-like deep thinking process. <i>arXiv preprint arXiv:2503.08569</i> .	927 928 929 930
877	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, and 1 others. 2024. Justice or prejudice? quantifying biases in llm-as-a-judge. <i>arXiv preprint arXiv:2410.02736</i> .	A The Use of LLMs for Writing	931
882	Sungduk Yu, Man Luo, Avinash Madasu, Vasudev Lal, and Phillip Howard. Is your paper being reviewed by an llm? investigating ai text detectability in peer review. In <i>Neurips Safe Generative AI Workshop 2024</i> .	We employed Google’s Gemini 2.5 Pro and OpenAI’s GPT-5 as writing assistance tools during the preparation of this manuscript. Their role was exclusively for language refinement, such as improving readability and rephrasing for clarity in an academic writing style. This usage aligns with standard academic practices for language polishing.	932 933 934 935 936 937 938
887	Sungduk Yu, Man Luo, Avinash Madusu, Vasudev Lal, and Phillip Howard. 2025. Is your paper being reviewed by an llm? benchmarking ai text detection in peer review. <i>arXiv preprint arXiv:2502.19614</i> .	B Additional Analysis	939
891	Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news . In <i>Advances in Neural Information Processing Systems</i> , volume 32.	B.1 Judgment-Text Co-Detection: An Application	940 941
892		In this section, we explore two real-world scenarios where LLM-generated judgment detection can support peer review authenticity checking. First, the few-shot detection setting simulates cases where a new conference is launched or the review form has	942 943 944 945 946

Method	Few-shot	Missing review
w/ RoBERTa-text	67.2	90.5
w/ <i>J-Detector</i>	64.4	86.2
w/ RoBERTa-text & <i>J-Detector</i>	74.6	99.3

Table 2: An application to leverage judgment and text feedback for AI-generated review detection in few-shot and missing review scenarios.

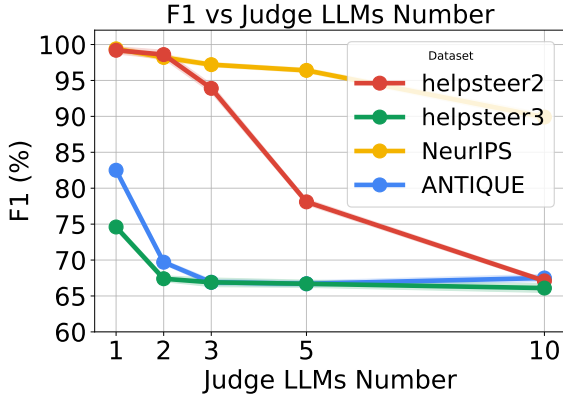


Figure 9: Detectability of LLM-generated judgment in multiple LLM judges setting.

changed. Here, we set the number of training samples to be 60. Second, the missing-text detection setting addresses the common case where reviews lack enough textual feedback. We simulate this setting by masking 15% of the text reviews.

The results in Table 2 show that combining the *J-Detector* with a text-based detector (RoBERTa-text) achieves the best performance in both settings (74.6% vs. 67.2% in few-shot, and 99.3% vs. 90.5% in missing-text), outperforming either method alone. This demonstrates that LLM-generated judgment detection provides complementary signals to text-based detectors and is highly valuable in real-world low-resource or judgment score-only scenarios for robust and reliable detection.

B.2 Judgment Detection with Multiple LLM Judges

In this section, we examine how the detectability of LLM-generated judgments changes when multiple LLM judges are involved. This setting reflects real-world scenarios where judgments may come from a diverse pool of LLMs. As shown in Figure 9, we randomly sample 2, 3, 5, or 10 LLMs from our *JD-Bench* and mix their judgments in both the training and testing sets. We observe a substantial drop in detection performance across all four datasets (e.g., the F1 score decreases from

99.8% to 66.9% on Helpsteer2). This suggests that detecting LLM-generated judgments becomes significantly more challenging when multiple LLM judges are present, as detectors must learn to recognize distinct patterns from different models. Notably, the performance drop is relatively small on the NeurIPS dataset, indicating stronger shared biases among LLM judges in that domain. One promising direction for future work is to explore effective LLM-generated judgment detection methods under multiple judges’ settings.

C Experiment Implementation Details

C.1 General Implementation Details

C.1.1 LLM-enhanced Feature Extraction

To extract the LLM-enhanced features (F_{LLM}) mentioned in Section 5, we utilize Qwen-2.5-7B-Instruct as the backbone extractor. For each candidate c , we use the following prompt template to obtain stylistic and judgment-aligned scores:

Prompt Template for Feature Extraction

```
You are an expert annotator. Please
evaluate the following content based on
specific attributes.
Content: {candidate_text}
Task: Assign a score from 1 to 5 for
the following dimensions: [Style, Wording,
Format, Soundness, Novelty].
Output format: {"style": x, "wording": x,
"format": x, ...}
```

C.1.2 Training and Evaluation Setup

For all experiments, we adopt a 5-fold cross-validation strategy to ensure the robustness of *J-Detector*.

- **Data Split:** Each dataset in JD-Bench is split into 80% for training and 20% for testing. To prevent data leakage, we ensure that candidates appearing in the training set do not appear in the test set.
- **Classifiers:** We implement *J-Detector* using LightGBM, RandomForest, and XGBoost via the scikit-learn library with default hyperparameters.
- **Runtime:** The feature extraction via LLM is a one-time offline process. Once features are extracted, the training of the lightweight classifiers takes less than 1 minute on a standard CPU (Intel i7-12700K).

1014 C.1.3 Dataset Statistics

1015 Table 3 provides a comprehensive summary of the
1016 datasets included in JD-Bench.

1017 C.2 Detailed Definition of Various Judgment 1018 Types

1019 Depending on the evaluation protocol, judgments
1020 can take multiple forms (Li et al., 2024): (i) *Score-*
1021 *based judgments*: $j \in \mathbb{R}$, such as a numerical rating
1022 on one or several dimensions; (ii) *Pairwise judg-*
1023 *ments*: $j \in \{(c_a \succ c_b), (c_b \succ c_a)\}$, indicating
1024 a preference between two candidates $c_a, c_b \in \mathcal{C}$;
1025 (iii) *Listwise judgments*: $j \in \pi(\mathcal{C})$, representing a
1026 permutation (ranking) π over a candidate set.

1027 C.3 JD-Bench Details

1028 To systematically study the detectability of LLM-
1029 generated judgments, we introduce **JD-Bench**, a
1030 large-scale benchmark that integrates diverse appli-
1031 cations, judgment types, and model sources. JD-
1032 Bench provides a unified testbed for evaluating
1033 both existing and newly proposed detectors under
1034 realistic settings.

1035 **Dataset Selection.** We construct JD-Bench by
1036 aggregating data from multiple domains and judg-
1037 ment types, ensuring broad coverage of evaluation
1038 practices:

- 1039 • **HelpSteer2** (Wang et al., 2024b): HelpSteer2 is an
1040 open-source dataset designed to train and evaluate
1041 reward models for helpfulness assessment of LLM-
1042 generated responses. It contains large-scale human-
1043 annotated pointwise judgments that assign numeri-
1044 cal scores to responses across diverse instruction-
1045 following tasks. The dataset covers multiple do-
1046 mains and languages, enabling robust generaliza-
1047 tion of reward models. Its fine-grained annotations
1048 make it a strong benchmark for pointwise/score-
1049 based evaluation.
- 1050 • **HelpSteer3** (Wang et al., 2025): HelpSteer3 ex-
1051 tends HelpSteer2 by collecting pairwise human
1052 preference data on LLM responses. Instead of ab-
1053 solute scores, annotators compare two candidate
1054 responses to the same prompt and indicate which
1055 is better, yielding high-quality comparative judg-
1056 ments. The dataset spans a wide range of tasks
1057 and languages, supporting cross-lingual preference
1058 modeling and fine-grained ranking evaluation.
- 1059 • **NeurIPS Review Dataset** (Yu et al., 2025): This
1060 dataset comprises a large collection of real aca-
1061 demic peer reviews from the NeurIPS conference,
1062 annotated with multi-dimensional scores such as

1063 soundness, novelty, clarity, and overall rating. It
1064 represents a domain where judgments are struc-
1065 tured, multi-faceted, and highly consequential. The
1066 dataset captures nuanced reviewing language and
1067 decision rationales, providing a challenging bench-
1068 mark for modeling human-like expert evaluation.
1069 It is especially valuable for studying judgment be-
1070 havior in formal and high-stakes settings.

- **ANTIQU** (Hashemi et al., 2020): ANTIQU is a
1071 benchmark for non-factoid question answering, fo-
1072 cused on ranking passages based on their relevance
1073 to user queries. It includes listwise relevance judg-
1074 ments collected from crowdworkers, where mul-
1075 tiple candidate documents are ordered according
1076 to their usefulness. The questions are open-ended
1077 and require deeper understanding rather than sim-
1078 ple fact retrieval, making the ranking task more
1079 challenging.

1080 Each dataset provides *human-labeled* judgments
1081 as a reliable reference. To complement these, we
1082 collect *LLM-generated* judgments following the
1083 judging principles outlined in the respective papers,
1084 ensuring consistency in evaluation criteria.

1085 **LLM Selection.** To obtain LLM-generated judg-
1086 ments, we employ a diverse set of both closed-
1087 source and open-source models across a wide range
1088 of sizes and model families. This diversity is es-
1089 sential to cover heterogeneous judgment patterns
1090 and to test detector generalization. Specifically,
1091 JD-Bench includes judgments from:

1092 • Closed-source models:

- 1093 – OpenAI series: GPT-4o, GPT-5-mini, o4-
1094 mini.
- 1095 – Anthropic series: Claude-Haiku-3.5,
1096 Claude-Haiku-3, Claude-3.5-Sonnet.
- 1097 – Google series: Gemini-2.0-Flash,
1098 Gemini-2.5-Flash, Gemini-1.5-Pro.

1099 • Open-source models:

- 1100 – LLaMA family: LLaMA-3.2-3B,
1101 LLaMA-3.1-8B, LLaMA-3.1-70B.
- 1102 – Qwen family: Qwen-2.5-7B, Qwen-2.5-
1103 72B, RM-R1, RISE-Judge.
- 1104 – Mistral family: Mistral-7B, Mistral-
1105 Small-24B.
- 1106 – DeepSeek series: DeepSeek-V3,
1107 DeepSeek-R1.

1108 This mixture of datasets and models results in
1109 a benchmark that is both large-scale and diverse:
1110

Table 3: Summary of JD-Bench datasets, including input types and evaluation aspects.

Dataset	Candidate Input	Judgment Type	Annotation Aspects
HelpSteer2	LLM Prompt & Response	Pointwise	Helpfulness, Correctness, etc.
HelpSteer3	LLM Prompt & Response Pair	Pairwise	Overall Preference
NeurIPS	Academic Paper	Multi-dim	Soundness, Novelty, Confidence
ANTIQUA	Question & Doc Collection	Listwise	Relevance Ranking

JD-Bench covers *multiple application scenarios, different judgment types* (score, pairwise, listwise), and *a wide spectrum of LLM families*, making it a comprehensive resource for advancing judgment detection research. Table 4 presents the statistics of JD-Bench.

Prompt for JD-Bench Construction

HelpSteer2 Prompt (Pointwise, 5-Dimension Scoring)

Given a prompt and a response, follow the rubric to make a judgment.

Rubric:
Judge the response on five aspects: **helpfulness**, **correctness**, **coherence**, **complexity**, and **verbosity**. Assign each aspect a scalar score in [0, 4].

Prompt: [PROMPT]

Response: [RESPONSE]

Please output a valid JSON object using the following schema: "Rationale": <explanation for the given scores>, "Helpfulness": <0-4>, "Correctness": <0-4>, "Coherence": <0-4>, "Complexity": <0-4>, "Verbosity": <0-4>

Formatted the abovementioned schema and produce the judgment JSON now.

HelpSteer3 Prompt (Pairwise Comparison)

Given a prompt and two responses, follow the rubric to make a comparative judgment.

Rubric: Compare **Response 1** and **Response 2** along five aspects: **helpfulness**, **correctness**, **coherence**, **complexity**, and **verbosity**. Assign a single comparative score in -3,-2,-1,0,1,2,3 using the scale: -3: R1 much better than R2; -2: R1 better than R2; -1: R1 slightly better than R2; 0: about the same; 1: R2 slightly better than R1; 2: R2 better than R1; 3: R2 much better than R1.

Prompt (conversation/context): [CONTEXT AS FLATTENED TEXT]

Response 1: [RESPONSE_1]

Response 2: [RESPONSE_2]

Please output a valid JSON object using the following schema: "Rationale": <explanation for the comparative score>, "Score": <-3|-2|-1|0|1|2|3>

Formatted the abovementioned schema and produce the judgment JSON now.

NeurIPS Review Prompt (Structured JSON Review)

You are an AI researcher reviewing a paper submitted to a prestigious AI conference. Thoroughly evaluate the paper, adhering to the provided guidelines, and return a detailed assessment in the specified JSON format.

Manuscript: [MANUSCRIPT TEXT OR CONCATENATED CHUNKS]

Reviewer Guidelines (dimensions to cover):

Summary: Briefly summarize contributions (no critique here).

Strengths & Weaknesses across: Originality, Quality, Clarity, Significance.

Provide **Questions** for authors (useful for rebuttal).

Discuss **Limitations** and potential societal impact.

Flag **Ethical Concerns** if applicable (per conference policy).

Assign numerical ratings: **Soundness**, **Presentation**, **Contribution** (1-4 each).

Provide an **Overall** score (1-10) and **Confidence** (1-5).

Output a valid JSON object with the following fields: "Summary": <summary for the paper>, "Questions": <questions for the author>, "Limitations": <limitations for the paper>, "Soundness": <1-4>, "Presentation": <1-4>, "Contribution": <1-4>, "Overall": <1-10>, "Confidence": <1-5>

Formatted the abovementioned schema and produce the review JSON now.

Dataset	HelpSteer2	HelpSteer3	NeurIPS	ANTIQUÉ
Application	Resp. Eval.	Resp. Eval.	Peer Review	Doc Ranking
Judgment Type	Pointwise	Pairwise	Pointwise	Listwise
Judgment Dims	Helpfulness, Correctness, Coherence, Complexity, Verbosity	Overall	Overall, Confidence, Soundness, Presentation, Contribution	Relevance
Rating Scale	0–4	–3–3	1–10 / 1–5 / 1–4	1–4
#Train / #Test	62,961 / 21,778	62,880 / 42,317	63,210 / 62,664	102,417 / 61,909

Table 4: Overview of datasets included in JD-Bench.

ANTIQUÉ Prompt (3-Way Relevance Ranking)

Given a prompt and three responses, follow the rubric to assess relevance and rank the responses.

Rubric (per-response relevance score in [1, 4]): **4**: Reasonable and convincing; on par with or better than a likely correct answer. **3**: Possibly an answer, but not sufficiently convincing; a better-quality answer likely exists. **2**: Not an acceptable answer; unreasonable or does not address the question, but still on-topic. **1**: Completely out of context or nonsensical.

Prompt: [QUERY]

Response 1: [RESPONSE_1]

Response 2: [RESPONSE_2]

Response 3: [RESPONSE_3]

Please output a valid JSON object using the following schema: "Rationale": <explanation for your judgment and ranking>, "Response1 Score": <1-4>, "Response2 Score": <1-4>, "Response3 Score": <1-4>, "Ranking": <list of indices indicating best→worst, e.g., [0,1,2]>

Formatted the abovementioned schema and produce the judgment JSON now.

C.4 J-Detector Details

C.4.1 Linguistic Features

We extract a comprehensive set of surface, lexical, syntactic, and discourse indicators from each candidate response using spaCy-based parsing pipelines.

- **Length & Structure:** word_count, char_count, sentence_count, avg_sentence_length, list_count (bullet or numbered lists), paragraph_count, punctuation_count, reference_count (e.g., URLs).

- **Lexical Diversity:** unique_words,

- vocab_diversity (unique/total word ratio), 1135
- average_word_length, noun_verb_ratio, 1136
- adjective_ratio, adverb_ratio, 1137
- pronoun_ratio, contraction_rate. 1138

- **Readability:** coleman_liau index. 1139

- **Syntactic Complexity:** syntax_tree_depth (maximum dependency depth), 1140
average_dependency_length, 1142
passive_voice_ratio (fraction of sentences with nsubjpass/csubjpass), 1143
subordinate_clause_rate (rate of mark tokens). 1144

- **Discourse/hedging:** hedging_frequency (occurrence of hedge words such as “may”, “possibly”), discourse_marker_rate (connectives such as “however”, “moreover”). 1147

These features are computed for each response independently. For pairwise or listwise datasets (e.g., HelpSteer3, ANTIQUÉ), we additionally compute *difference features* such as $r_1 - r_2$ on each scalar dimension when comparing two responses. 1151

C.4.2 LLM-Enhanced Features 1156

Beyond surface-level indicators, we harness powerful large language models (e.g., Qwen3-8B) to derive task-aligned evaluation features. For each dataset, the model is prompted with the original instruction or query together with its candidate responses, and asked to generate structured JSON judgments that include detailed rationales and aspect-specific scores. 1157

Pointwise Setting (e.g., HelpSteer2). Each response is scored independently along eight stylistic and content dimensions: 1158

- Style, Format, Wording 1168

- Helpfulness, Correctness, Coherence 1169

1170	• Complexity, Verbosity		
1171	The model outputs both a natural language rationale and numeric scores (0–4) per dimension plus an overall_score.		
1172		– <i>Pairwise comparisons</i> : keys such as pairwise, pairs, comparisons, or prefs.	1218
1173		– <i>Ranking lists</i> : an explicit ranking field if available.	1219
1174	Pairwise Setting (e.g., HelpSteer3). Two responses are jointly compared under criteria such as <i>helpfulness</i> , <i>correctness</i> , <i>coherence</i> , <i>complexity</i> , and <i>verbosity</i> . The LLM produces a signed comparison score from -3 (Response 1 \gg Response 2) to $+3$ (Response 2 \gg Response 1) and a supporting rationale.	– <i>Metadata</i> : optional question/prompt/task descriptions to provide minimal context.	1220
1175			1221
1176			1222
1177		The resulting text is tokenized and directly used as the classifier input.	1223
1178			1224
1179			1225
1180			1226
1181	Listwise Setting (e.g., ANTIQUE). Three responses are simultaneously ranked by relevance. The LLM assigns a 1–4 relevance score to each response and outputs an ordered ranking list $[0, 1, 2]$ to indicate relative quality.	• Judgment + Candidate In this richer setting, we augment the above judgment text with the <i>candidate contents</i> being judged. Candidate responses are extracted from dataset fields such as:	1227
1182		– <code>examples[*].docs</code> for passage-style corpora (e.g., ANTIQUE);	1228
1183		– <code>examples[*].context</code> for conversational datasets (e.g., HelpSteer3), where only assistant turns are kept;	1229
1184		– top-level docs, candidates, or answers if present.	1230
1185			1231
1186	Long-form Paper Evaluation (e.g., NeurIPS Submissions). For full papers, we ask the model to return review-like signals: style, format, wording (0–4), rating (1–10), confidence (1–5), soundness/presentation/contribution (1–4 each), together with detailed reasoning.		1232
1187			1233
1188			1234
1189			1235
1190			1236
1191			1237
1192	These LLM-enhanced features provide semantically rich, high-level signals that complement the surface-level linguistic statistics, enabling our detector to exploit both human-interpretable cues and task-specific, model-derived evaluations.		1238
1193			1239
1194			1240
1195			1241
1196			1242
1197	C.5 SLM-based Method Details		1243
1198	To benchmark the ability of small language models (SLMs) to discriminate between human and LLM-generated judgments, we adapt text classification pipelines with two input configurations: <i>judgment-only</i> (w/o candidates) and <i>judgment+candidate</i> (w/ candidates). Both settings train a binary classifier to predict whether a group of judgments originates from a human annotator (label 0) or an LLM (label 1). We employ roberta-base and allenai/longformer-base-4096 as backbones, with max sequence lengths 512 and 4096, respectively.	JudgmentText === Candidates === Candidate ₁ ... Candidate _n	1244
1200			
1201			
1202			
1203			
1204			
1205			
1206			
1207			
1208			
1209			
1210			
1211			
1212			
1213			
1214			
1215	• Judgment-Only Inspired by SLM-based text detection, this setting feeds only the <i>judgment artifacts</i> into the model. Each group is represented by a textualized summary of available signals, including:		
1216	– <i>Numeric scores</i> : fields such as rating, score, confidence, soundness, presentation, contribution, etc.		
1217			
			1245
			1246
			1247
			1248
			1249
			1250
			1251
			1252
			1253
			1254
			1255
			1256
			1257

Dataset Setting	LLM-Generated Feature Dimensions
HelpSteer2 (pointwise)	Style, Format, Wording, Helpfulness, Correctness, Coherence, Complexity, Verbosity, Overall
HelpSteer3 (pairwise)	Helpfulness, Correctness, Coherence, Complexity, Verbosity, Pairwise Score ($-3 - +3$)
ANTIQUA (listwise)	Response relevance scores (1–4), Ranking order, Rationale
NeurIPS (pointwise)	Style, Format, Wording, Rating (1–10), Confidence (1–5), Soundness, Presentation, Contribution

Table 5: Example LLM-enhanced feature dimensions by dataset.

Mode	Input Composition	Example Fields Used
w/o candidates	Judgments only	ratings, scores, pairwise, ranking, task
w/ candidates	Judgments + trimmed candidate texts	docs, context (assistant turns), answers

Table 6: Two input modes for SLM-based judgment detection.

1258	• Judgment-only signals: helpfulness, correctness, coherence, complexity, verbosity, ranking, and pairwise preferences.	• pairwise: sample k random Human–LLM pairs.	1291
1259			1292
1260			
1261	• Optional candidates: trimmed prompt/response or passage text to provide weak context.	3. Feed each pair to an LLM agent that proposes actions to maintain a Feature Bank :	1293
1262			1294
1263	We prompt the detector LLM with an instruction template asking it to decide whether the judgments were written by a <i>Human</i> or by an <i>LLM</i> , based on style, consistency, and calibration artifacts:	Add: {"name": "..."} 1295	1295
1264		Delete: {"name": "..."} 1296	1296
1265		Merge: {"name": "..."} 1297	1297
1266			1298
1267	{	4. Typical mined features include:	1299
1268	"Rationale": "<brief explanation>",	• Length or verbosity bias;	1300
1269	"Prediction": "Human" "LLM"	• Overly smooth or formulaic score patterns;	1301
1270	}	• Deterministic tone and calibration artifacts.	1302
1271	Two modes are supported:	The resulting Feature Bank $\mathcal{F}_{\text{sample}}$ captures diagnostic cues distilled by the LLM itself and is later injected into the final detection prompt.	1303
1272	• judgment_only: only judgment artifacts are provided.		1304
1273			1305
1274	• enable_candidate: judgment artifacts plus trimmed candidate texts.	C.6.3 Distribution-Level LLM-Based Analysis	1306
1275		Beyond individual samples, we analyze dataset-wide statistics to extract global signals of LLM-generated judgments:	1307
1276	This baseline does not use any explicit feature engineering but leverages the LLM’s implicit ability to reason about stylistic and distributional cues.		1308
1277			1309
1278			
1279	C.6.2 Sample-Level LLM-Based Analysis	1. Compute per-label histograms and descriptive statistics for all available judgment dimensions (e.g., helpfulness, correctness, coherence, complexity).	1310
1280	We further design an agentic feature mining procedure to expose regularities in Human vs. LLM judgments at the <i>instance level</i> . Given a training set of groups, we:		1311
1281			1312
1282			1313
1283			
1284	1. Flatten them into a table of <i>prompt</i> , <i>response</i> , <i>label</i> , <i>scores</i> , and derived metrics such as length and average score.	2. Analyze correlations :	1314
1285		• Length–score Spearman correlations within Human/LLM groups;	1315
1286		• Cross-dimension correlations (e.g., helpfulness vs. coherence).	1316
1287	2. Mine Human–LLM pairs using two strategies:		1317
1288			1318
1289	• scoring: select k pairs with the largest average-score gaps under the same prompt.	3. Summarize these findings as structured text and feed them to an LLM to propose additional high-level features, such as:	1319
1290		• Consistent score calibration (LLM often shows smaller variance);	1320
			1321
			1322
			1323

- Stronger length–score coupling in LLM judgments;
- Reduced inter-dimension diversity compared to human raters.

The discovered global patterns augment the feature bank as $\mathcal{F}_{\text{dist}}$, complementing sample-level cues with distributional regularities.

C.6.4 Final Detection

The final detector integrates:

- A **Feature Bank** $\mathcal{F} = \mathcal{F}_{\text{sample}} \cup \mathcal{F}_{\text{dist}}$;
- Group-level summaries (judgments + optional candidates).

An LLM receives this structured prompt and outputs the final label prediction:

$$\hat{y} = f_{\text{LLM}}(\text{summary}(G), \mathcal{F}),$$

where f_{LLM} denotes the LLM-based reasoning process conditioned on both the mined features and the group payload.

In practice, the multilevel detector (sample + distribution) consistently improves accuracy by guiding the LLM with both fine-grained instance cues and global dataset regularities.

D Theoretically Analysis on LLM-generated Judgment Detectability

We model the detectability of whether a group of judgments G (scores, pairwise preferences, or listwise rankings) was produced by a human or an LLM. Let m denote the group size, d the number of attribute dimensions, and S the effective rating scale cardinality:

$$S = \begin{cases} L, & \text{for } L\text{-level scoring;} \\ 2x + 1, & \text{for pairwise judgments with } x \in \mathbb{Z}_{\geq 1} \text{ superiority levels per side (including tie);} \\ k!, & \text{for a full ranking over } k \text{ candidates.} \end{cases}$$

The per-judgment information is $\log S$ nats.²

Let P_H and P_M be the conditional distributions over judgment outcomes induced by humans and LLMs, respectively. Denote $\Delta = \text{TV}(P_H, P_M)$ as their total variation distance.

²For listwise $k!$, Stirling’s approximation gives $\log(k!) \approx k \log k - k$. For continuous pairwise margins, discretization into B bins yields $S = B$.

From sample complexity to group detectability. With n i.i.d. observations, the total variation between product distributions grows as

$$\text{TV}(P_H^{\otimes n}, P_M^{\otimes n}) = 1 - \exp\{-nI_c(P_H, P_M) + o(n)\},$$

where I_c is the Chernoff information, scaling quadratically with Δ . In our setting, the effective observation budget is

$$n_{\text{eff}} = m \cdot d \cdot \log S,$$

which accounts for group size, dimensionality, and rating resolution.

Detectability index. Thus, the detectability index becomes

$$\text{Det}(G) = 1 - \exp\{-\beta m d \log S \Delta^2\},$$

where $\beta > 0$ is dataset- and model-dependent. The detectability increases monotonically with four factors: (i) rating scale S , (ii) attribute dimensions d , (iii) group size m , and (iv) distribution gap Δ .

Instantiation by type. For L -level scores, use $S = L$. For pairwise preferences, use $S = L_{\text{pair}}$ (e.g., 7 for $\{-3, \dots, 3\}$). For listwise ranking over k items, use $S = k!$ (or $\log S \approx k \log k - k$). For mixed-type groups, sum $m d \log S$ across instances.

Method	Uses Candidates?	Feature Bank	Level of Analysis
LLM-as-a-Judge	Optional	None	Per-group
Sample-level	Optional	$\mathcal{F}_{\text{sample}}$	Instance-level
Distribution-level	Optional	$\mathcal{F}_{\text{sample}} + \mathcal{F}_{\text{dist}}$	Global + per-group

Table 7: Comparison of the three LLM-based detection strategies.