
Diagnosing LLM Judge Reliability: Conformal Prediction Sets and Transitivity Violations

Anonymous Authors¹

Abstract

LLM-as-judge frameworks are increasingly used for automatic NLG evaluation, yet their per-instance reliability remains poorly understood. We present a two-pronged diagnostic toolkit applied to SummEval: **(1)** a transitivity analysis that reveals widespread per-input inconsistency masked by low aggregate violation rates ($\bar{p} = 0.8\text{--}4.1\%$), with 33–67% of documents exhibiting at least one directed 3-cycle; and **(2)** split conformal prediction sets over 1–5 Likert scores providing theoretically-guaranteed $\geq (1-\alpha)$ coverage, with set width serving as a per-instance reliability indicator ($r_s = +0.576$, $N=1,918$, $p < 10^{-100}$, pooled across all judges). Critically, prediction set width shows consistent cross-judge agreement ($\bar{r} = 0.32\text{--}0.38$), demonstrating it captures *document-level difficulty* rather than judge-specific noise. Across four judges and four criteria, both diagnostics converge: *criterion matters more than judge*, with *relevance* judged most reliably (avg. set size ≈ 3.0) and *coherence* moderately so (avg. set size ≈ 3.9), while *fluency* and *consistency* remain unreliable (avg. set size ≈ 4.9). We release all code, prompts, and cached results.

1. Introduction

Automatic evaluation of natural language generation (NLG) has become a cornerstone of modern NLP research. LLM-as-judge systems, where a large language model scores or ranks system outputs, have gained rapid adoption as scalable proxies for human annotation (Zheng et al., 2023; Liu et al., 2023; Fu et al., 2023). A single LLM call can replace expensive crowd-sourced annotation pipelines, and practitioners increasingly treat these scores as ground truth.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Yet a critical question goes largely unasked: *when should you trust an LLM judge?* Aggregate metrics such as system-level Kendall’s τ or Pearson correlation with human scores look impressive, but they average over hundreds of instances. A judge that is right 90% of the time can be spectacularly wrong on the 10% that matters most. **This paper develops and evaluates two complementary diagnostics for per-instance reliability.**

Contribution 1: Transitivity diagnostic. Directed cycles ($A \succ B$, $B \succ C$, $C \succ A$) arise naturally when alternatives are near-equally preferred (Young, 1988). We measure 3-cycle violation rates across four judges on SummEval (Fabbri et al., 2021), crucially *disaggregated by input document*. Aggregate rates are low ($\bar{p} < 5\%$), yet 33–67% of documents exhibit at least one violation, with per-document rates reaching 30.4% for Mistral-Small-3.1. We also test whether MFAS ranking repair (Ailon et al., 2008) improves human ranking agreement—it does not, confirming violations are sparse noise rather than systematic bias.

Contribution 2: Conformal prediction diagnostic. We apply split conformal prediction (Vovk et al., 2005; Angelopoulos & Bates, 2021) to direct Likert scores, producing prediction sets with finite-sample, distribution-free coverage guarantees. The central finding is that prediction set *width* is a statistically robust per-instance reliability indicator: pooling 1,918 observations across all judges and criteria yields Spearman $r_s = +0.576$ ($p < 10^{-100}$) between width and actual judge–human disagreement. Moreover, different judges assign wide sets to the *same documents* ($\bar{r} = 0.32\text{--}0.38$ across judge pairs for fluency, consistency, and relevance), demonstrating that width measures inherent document difficulty rather than judge-specific noise.

Unified finding. Both diagnostics independently identify the same axis of variation: *criterion* explains reliability more than *judge*. This is actionable: a practitioner deploying LLM judges should trust coherence and relevance scores more than fluency and consistency scores, regardless of which judge is used.

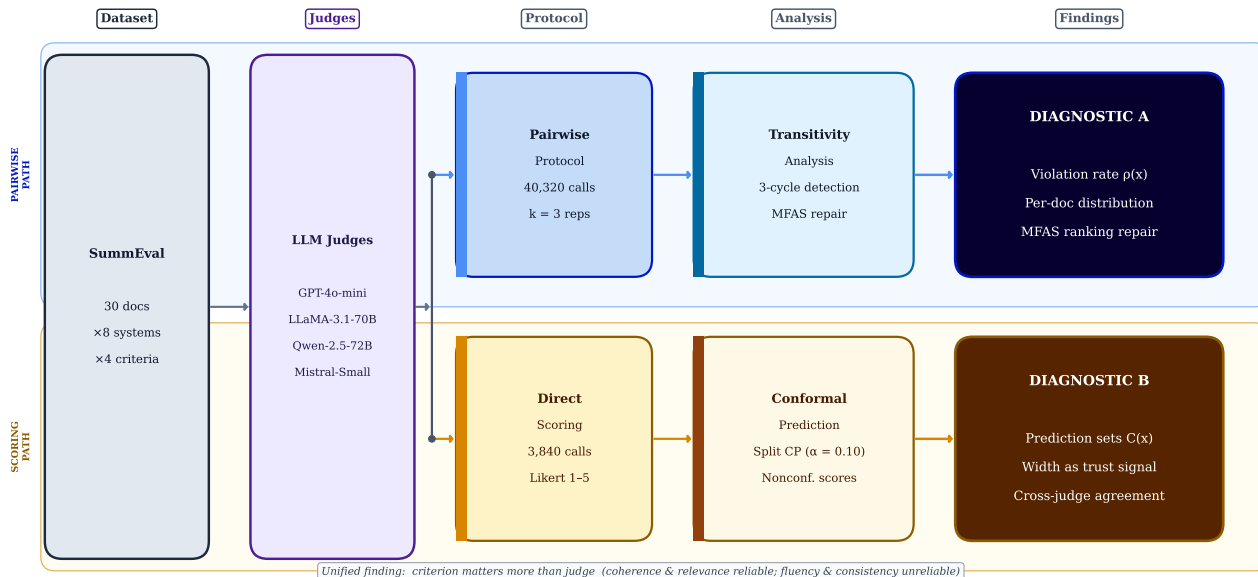


Figure 1. **Two-pronged diagnostic pipeline.** SummEval documents are evaluated by four LLM judges under two protocols. The *pairwise protocol* (40,320 API calls) feeds the transitivity diagnostic, which measures directed 3-cycle violation rates $\rho(x)$ per input and tests whether MFAS ranking repair improves agreement with human rankings. The *direct scoring protocol* (3,840 API calls) feeds the conformal diagnostic, which produces prediction sets $\mathcal{C}(x)$ with guaranteed $\geq (1-\alpha)$ coverage; set width serves as a per-instance trust signal.

2. Related Work

LLM-as-judge reliability. Zheng et al. (2023) introduced MT-Bench and Chatbot Arena, establishing LLM judges as scalable evaluation tools. Liu et al. (2023) showed that GPT-4-based G-Eval correlates strongly with human judgment on SummEval. Known biases include position (Wang et al., 2023), verbosity (Saito et al., 2023), and self-enhancement effects. Fernandes et al. (2023) and Koo et al. (2023) audit LLM judges at scale, finding systematic weaknesses on specific input types, consistent with our per-document view. Concurrent to our work, Ye et al. (2024) study fine-grained reliability of LLM judges across skill categories, but without formal uncertainty guarantees.

Transitivity and ranking consistency. Condorcet cycles in pairwise preferences have been studied in social choice theory since de Condorcet (1785) and are known to be ubiquitous when alternatives are near-equal (Young, 1988; Moon, 1968). MFAS-based ranking repair has been applied to preference aggregation (Ailon et al., 2008) and recently to LLM-generated ranked lists (Qin et al., 2024). We are the first to measure directed 3-cycle rates in LLM judges at the per-document level and connect them to conformal uncertainty.

Conformal prediction in NLP. Split conformal prediction (Vovk et al., 2005; Papadopoulos et al., 2002) provides distribution-free coverage guarantees; Angelopoulos &

Bates (2021) give a modern treatment. Applications to NLP include uncertainty for machine translation (Fomicheva et al., 2020), text classification (Maltoudoglou et al., 2020), and question answering (Quach et al., 2024). Kumar et al. (2023) apply conformal methods to LLM generation quality, while Kuhn et al. (2023) propose semantic entropy as a complementary uncertainty signal. Our work is the first to apply conformal prediction to *LLM-as-judge scores* and to interpret prediction set width as a per-instance deployment signal.

3. Methods

Figure 1 illustrates our pipeline. Both diagnostics share the same four judges, dataset, and criteria, enabling direct comparison of their findings.

3.1. Transitivity Diagnostic

Tournament formulation. For each input document x and set of n system outputs $\mathcal{S} = \{s_1, \dots, s_n\}$, a pairwise judge produces a tournament $G = (\mathcal{S}, E)$ where $(s_i, s_j) \in E$ iff the judge prefers s_i over s_j . A *transitivity violation* is a directed 3-cycle: $s_i \succ s_j, s_j \succ s_k, s_k \succ s_i$.

Per-document violation rate. We define the violation rate for document x as:

$$\rho(x) = \frac{\# \text{ directed 3-cycles in } G_x}{\binom{n}{3}}, \quad (1)$$

normalized by the total number of possible triples. We report the aggregate mean $\bar{\rho}$, the fraction of documents with $\rho(x) > 0$, and the full distribution $\{\rho(x)\}$.

Ranking methods. We compare five ranking methods: **Win Rate** (fraction of pairwise wins), **Bradley-Terry** (Bradley & Terry, 1952) (maximum-likelihood strength scores), **Schulze** (Schulze, 2011) (beatpath method), **MFAS-ILP** (exact Minimum Feedback Arc Set via integer linear programming), and **MFAS-Copeland** (Copeland scores as a fast MFAS approximation). We measure agreement with human rankings via Kendall’s τ .

Repetitions. Each pair is queried $k=3$ times per judge to measure win-rate confidence (0, 1/3, 2/3, or 1).

3.2. Conformal Prediction Diagnostic

Setup. We use split conformal prediction (Vovk et al., 2005) in the *direct scoring* setting: the judge assigns a Likert score $\hat{y} \in \{1, \dots, 5\}$, and the calibration target is the rounded average human score $y^* \in \{1, \dots, 5\}$.

Nonconformity score and prediction set. We use the absolute residual as the nonconformity score: $s_i = |\hat{y}_i - y_i^*|$. Given calibration set $\{(x_i, y_i^*)\}_{i=1}^n$, the conformal threshold is:

$$\hat{q} = s_{(\lceil (1-\alpha)(n+1) \rceil)}, \tag{2}$$

the appropriate empirical quantile ensuring marginal coverage $\mathbb{P}(y^* \in \mathcal{C}(x)) \geq 1 - \alpha$ (Tibshirani et al., 2019). The prediction set for a new instance with judge score \hat{y} is:

$$\mathcal{C}(x) = \{y \in \{1, \dots, 5\} : |\hat{y} - y| \leq \hat{q}\}. \tag{3}$$

Set width $w(x) = |\mathcal{C}(x)|$ ranges from 1 (maximally confident) to 5 (full uncertainty).

Evaluation protocol. We evaluate (1) empirical coverage vs. target $1-\alpha$; (2) average set size as an informativeness measure; (3) Spearman $r_s(w, |\hat{y} - y^*|)$ to quantify how well set width predicts actual judge error; and (4) inter-judge width agreement, the Spearman correlation between widths assigned by two different judges to the same document, which tests whether width reflects document-level difficulty or judge-specific noise. All metrics are averaged over 20 random 50/50 calibration/test splits for stable estimates.

4. Experimental Setup

Dataset. SummEval (Fabbri et al., 2021) contains 100 documents \times 16 systems (= 1,600 outputs) with human Likert scores (1–5) on *coherence*, *consistency*, *fluency*, and *relevance*, averaged over three annotators. We subsample to 30 documents \times 8 systems (systems 0, 2, 4, 6, 9, 11, 13,

Judge	Agg. $\bar{\rho}$	% docs ≥ 1	Max ρ	Med. ρ
LLAMA-3.1-70B	0.008	33.3%	3.6%	0.0%
QWEN-2.5-72B	0.022	50.0%	14.3%	1.8%
MISTRAL-SMALL	0.041	50.0%	30.4%	1.8%
GPT-4O-MINI	0.014	46.7%	7.1%	0.0%

Table 1. Per-document transitivity violation statistics (*coherence*). Aggregate rates $\bar{\rho}$ appear low (<5%), yet 33–50% of documents exhibit at least one directed 3-cycle, and per-document rates reach 30.4% for MISTRAL-SMALL. Median = 0 for all judges: most documents are violation-free, but a minority expose severe judge inconsistency.

15) for cost efficiency, rounding averaged human scores to the nearest integer for conformal calibration.

Judges. We evaluate four instruction-tuned LLMs accessed via OpenRouter: GPT-4O-MINI (gpt-4o-mini), LLAMA-3.1-70B (meta-llama/llama-3.1-70b-instruct), QWEN-2.5-72B (qwen/qwen-2.5-72b-instruct), and MISTRAL-SMALL-3.1 (mistralai/mistral-small-3.1-24b-instruct). All responses are cached in SQLite.

5. Results

5.1. Transitivity Violations: Aggregate Rates Mask Per-Document Heterogeneity

Table 1 reveals a fundamental measurement problem. Aggregate violation rates $\bar{\rho}$ range from 0.8% to 4.1% across judges, statistics that would reassure any practitioner. But the per-document view tells a different story: between 33% and 50% of documents carry at least one directed 3-cycle, and the worst document for MISTRAL-SMALL has $\rho(x) = 30.4\%$, meaning nearly a third of all triples form preference cycles on that input. This is not random noise. Figure 2 shows that the per-document distribution is heavily right-tailed: the median is zero across all judges, but a small fraction of documents, typically those where system outputs differ subtly in quality, drive the aggregate statistic. A practitioner who relies on $\bar{\rho}$ alone will miss this concentrated inconsistency.

Appendix C (Table 6) provides violation statistics across all four criteria. *Fluency* and *consistency* show the highest fraction of documents with ≥ 1 violation (up to 66.7% for LLaMA–fluency), converging with the conformal results below.

MFAS ranking repair does not help. Table 2 compares five ranking methods against the SummEval human gold standard on *coherence*. MFAS-ILP achieves the highest Kendall’s τ on one judge but matches or falls below Win Rate on the other three. No method consistently dominates.

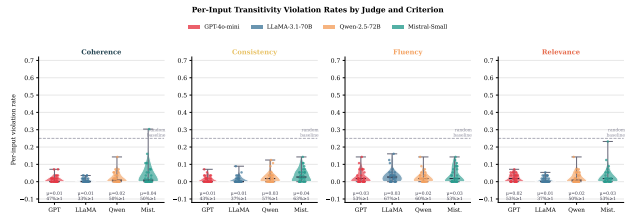


Figure 2. Per-document violation rate distributions. Each violin shows the distribution of $\rho(x)$ across 30 documents for one judge-criterion pair. Dashed horizontal line: random-baseline rate (0.25). All distributions are right-tailed with median = 0, but the upper tails, where a single document can expose > 30% violation rates, are practically significant. Fluency consistently shows the widest tails.

This non-result is informative: when violations are sparse and concentrated in a few documents, the tournament graph is nearly acyclic and MFAS repair finds little to improve. The diagnostic value of transitivity analysis lies in identifying unreliable documents, not in repairing global rankings.

Method	LLaMA	Qwen	Mistral	GPT
Win Rate	0.571	0.643	0.071	0.714
Bradley-Terry	0.571	0.643	0.071	0.714
Schulze	0.643	0.643	0.071	0.643
MFAS-Copeland	0.571	0.643	0.071	0.714
MFAS-ILP [†]	0.643	0.643	0.071	0.643

Table 2. Ranking agreement (Kendall’s τ) vs. SummEval human gold standard (coherence). MFAS does not consistently outperform Win Rate or Bradley-Terry. When violations are sparse, standard aggregation already captures the signal; MFAS repair adds noise rather than structure.

[†] Exact ILP solution; **bold** = best per judge.

5.2. Conformal Prediction Sets: Guaranteed Coverage with a Trust Signal

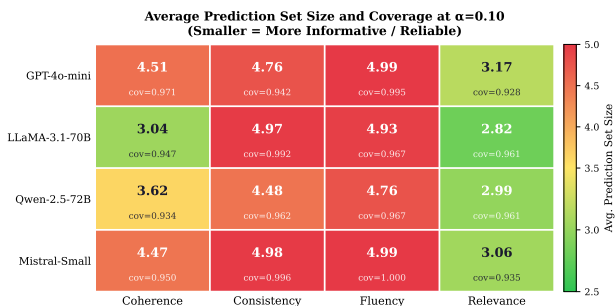


Figure 3. Average prediction set size at $\alpha=0.10$ (green = small = reliable; red = large = unreliable). Each cell shows average set size (larger text) and empirical coverage (smaller text). The criterion axis drives variation far more than the judge axis: coherence and relevance (left two columns) are reliably judged (≈ 3.0), while fluency and consistency are near-maximally uncertain (≈ 5.0). All 16 cells meet the 90% coverage guarantee.

Judge		Coh.	Rel.	Con.	Flu.	$r_s(w,e)$
GPT-4o-MINI	Size	4.51	3.17	4.76	4.99	+0.18 / +0.25
	Cov.	.971	.928	.942	.995	
LLAMA-3.1-70B	Size	3.04	2.82	4.97	4.93	-0.03 / +0.08
	Cov.	.947	.961	.992	.967	
QWEN-2.5-72B	Size	3.62	2.99	4.48	4.76	+0.14 / +0.65
	Cov.	.934	.961	.962	.967	
MISTRAL-SMALL	Size	4.47	3.06	4.98	4.99	-0.02 / +0.17
	Cov.	.950	.935	.996	1.00	

Table 3. Conformal results at $\alpha=0.10$ (90% target coverage). Columns ordered: reliable criteria then unreliable criteria. $r_s(w,e)$ = Spearman correlation of set width with absolute error, reported as Coh./Con. (coherence and consistency). QWEN-2.5-72B achieves the strongest width-error signal for consistency ($r_s = +0.65$).

Coverage guarantee satisfied. The conformal guarantee holds for all 16 judge \times criterion combinations across all four α levels tested ($\alpha \in \{0.05, 0.10, 0.15, 0.20\}$). Figure 6 confirms that empirical coverage tracks the theoretical $1-\alpha$ line, meeting or exceeding it at every operating point. Full coverage and set-size results are in Appendix B.

Criterion drives reliability, not judge. Table 3 and Figure 3 show the dominant pattern: the criterion column explains far more variance in average set size than the judge row. Coherence and relevance receive small sets (2.82–4.51 labels) while fluency and consistency receive near-maximal sets (4.47–4.99) across all four judges.

This is not a failure of the method. SummEval neural summaries are uniformly fluent, leaving little variance for a judge to exploit (Fabbri et al., 2021). Consistency requires cross-document factual reasoning that 24–72B scale models perform inconsistently. The wide prediction sets are the method correctly reporting that it cannot reliably grade these criteria.

Set width predicts judge error. Figure 4 shows pooled reliability diagrams for each criterion: mean absolute error (MAE) vs. prediction set width, pooled across all four judges. Thirteen of sixteen judge \times criterion combinations show a perfectly monotonic width-error relationship. Pooling all 1,918 observations, Spearman $r_s = +0.576$ ($p < 10^{-100}$). Per-criterion pooled correlations are $r_s = +0.34$ (consistency), $r_s = +0.16$ (coherence), $r_s = +0.15$ (fluency), and $r_s = -0.01$ (relevance, n.s.).

The three non-monotonic cases, LLaMA fluency, LLaMA relevance, and Mistral coherence, share a common structure: only two distinct width values are observed (widths 2 and 3), leaving insufficient variance for a meaningful per-judge correlation. The width-error relationship is real but requires pooling to emerge against this discrete noise floor.

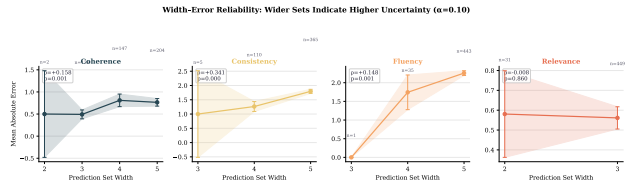


Figure 4. Pooled reliability diagrams (all four judges, $\alpha=0.10$). x -axis: prediction set width; y -axis: mean absolute error (MAE) vs. human score. Error bars: 95% CI. Annotations: sample count per width. Spearman r_s and p -value shown per panel. *Consistency* shows the clearest signal ($r_s = +0.34, p < 0.0001$); *relevance* is the exception ($r_s \approx 0, p = 0.86$).

Judge Pair	Coh.	Con.	Flu.	Rel.
GPT / LLaMA	+0.12	+0.22†	+0.45‡	+0.28‡
GPT / Qwen	+0.27‡	+0.38‡	+0.23†	+0.70‡
GPT / Mistral	+0.31‡	+0.32‡	+0.81‡	+0.29‡
LLaMA / Qwen	-0.08	+0.17	+0.21†	+0.20†
LLaMA / Mistral	+0.08	+0.56‡	+0.23†	+0.25†
Qwen / Mistral	-0.09	+0.24‡	+0.33‡	+0.44‡
<i>Mean</i>	<i>+0.10</i>	<i>+0.32</i>	<i>+0.38</i>	<i>+0.36</i>

† $p < .05$ ‡ $p < .01$ or better **bold** $r > 0.50$

Table 4. Inter-judge agreement on prediction set width (Spearman $r, \alpha=0.10$). Positive values indicate different judges assign wider sets to the same documents, width reflects *document-level difficulty*, not judge-specific noise.

Width reflects document difficulty, not judge noise. A potential confound: perhaps wide sets indicate that a specific judge is unreliable, rather than that the document is inherently hard to evaluate. Table 4 addresses this directly by measuring inter-judge agreement on prediction set width: the Spearman correlation between widths assigned by two *different* judges to the same documents.

For *fluency*, *consistency*, and *relevance*, 15 of 18 judge pairs show significant positive width agreement ($p < 0.05$), with mean correlations of $\bar{r} = 0.38, 0.32$, and 0.36 respectively. The standout pairs, GPT/Mistral fluency ($r = +0.81$) and GPT/Qwen relevance ($r = +0.70$), confirm that different model families converge on the same documents as hard to judge. *Coherence* is the exception ($\bar{r} = 0.10$), consistent with its smaller and less variable set sizes.

Figure 5 visualizes the full inter-judge agreement matrices, making the qualitative pattern clear: coherence shows near-zero off-diagonal entries while fluency and relevance show predominantly warm (positive) off-diagonal colors.

5.3. Convergent Evidence: Criterion Drives Reliability

Both diagnostics, applied independently, reach the same conclusion. *Fluency* and *consistency* show the highest fraction of documents with ≥ 1 transitivity violation (LLaMA: 66.7% for fluency; Qwen: 60%; Table 6) and the widest

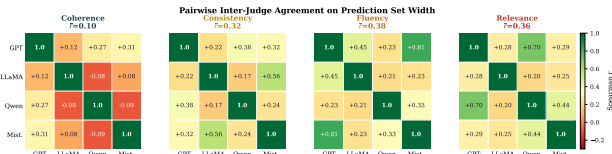


Figure 5. Inter-judge width agreement matrices (Spearman $r, \alpha=0.10$). Rows/columns: the four judges. Diagonal forced to 1.0. *Coherence* (leftmost) shows predominantly near-zero off-diagonal entries; *fluency* and *relevance* show consistently positive agreement, confirming that prediction width tracks document-level difficulty across model families.

prediction sets (avg. size > 4.9 for three of four judges). *Coherence* and *relevance* show lower violation rates and narrower prediction sets (avg. size < 4.0 for three of four judges).

This cross-diagnostic agreement is strong evidence that criterion difficulty is a fundamental property of the evaluation task, not an artifact of any single measurement approach. A document that causes preference cycles in pairwise ranking also receives a wide prediction interval in direct scoring, both signals point to the same underlying difficulty.

6. Discussion

The masked-heterogeneity problem. A practitioner who reports only $\bar{\rho} < 5\%$ or system-level $\tau > 0.5$ will conclude that their LLM judge is reliable. Our results show this conclusion is premature: nearly half of all documents expose judge inconsistency at the per-document level. We recommend that evaluation studies report, at minimum: (a) the fraction of documents with ≥ 1 violation; and (b) the distribution of per-document violation rates, not just the aggregate mean.

Prediction sets as a deployment signal. Our conformal results have an immediate practical implication: before accepting an LLM judge score, compute the prediction set.

- If $|\mathcal{C}(x)| \leq 2$: the judge is likely reliable for this instance, proceed.
- If $|\mathcal{C}(x)| = 5$ (full scale): the judge expresses maximum uncertainty, consider human annotation.

This selective escalation strategy is principled: the coverage guarantee ensures $\mathcal{C}(x)$ contains the human score with at least $1-\alpha$ probability. The cross-judge agreement results further justify this approach: a wide set from any judge is a warning about the document, not about that specific model.

Why MFAS does not help. When violations are sparse ($\bar{\rho} < 5\%$) and concentrated in a minority of documents, the

275 tournament graph is nearly acyclic. MFAS repair finds at
 276 most a handful of edges to reverse, insufficient to change
 277 aggregate rankings. Detecting these problematic documents
 278 for human follow-up is a more effective use of the transitivity
 279 signal than attempting automated repair.

280
 281
 282 **The coherence exception.** *Coherence* is the one criterion
 283 where inter-judge width agreement is weak ($\bar{r} = 0.10$). We
 284 hypothesize two reasons: (1) neural summaries in SummEval
 285 vary substantially in coherence, making it a more
 286 discriminable dimension; and (2) different model families
 287 may have different internal representations of “coherence,”
 288 leading to idiosyncratic scoring patterns that do not general-
 289 ize across judges.

290 291 7. Limitations

292
 293 **Scale and generalization.** We use 30 documents \times 8 sys-
 294 tems from SummEval. Results may differ on larger subsets,
 295 other summarization datasets, or non-summarization tasks
 296 (dialogue, translation, etc.).

297 **Marginal vs. conditional coverage.** Split conformal guar-
 298 antees *marginal* coverage $\mathbb{P}(y^* \in \mathcal{C}(x)) \geq 1 - \alpha$, not per-
 299 document conditional coverage. Harder documents may
 300 receive tighter-than-justified sets in practice; conditional
 301 conformal methods (Angelopoulos & Bates, 2021) could
 302 address this.

303
 304 **Fixed nonconformity score.** We use the absolute residual
 305 $|\hat{y} - y^*|$. Learned nonconformity scores (e.g., based on
 306 judge confidence or LLM log-probabilities) could produce
 307 tighter, more informative sets.

308 **Prompt sensitivity.** Each judge uses a single prompt tem-
 309 plate per criterion. Different prompts may yield different
 310 violation rates and set widths; we leave prompt-robustness
 311 analysis to future work.

312
 313 **Human score rounding.** SummEval provides averaged
 314 annotations; rounding to integers introduces a small dis-
 315 cretization error in calibration targets.

316 317 8. Conclusion

318
 319 We presented two complementary, low-cost diagnostics for
 320 LLM judge reliability. Transitivity analysis reveals that
 321 per-document inconsistency is dramatically higher than ag-
 322 gregate statistics suggest, a pattern invisible to standard
 323 evaluation metrics. Conformal prediction sets provide finite-
 324 sample coverage guarantees and a practical deployment sig-
 325 nal: prediction set width predicts actual judge error (pooled
 326 $r_s = +0.576$, $N=1,918$) and tracks document-level diffi-
 327 culty rather than judge-specific noise.

328
 329 Both diagnostics independently converge on the same find-

ing: *criterion matters more than judge*. Coherence and
 relevance can be judged reliably by any of the four models
 tested; fluency and consistency should be treated with skep-
 ticism regardless of the model. We recommend that LLM
 evaluation pipelines adopt per-instance uncertainty report-
 ing as standard practice, and we release all code, prompts,
 and cached API responses to support reproducibility.

Impact Statement

This work develops diagnostic tools for assessing the reli-
 ability of LLM-as-judge systems, with the goal of reducing
 uncritical reliance on automated evaluation. All experiments
 use publicly available data (SummEval) and commercially
 available LLMs accessed via a standard API. We do not
 collect human annotations, and the research poses no direct
 risk of harm to individuals. The broader impact is positive:
 surfacing systematic failure modes of LLM judges helps
 practitioners deploy automated evaluation more responsibly
 and identify when human oversight is needed.

References

- Ailon, N., Charikar, M., and Newman, A. Aggregating in-
 consistent information: Ranking and clustering. *Journal*
of the ACM, 55(5):1–27, 2008.
- Angelopoulos, A. N. and Bates, S. A gentle introduction
 to conformal prediction and distribution-free uncertainty
 quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- Bradley, R. A. and Terry, M. E. Rank analysis of incom-
 plete block designs: I. the method of paired comparisons.
Biometrika, 39(3/4):324–345, 1952.
- de Condorcet, M. J. A. N. C. Essai sur l’application de
 l’analyse à la probabilité des décisions rendues à la plu-
 ralité des voix. *Imprimerie Royale*, 1785.
- Fabrizi, A. R., Kryściński, W., McCann, B., Xiong, C.,
 Socher, R., and Radev, D. SummEval: Re-evaluating
 summarization evaluation. *Transactions of the Associa-
 tion for Computational Linguistics*, 9:391–409, 2021.
- Fernandes, P., Deutsch, D., Finkelstein, M., Riley, P., Mar-
 tins, A. F., Neubig, G., Garg, A., Clark, J. H., Freitag,
 M., and Firat, O. The devil is in the errors: Leveraging
 large language models for fine-grained machine transla-
 tion evaluation. In *Proceedings of the Eighth Conference*
on Machine Translation, pp. 1066–1083. Association for
 Computational Linguistics, 2023.
- Fomicheva, M., Sun, S., Yankovskaya, L., Blain, F.,
 Guzmán, F., Fishel, M., Aletras, N., Chaudhary, V., and
 Specia, L. Unsupervised quality estimation for neural
 machine translation. In *Transactions of the Association*

- 330 for *Computational Linguistics*, volume 8, pp. 539–555,
 331 2020.
- 332 Fu, J., Ng, S.-K., Jiang, Z., and Liu, P. GPTScore: Evaluate
 333 as you desire. *arXiv preprint arXiv:2302.04166*, 2023.
- 334 Koo, R., Lee, M., Raheja, V., Park, J. I., Kim, Z. M.,
 335 and Kang, D. Benchmarking cognitive biases in
 336 large language models as evaluators. *arXiv preprint*
 337 *arXiv:2309.17012*, 2023.
- 338 Kuhn, L., Gal, Y., and Farquhar, S. Semantic uncertainty:
 339 Linguistic invariances for uncertainty estimation in nat-
 340 ural language generation. In *The Eleventh International*
 341 *Conference on Learning Representations*, 2023.
- 342 Kumar, B., Lu, C., Gupta, G., Palepu, A., Bellamy, D.,
 343 Raskar, R., and Beam, A. Conformal prediction with large
 344 language models for multi-choice question answering.
 345 *arXiv preprint arXiv:2305.18404*, 2023.
- 346 Liu, Y., Iyer, D., Xu, Y., Wang, S., Xu, R., and Zhu, C. G-
 347 Eval: NLG evaluation using GPT-4 with better human
 348 alignment. In *Proceedings of the 2023 Conference on*
 349 *Empirical Methods in Natural Language Processing*, pp.
 350 2511–2522. Association for Computational Linguistics,
 351 2023.
- 352 Maltoudoglou, L., Paisios, A., and Sakkas, H. BERT-based
 353 conformal predictor for intent classification. In *Proceed-*
 354 *ings of the Ninth Symposium on Conformal and Probabilistic*
 355 *Prediction and Applications*, pp. 178–193, 2020.
- 356 Moon, J. W. *Topics on Tournaments*. Holt, Rinehart and
 357 Winston, 1968.
- 358 Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman,
 359 A. Inductive confidence machines for regression. In
 360 *European Conference on Machine Learning*, pp. 345–356.
 361 Springer, 2002.
- 362 Qin, Z., Jagerman, R., Hui, K., Zhuang, H., Wu, J., Shen,
 363 L., Liu, T., Liu, J., Metzler, D., Wang, X., and Bendersky,
 364 M. Large language models are effective text rankers
 365 with pairwise ranking prompting. In *Findings of the*
 366 *Association for Computational Linguistics: NAACL 2024*,
 367 2024.
- 368 Quach, V., Fisch, A., Schuster, T., Yala, A., Sohn, J. H.,
 369 Jaakkola, T. S., and Barzilay, R. Conformal language
 370 modeling. In *The Twelfth International Conference on*
 371 *Learning Representations*, 2024.
- 372 Saito, K., Sugawara, S., and Inui, K. Verbosity bias in
 373 preference labeling by large language models. *arXiv*
 374 *preprint arXiv:2310.10076*, 2023.
- Schulze, M. A new monotonic, clone-independent, reversal
 symmetric, and condorcet-consistent single-winner elec-
 tion method. *Social Choice and Welfare*, 36(2):267–303,
 2011.
- Tibshirani, R. J., Barber, R. F., Candès, E., and Ramdas, A.
 Conformal prediction under covariate shift. *Advances in*
Neural Information Processing Systems, 32, 2019.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic*
Learning in a Random World. Springer, 2005.
- Wang, C., Yang, Y., Dang, C., and Che, W. Large language
 models are not yet human-level evaluators for abstractive
 summarization. In *Findings of the Association for Com-*
putational Linguistics: EMNLP 2023, pp. 4215–4233.
 Association for Computational Linguistics, 2023.
- Ye, S., Kim, D., Jang, S., Shin, H., Baek, Y., Song, J., Park,
 D., and Seo, M. FLASK: Fine-grained language model
 evaluation based on alignment skill sets. In *Proceedings*
of the 2024 Conference of the North American Chapter
of the Association for Computational Linguistics, 2024.
- Young, H. P. Condorcet’s theory of voting. *American*
Political Science Review, 82(4):1231–1244, 1988.
- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z.,
 Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P., Zhang,
 H., Gonzalez, J. E., and Stoica, I. Judging LLM-as-a-
 Judge with MT-Bench and Chatbot Arena. In *Advances*
in Neural Information Processing Systems, volume 36,
 2023.

A. Prompt Templates

Pairwise Preference Prompt

You are evaluating the {criterion} of two summaries of the following article.

Article: {document}

Summary A: {system.a}

Summary B: {system.b}

Which summary is better in terms of {criterion}? Answer with exactly 'A' or 'B'. No explanation.

Direct Scoring Prompt

You are evaluating the {criterion} of the following summary of an article.

Article: {document}

Summary: {system.output}

Rate the {criterion} on a scale of 1--5 where:

1 = Very Poor, 2 = Poor, 3 = Fair, 4 = Good, 5 = Excellent

Respond with a single integer between 1 and 5. No explanation needed.

B. Full Conformal Results

Table 5 reports conformal coverage and average set size at all four α levels tested. Every entry exceeds the $1-\alpha$ target, confirming the theoretical guarantee holds robustly across all operating points.

Table 5. Full conformal prediction results across all α levels. Mean coverage and set size across 20 random splits. All coverages meet or exceed the $1-\alpha$ target (bolded).

Judge	Criterion	$\alpha = 0.05$		$\alpha = 0.10$		$\alpha = 0.15$		$\alpha = 0.20$	
		Cov.	Size	Cov.	Size	Cov.	Size	Cov.	Size
LLAMA-3.1-70B	Coherence	.993	4.29	.947	3.04	.945	2.98	.945	2.98
	Consistency	.992	4.97	.992	4.97	.889	4.26	.889	4.26
	Fluency	.971	4.95	.967	4.93	.967	4.93	.967	4.93
	Relevance	.964	2.98	.961	2.82	.961	2.82	.961	2.82
QWEN-2.5-72B	Coherence	1.000	4.55	.934	3.62	.905	2.99	.905	2.99
	Consistency	.963	4.53	.962	4.48	.962	4.48	.962	4.48
	Fluency	.967	4.76	.967	4.76	.967	4.76	.967	4.76
	Relevance	.975	3.63	.961	2.99	.961	2.99	.961	2.99
MISTRAL-SMALL	Coherence	.980	4.83	.950	4.47	.950	4.47	.901	4.10
	Consistency	.996	4.98	.996	4.98	.996	4.98	.970	4.94
	Fluency	1.000	4.99	1.000	4.99	1.000	4.99	1.000	4.99
	Relevance	.995	4.15	.935	3.06	.931	2.93	.931	2.93
GPT-4O-MINI	Coherence	.979	4.59	.971	4.51	.899	3.46	.876	2.98
	Consistency	.996	4.99	.942	4.76	.899	4.54	.899	4.54
	Fluency	.995	4.99	.995	4.99	.995	4.99	.984	4.97
	Relevance	1.000	4.37	.928	3.17	.919	2.96	.919	2.96

C. Violation Rates Across All Criteria

Table 6 reports aggregate and per-document transitivity violation rates for all four judges across all four criteria. *Fluency* consistently exhibits the highest fraction of documents with at least one violation.

Judge	Criterion	Agg. $\bar{\rho}$	% docs ≥ 1
LLAMA-3.1-70B	Coherence	0.008	33.3%
	Consistency	0.012	36.7%
	Fluency	0.033	66.7%
	Relevance	0.011	36.7%
QWEN-2.5-72B	Coherence	0.022	50.0%
	Consistency	0.025	56.7%
	Fluency	0.024	60.0%
	Relevance	0.023	50.0%
MISTRAL-SMALL	Coherence	0.041	50.0%
	Consistency	0.036	63.3%
	Fluency	0.034	53.3%
	Relevance	0.029	53.3%
GPT-4O-MINI	Coherence	0.014	46.7%
	Consistency	0.013	43.3%
	Fluency	0.026	53.3%
	Relevance	0.020	53.3%

Table 6. Aggregate and per-document violation rates across all four criteria. *Fluency* consistently shows the highest fraction of documents with ≥ 1 violation, reaching 66.7% for LLAMA-3.1-70B.

D. Coverage vs. α Curves

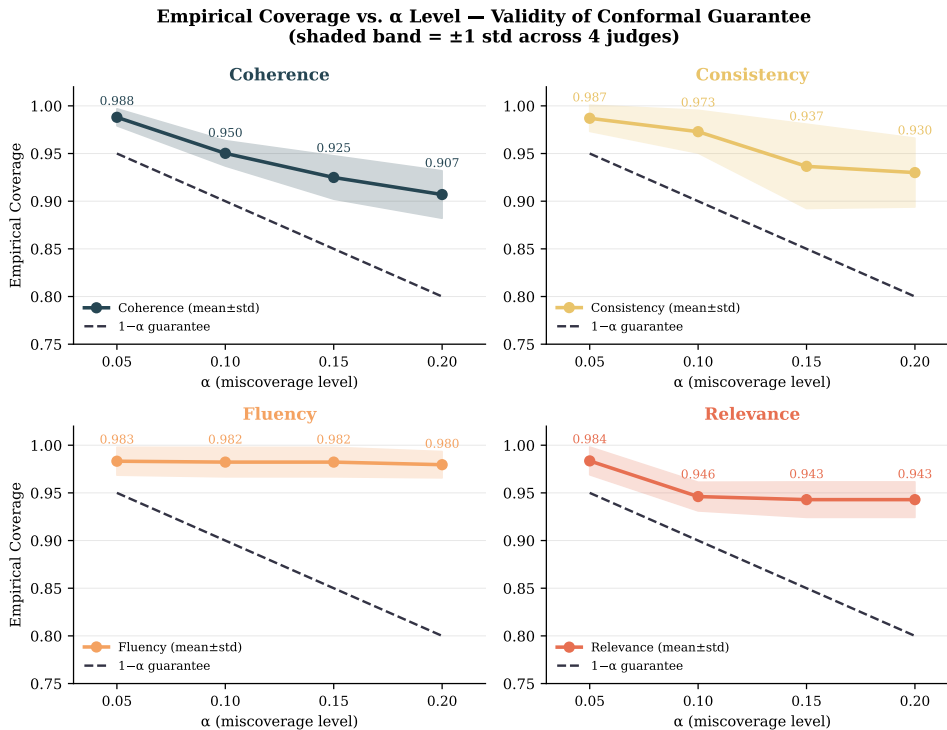


Figure 6. Empirical coverage vs. α . Shaded bands: ± 1 std across the four judges. Dashed line: theoretical guarantee $1-\alpha$. Coverage meets or exceeds the guarantee at every operating point for all four criteria.