

# Who Judges the Judge? LLM Jury-on-Demand: Building Trustworthy LLM Evaluation Systems

Xiaochuan Li\*, Ke Wang\*, Girija Gouda,  
Shubham Choudhary, Yaqun Wang, Linwei Hu,  
Joel Vaughan, Freddy Lecue

Wells Fargo Bank, N.A., USA

December 2, 2025

## Abstract

As Large Language Models (LLMs) become integrated into high-stakes domains, there is a growing need for evaluation methods that are both scalable for real-time deployment and reliable for critical decision-making. While human evaluation is reliable, it is slow and costly. Single LLM judges are biased, and static juries lack adaptability. To overcome these limitations, we propose LLM Jury-on-Demand - a dynamic, learning-based framework for scalable and context-aware evaluation. Our method trains a set of reliability predictors to assess when LLM judges will agree with human experts, leveraging token distributions, embeddings, and structural input features. This enables a fully adaptive evaluation where, for each data point, an optimal jury of the most reliable judges is dynamically selected and their scores are aggregated using their reliability as weights. Experiments on summarization and RAG benchmarks show that our dynamic jury system achieves significantly higher correlation with human judgment than both single-judge and static-jury baselines. These results highlight the promise of adaptive, learning-based juries for building scalable, more reliable and trustworthy evaluation systems for modern LLMs in high-stakes domains.<sup>1</sup>

## 1 Introduction

Large Language Models (LLMs) such as the GPT series, Llama, and Gemini have demonstrated transformative capabilities, leading to their rapid integration into critical, real-world applications (Brown et al., 2020; Touvron et al.,

---

\*These authors contributed equally to this work.

<sup>1</sup>The views expressed in this paper are solely those of the authors and do not necessarily reflect the views of their affiliated institutions.

2023; Team et al., 2023). As these models are deployed in high-stakes domains, ensuring their outputs are reliable, safe, and aligned with human expectations has become a paramount concern (Shukla, 2025; Wang et al., 2023). The gold standard for assessing these qualities would be human evaluation, where experts provide nuanced judgments. However, this process is notoriously slow, expensive, and difficult to scale, making it impractical for the rapid development cycles of modern AI (Calderon et al., 2025). To overcome this scalability bottleneck, the field historically relied on reference-based automated metrics like BLEU and ROUGE, which measure lexical overlap between the generated output and a ground-truth reference text (Papineni et al., 2002; Lin, 2004). These methods are now widely considered insufficient for capturing multifaceted attributes like completeness, relevance, or groundedness in the sophisticated outputs of modern generative models (Zhang et al., 2019; Cao et al., 2025).

To address this evaluation gap, researchers have increasingly adopted the LLM-as-a-Judge paradigm, which leverages powerful language models like GPT-4 to serve as scalable, automated evaluators (Zheng et al., 2023; Li et al., 2024b; Gu et al., 2024). While promising, this approach introduces a critical trade-off where the scalability of a single LLM judge comes at the cost of reliability. The papers Schroeder & Wood-Doughty (2024), Li et al. (2024a) and Baumann et al. (2025) contain substantial evidence showing that single judges can be prone to systematic biases and inconsistencies, limiting their trustworthiness. A logical evolution has been to employ a “jury” of multiple LLMs to improve robustness (Feng et al., 2025; Verga et al., 2024). However, these jury systems typically rely on static aggregation methods, such as simple averaging. This fails to address a more fundamental issue, as a judge’s expertise is context-dependent and its reliability can change dramatically based on the text being evaluated. This leaves a critical gap for a truly adaptive evaluation system.

In this paper, we introduce LLM Jury-on-Demand, a novel framework that bridges this gap by creating a dynamic, learning-based evaluation system. Our work moves beyond static juries by training a system to predict the reliability of each potential judge based on a rich set of features extracted from the text. This allows our framework to perform a fully adaptive evaluation where, for each data point, an optimal jury of the most reliable judges is dynamically selected, and their scores are aggregated using their reliability as weights. Our main contributions are threefold:

- A new framework for adaptive LLM evaluation that demonstrates superior correlation with human judgment compared to single-judge and static-jury baselines.
- A method to predict LLM judge reliability at the instance level using text-based features.
- Extensive experiments and analyses across multiple tasks and datasets to validate the effectiveness of the proposed approach.

## 2 Related Work

The evaluation of large language models is a rapidly evolving field, as captured in recent surveys mapping the transition from static benchmarks to more dynamic and automated evaluation frameworks (Cao et al., 2025). Our work builds upon three key research areas: the LLM-as-a-Judge paradigm, the evolution from single judges to multi-model juries, and the broader concept of LLM performance prediction.

The LLM-as-a-Judge approach has become a scalable alternative to human annotation (Zheng et al., 2023), with surveys documenting its widespread application and promising correlation with human preferences (Li et al., 2024b; Gu et al., 2024). However, this paradigm has significant limitations. LLM judges exhibit biases, such as a preference for longer answers and sensitivity to the order in which responses are presented (Schroeder & Wood-Doughty, 2024), and their judgments can be skewed by their own intrinsic style or pre-training data, which compromises the fairness and reliability of the evaluation (Li et al., 2024a). These challenges motivate the need for more robust frameworks that can mitigate the inherent biases of a single judge.

To address these limitations, a growing body of work has explored using a “jury” of multiple LLMs, based on the insight that collaboration among diverse models can lead to more stable and reliable assessments (Feng et al., 2025). Initial work shows that simple ensembles, such as averaging the scores from a panel of smaller models, can outperform a single, larger model at a lower cost (Verga et al., 2024; Rahmani et al., 2024). More advanced methods have explored multi-agent frameworks where judges engage in peer-review or debate-like discussions to arrive at a consensus (Chu et al., 2024; Zhao et al., 2024). While a significant step forward, they typically rely on either static aggregation methods like simple voting or averaging or require complex and often unscalable conversational interactions. They do not account for the fact that a judge’s expertise varies across different contexts, leaving a critical gap for systems that can adapt the jury’s composition and weight to the specific context of the text being evaluated.

Our work is also grounded in LLM performance prediction. Studies have shown that it is possible to train a model to predict an LLM’s performance on a given task by using features derived from the model and the task itself (Ye et al., 2023). Some approaches have even trained “assessor” models to predict when another model is likely to answer a question correctly, a concept that parallels our goal of predicting reliability (Schellaert et al., 2025). While these works validate the fundamental premise that LLM performance has learnable patterns, they typically focus on predicting general, task-level success rather than the instance-level reliability of an LLM acting as an evaluator. Our framework innovates by applying this concept to jury-based evaluation, enabling the dynamic selection and weighting of judges on a per-instance basis.

### 3 Methodology

In this section, we detail the architecture and components of LLM Jury-on-Demand framework. Our framework is designed to produce more reliable automated evaluations by shifting from a static to a dynamic, learning-based process. The central hypothesis of our work is that an LLM judge's reliability is not fixed, but varies based on the specific characteristics of the text it evaluates. Our system models this variance by learning to predict when each judge is likely to agree with human experts.

#### 3.1 Framework Overview

The LLM Jury-on-Demand framework operates through a multi-stage pipeline, as shown in Fig. 1. The process begins with a set of input texts related to the evaluation task. Crucially, the set of texts we analyze depends on the task itself. For summarization tasks, the system analyzes both the original source text and the model-generated output summary. For Retrieval-Augmented Generation (RAG) tasks, it analyzes the source text, the retrieved context, and the final generated answer. This context-rich approach allows the system to capture a more complete picture of the evaluation challenge.

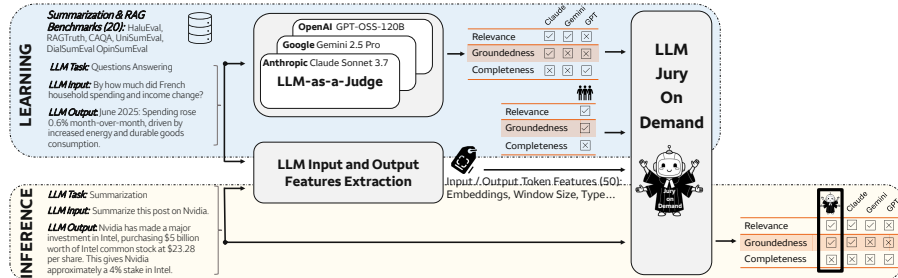


Figure 1: Overview of the LLM Jury-on-Demand inference pipeline. The system extracts features from input texts to predict judge reliability, dynamically assembles a jury of the top  $K$  most reliable judges for each instance, and calculates a final weighted score.

These texts are first processed by a feature extraction module, which extracts a wide range of textual and semantic signals. The resulting feature vector for each instance serves as input to a suite of pre-trained reliability prediction models. In the final stage, the system leverages these reliability predictions to perform a fully adaptive, per-instance evaluation. For each data point, a jury of a pre-tuned size  $K$  is dynamically assembled by selecting the judges with the highest predicted reliability for that specific instance. The final score is then computed as a weighted average of these selected judges' raw scores, using their reliability predictions as the weights.

### 3.2 Feature Engineering for Reliability Prediction

The foundation of our system is its ability to represent the evaluation context through a rich set of predictive features. We hypothesize that signals related to a text's size, complexity, and semantic content can reveal the scenarios in which different LLM judges excel or struggle. The features are extracted from all available texts (source, context, and output, as applicable to the task) and concatenated into a single feature vector for each data point. Many of these features are computed using Natural Language Toolkit (NLTK) (Bird et al., 2009). A complete list of features is in the Appendix A. Key feature categories are:

**Text Size Features:** These include basic structural metrics such as word count, sentence count, paragraph count, and the compression ratio between the source and generated text.

**Special Words Features:** These count the occurrences of specific word types that can indicate the style or complexity of the text. Examples include counts of difficult words (words that have more than two syllables), named entities and modality verbs (e.g. “could”, “should”) etc.

**Text Complexity Features:** These quantify readability and ambiguity using established linguistic formulas. Examples include the Flesch reading ease index (Kincaid et al., 1975), lexical diversity (the variety of words used), and other measures of syntactic and semantic ambiguity.

**Embedding-Related Features:** Embeddings encode text into a dense vector representation (Mikolov et al., 2013; Pennington et al., 2014; Devlin et al., 2019). These capture the semantic meaning and topic of the text. Top 10 PCA components (Jolliffe, 2011) of text embeddings are used as features. Additionally, we compute cosine similarity between each text component’s embedding and a set of predefined topic embeddings (e.g. finance, technology), using these similarity scores as additional topical relevance features.

### 3.3 Learning to Predict Judge Reliability

Our framework learns judge reliability by training a dedicated machine learning model for each specific evaluation context. That is, for each potential judge, for each task (e.g., summarization), and for each evaluation metric (e.g., completeness), we train a distinct model. The purpose of this model is not to predict the evaluation score itself, but rather to predict the probability that the corresponding judge will be reliable on that metric for a given data point.

We frame this as a binary classification task to predict whether a judge's score will be “good” or “bad”. To generate the ground-truth labels for training, we compare each judge's score against a gold-standard human expert score. First, we apply min-max normalization to all human and model scores to scale them to a  $[0, 1]$  range. A judge's evaluation is then labeled as “good” (1) if its normalized score falls within a predefined tolerance hyperparameter,  $\tau$ , of the normalized human score. Otherwise, it is labeled as “bad” (0).

For each classification model, we use XGBoost, a gradient-boosted tree algorithm known for its strong performance on tabular data (Chen & Guestrin, 2016). This approach is grounded in the broader research area of LLM performance prediction, which has shown that model performance can be learned from features (Ye et al., 2023; Schellaert et al., 2025). At inference time, our trained models output a probability score between 0 and 1, which we use as the predicted reliability.

### 3.4 Assembling and Scoring the Jury

The core of our framework is its ability to assemble an expert jury and use its members' dynamically weighted opinions to compute a final score. This process, which is applied for each individual data point, involves two key mechanics: a reliability-based jury selection algorithm and an instance-level dynamic weighting scheme.

**Jury Selection.** For each data point, we first use the suite of pre-trained reliability models (described in Sec. 3.3) to generate a reliability score for each of the  $N$  judges in our pool. A jury of a pre-tuned size  $K$  is then formed by simply selecting the  $K$  judges with the highest predicted reliability scores for that specific instance. This approach allows the jury's composition to be completely dynamic, adapting to the unique characteristics of each text.

**Dynamic Score Aggregation.** Once the instance-specific jury is selected, the final evaluation score is calculated as a weighted average of the raw scores from those  $K$  jury members. The weights are derived directly from the instance-specific reliability scores,  $[r_1, r_2 \dots, r_K]$ . Specifically, the weight for judge  $i$  in the jury is calculated as  $w_i = r_i / (\sum_{j=1}^K r_j)$ . This dynamic, per-instance selection and weighting process allows our system to prioritize the opinions of the most trustworthy judges for any given text, which stands in contrast to prior systems that rely on static juries and static aggregation methods like simple averaging (Verga et al., 2024).

These two mechanics are the building blocks for our system's training and inference pipelines.

### 3.5 System Training and Inference

Our framework involves a one-time training and tuning phase to establish an optimal configuration, which is then used in a repeatable inference pipeline to evaluate new data points.

**Training and Hyperparameter Tuning.** The goal of the training phase is to find the optimal hyperparameters that will generalize best. This involves finding the single best jury size  $K$  and the optimal tolerance level  $\tau$  for each individual judge's reliability model. To do this, we first train a large pool of reliability predictor models on our training data, covering all possible combinations of judges and tolerance values. We then conduct a search over the hyperparameter space of possible jury sizes ( $K$ ) and per-judge tolerance configurations. Each configuration is evaluated on a held-out validation set. For every

data point in the validation set, we apply the *Jury Selection* and *Dynamic Score Aggregation* mechanics described in Sec. 3.4. The configuration that yields the highest Kendall's Tau correlation with human scores across the validation set is selected as the optimal configuration for the final system.

**Inference Pipeline.** With the optimal configuration locked in (i.e., a fixed  $K$  and a set of optimal reliability models), the system is ready to evaluate a new, unseen data point. For a new instance, the system first uses the optimal reliability models to predict an instance-specific reliability score for every potential judge in the pool. It then applies the *Jury Selection* and *Dynamic Score Aggregation* mechanisms to select the top  $K$  judges and compute the final, weighted score.

## 4 Experimental Setup

To validate the effectiveness of our LLM Jury-on-Demand framework, we conducted a series of experiments designed to measure its performance against standard evaluation methods. This section details the datasets, evaluation protocol, and implementation specifics of our experiments.

### 4.1 Datasets and Tasks

Our evaluation spans two challenging natural language generation tasks: summarization and retrieval-augmented generation (RAG). For each task, we focus on evaluating three core metrics: groundedness, relevance, and completeness. These metrics are essential for assessing the quality and trustworthiness of generated text. Detailed definitions for each metric are provided in Appendix B.

To train our jury framework, we chose datasets with human annotations for these dimensions. We reviewed a diverse set of datasets for training, selecting 3-4 datasets per task by metric dimension to ensure coverage across various domains. To prevent any single dataset or any single metrics category from dominating a particular evaluation task, we applied stratified down-sampling where necessary. The details are provided in Appendix C.

### 4.2 Evaluation Protocol

Our experimental protocol is designed to ensure a fair and rigorous comparison between our proposed system and relevant baselines.

**Evaluation Metric.** The primary metric for our experiments is the Kendall's Tau correlation coefficient (KENDALL, 1938). This non-parametric statistic measures the ordinal association between two sets of rankings. In our context, it quantifies how well a system's evaluation scores align with the rankings provided by human experts. A Kendall's Tau value close to 1 indicates strong agreement with human judgment.

**Judge Prompting Protocol.** To generate the raw scores for our experiment, each potential judge model was prompted with a carefully structured

template. This template includes a system prompt to set the judge's persona (“You are a helpful, respectful and honest assistant”) and a user prompt that defines the task, the specific metric (e.g., Completeness), and the required scoring format. The scoring scale was adapted to the task: for summarization, all judges were instructed to provide a single integer score from 1 (lowest quality) to 5 (highest quality); for RAG, a scale of 1 (lowest quality) to 3 (highest quality) was used. The full prompt structure is provided in the Appendix D. Additionally, we conducted an ablation study to analyse the effect on the system’s resilience to slight prompt variations as described in the Appendix H.4.

**Baselines.** To benchmark our system's performance, we established a judge pool consisting of 10 diverse LLMs. This pool serves as the foundation for all evaluation methods compared in our study and includes the following models: Claude 3.7 SONNET (Cla), Gemini 2.5 Pro (Comanici et al., 2025), Gemini 2.5 Flash, Gemini 2.0 Flash, GPT-OSS-20B (Agarwal et al., 2025), GPT-OSS-120B, Gemma 3-27B-IT (Team et al., 2025), Phi-4-reasoning (Abdin et al., 2025), LLAMA-3.2-3B-Instruct (Grattafiori et al., 2024), and DeepSeek-R1 (Guo et al., 2025). From this pool, we formed two categories of baselines:

1. **Single-Judge Baselines:** The performance of each of the 10 judges when used as a standalone evaluator.
2. **Static-Jury Baselines:** We compare against four distinct static jury formulations to rigorously test the benefits of dynamic selection:
  - **Static Jury (Average-All):** The performance of a non-adaptive, naive jury that uses all 10 judges in the pool. For each data point, the final score for this baseline is the simple average of the raw scores from all 10 judges.
  - **Static Jury (Average-Top-K):** This baseline identifies the Top- $K$  best-performing single judges based on their Kendall’s Tau on the validation set. The final score is the simple average of these Top- $K$  judges. The value of  $K$  is tuned on the validation set.
  - **Static Jury (Weighted-Regression):** A regression-based jury using all 10 judges. We train a linear regression model without intercept using human annotation scores as labels and single judge scores as features on the training set.
  - **Static Jury (Weighted-Tau):** A performance-weighted average of all 10 judges. The weights are derived from each judge’s validation Kendall’s Tau, normalized using a softmax function.

### 4.3 Implementation Details

This section provides the specific procedures used to configure and train our system for the experiments.

First, we prepared the data for each task and metric by combining all relevant source datasets. For example, to evaluate the completeness metric for



the summarization task, we aggregated the data from SummEval, TLDR, and UniSumEval into a single, combined dataset. We then performed a global 60-20-20 split on this combined data to create our final training, validation, and holdout test sets, ensuring data from all sources were represented in each split. As described in Sec. 3.3, all human and model scores were then normalized to a  $[0, 1]$  range to ensure consistency.

Next, we determined the system’s optimal configuration through a comprehensive hyperparameter tuning process on the combined validation set, as outlined in Sec. 3.5. We defined a search space for three key hyperparameter categories: the jury size  $K$  (ranging from 2 to 9), the per-judge tolerance values  $\tau$  used for training the reliability predictors, and the internal parameters of the XGBoost models. We exclude  $K = 1$  as it reduces the jury to a single judge selector, effectively duplicating the Single-Judge baseline paradigm. We also exclude  $K = 10$  (the full pool), as this configuration represents a static ensemble identical to the Average-All baseline, negating the benefit of dynamic selection. The sets of tolerance values were chosen to reflect the varying scales of the original human annotation scores across the different source datasets. For each candidate configuration, we evaluated its performance by applying the full per-instance jury selection and weighting pipeline (Sec. 3.4) to every data point in the combined validation set. The final optimal configuration was chosen based on which set of hyperparameters yielded the highest Kendall’s Tau correlation on this combined validation set. This finalized configuration was then used for the final, unbiased evaluation on the locked-away test set.

Additionally, to assess robustness beyond validation-based tuning, we conduct ablation studies on jury size ( $K$ ) and tolerance ( $\tau$ ), detailed in Appendix H.2 and H.3. Jury size is varied from 1 to 9 across two representative tasks, Summarization-Completeness and RAG-Groundedness, over 10 independent runs. As shown in Fig. 17, performance for both tasks follows a similar pattern: accuracy improves as jury size increases, reaches an optimal range (around  $K = 5-8$ ), and then declines slightly at very large sizes, indicating diminishing returns beyond the peak.

## 5 Results and Analysis

This section evaluates our LLM Jury-on-Demand framework against the single-judge and static-jury baselines. To ensure robustness, all experiments were repeated 10 times with different random seeds for data partitioning, and we report the mean and standard deviation of the results. Our analysis presents the main performance, both overall and at a granular dataset-level, followed by a diagnostic analysis of our system’s internal mechanics, including feature importance and judge selection frequency.

## 5.1 Performance Analysis

We assess performance by comparing the Kendall’s Tau correlation of each method with human judgment on the test sets for each task-metric combination. The distribution of results from our 10 independent runs is summarized in Fig. 2.

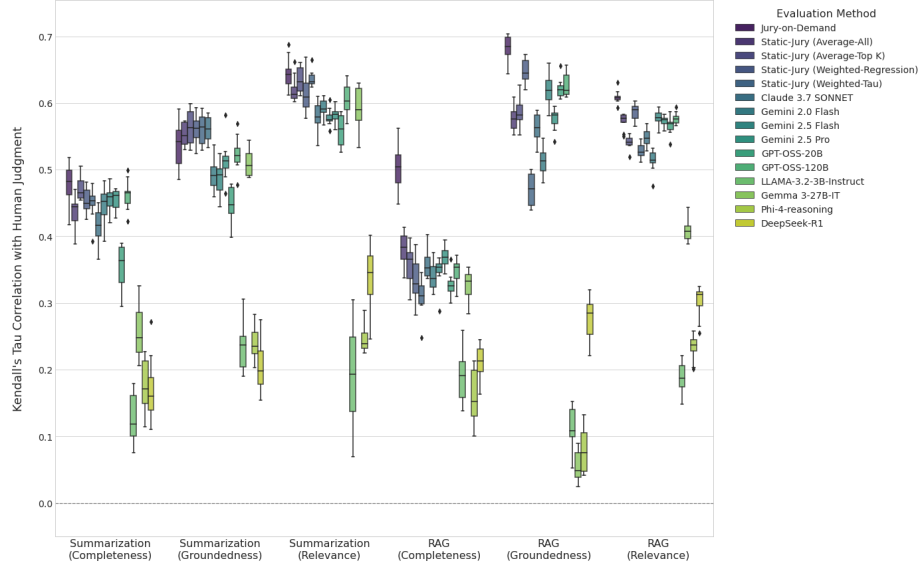


Figure 2: Overall performance comparison over 10 runs. Boxplot of Kendall’s Tau correlation between each evaluation method’s scores and human judgments, aggregated across all datasets for the 6 task-metric combinations. Our Jury-on-Demand system achieves the highest median correlation in nearly all categories and shows the most robust performance.

The results clearly demonstrate that our Jury-on-Demand framework consistently outperforms all baselines across every task and metric. In all six categories, our method achieves the highest mean correlation with human judgment, validating its effectiveness. For instance, in the challenging RAG-Groundedness task, our system achieves a mean Kendall’s Tau of  $0.68 (\pm 0.02)$ . This represents a significant improvement not only over the simple Static Jury (Average-All) ( $0.58 \pm 0.02$ ) but also over the stronger optimized baselines, such as Static Jury (Weighted-Regression) ( $0.65 \pm 0.02$ ) and the strongest single judge for that task, GPT-OSS-120B ( $0.63 \pm 0.02$ ).

To provide a more detailed view of performance, we now analyze the results for the RAG-Completeness task at the individual dataset level. The granular results in Table 1 confirm the findings from the overall analysis: our Jury-on-Demand system outperforms all four static jury baselines and the best-performing single judge on every individual dataset for this task. The performance lift is particularly pronounced on the ASQA dataset, indicating that our

system is robust and its advantages are not an artifact of data aggregation but hold true at a more granular level. The full breakdown of results for all tasks and datasets is available in Appendix E.1. We also report statistical significance and effect sizes. Specifically, we perform one-sided Wilcoxon signed-rank tests (Wilcoxon, 1945) and compute Cliff’s delta (Cliff, 1993), a non-parametric effect size metric that quantifies the difference between two groups. The results are presented in Appendix E.1.

Finally, the results also highlight the inherent unreliability of relying on a single LLM as an evaluator. As shown in Fig. 2, the performance of single judges is highly variable. The “best” single judge changes from one task to another; for example, Claude 3.7 SONNET is often the strongest single judge for Summarization-Groundedness, but it is one of the weakest for Summarization-Completeness. This instability proves that there is no single “best” LLM judge, making a static choice of evaluator a risky and unreliable strategy.

The comparison between our Jury-on-Demand and the various static jury baselines isolates the benefit of our core contribution. While baselines like Static Jury (Weighted-Regression) are competitive, our dynamic system is consistently superior. This demonstrates that the primary performance gain comes not just from using a jury, but from the ability to dynamically select and weight its members based on the context of the input. The stability of this outperformance, demonstrated across 10 runs, underscores the reliability of our dynamic approach.

Table 1: Granular Performance on RAG-Completeness. Mean Kendall’s Tau correlation ( $\pm$  std. dev.) on the individual test sets across 10 runs. Our Jury-on-Demand system consistently outperforms all static baselines and the best single judge.

Dataset	Jury-on-Demand	Static (Avg-All)	Static (Avg-TopK)	Static (W-Reg)	Static (W-Tau)	Best Single
ALCE	$0.47 \pm 0.07$	$0.38 \pm 0.09$	$0.28 \pm 0.08$	$0.34 \pm 0.11$	$0.23 \pm 0.10$	$0.40 \pm 0.07$ (GPT-OSS-20B)
ASQA	$0.54 \pm 0.05$	$0.38 \pm 0.05$	$0.38 \pm 0.04$	$0.36 \pm 0.04$	$0.34 \pm 0.03$	$0.42 \pm 0.03$ (Claude 3.7 SONNET)
QASPER	$0.44 \pm 0.08$	$0.41 \pm 0.08$	$0.27 \pm 0.08$	$0.35 \pm 0.07$	$0.24 \pm 0.11$	$0.43 \pm 0.07$ (GPT-OSS-120B)

## 5.2 Analyzing the Interaction Between Judges, Tasks, and Data Attributes

We begin by analyzing judge selection patterns within juries across different tasks and datasets. Fig. 3 summarizes the selection frequency of each judge for the RAG groundedness and RAG completeness tasks. The results reveal distinct preferences: Claude 3.7 Sonnet and DeepSeek R1 are frequently selected

for completeness evaluation but are rarely chosen for groundedness. In contrast, Gemini 2.5 Flash is commonly selected for groundedness but appears less frequently in completeness evaluations. GPT OSS 20B and GPT OSS 120B are consistently selected across both metrics. A comprehensive comparison of judge selection across all tasks is in Appendix E.2.

We now examine how data properties impact judge performance. For illustration, we focus on two tasks: RAG groundedness and summarization completeness. The analysis for summarization completeness is in Appendix E.2. We select properties that rank among the most important features in the XGBoost model. For the summarization task, the key property is the compression ratio (i.e., the length of the summary divided by the length of the article). For the RAG task, the selected property is the character count of each response.

To better illustrate the findings, we focus on three judges for each task. We begin with the RAG groundedness task. Fig. 4 shows model performance across bins of low, medium, and high response character counts. Two main observations emerge:

1. All judges perform worse as the character count increases.
2. Gemini 2.0 Flash performs comparably to the others when the character count is low, but its performance drops significantly at higher character counts, especially compared to Gemini 2.5 Flash and GPT OSS 20B.

To explain these trends, we examine the distribution of annotation scores in the low and high character count regions (see Fig. 5). In the high character count region, many responses receive a score of 1 (moderately ungrounded), whereas in the low character count region, most scores are either 0 (severely ungrounded) or 2 (fully grounded). It is easier for judges to distinguish between scores 0 and 2, but more difficult to differentiate between 1 and 2. Weaker judges, such as Gemini 2.0 Flash, particularly struggle with identifying ungrounded content. Fig. 6 and Fig. 7 present the score confusion matrices for Gemini 2.0 Flash and Gemini 2.5 Flash for long and short responses. Compared to Gemini 2.5 Flash, Gemini 2.0 Flash assigns a disproportionately high number of score 2s, which explains its poor performance in the high character count region and relatively good performance in the low character count region. We also reviewed specific examples to understand why Gemini 2.0 Flash assigns score 2 even when the content is ungrounded. This analysis is provided in the Appendix E.3.

Finally, the jury selection dynamically aligns with model performance, as shown in Fig. 4. This is most evident in the high character count bin, where the performance of Gemini 2.0 Flash drops significantly. Correspondingly, its selection percentage in our dynamic jury plummets to its lowest point, demonstrating that the reliability predictors correctly identify and avoid this weaker judge when it is unreliable. Conversely, Gemini 2.5 Flash maintains the highest performance in this high-count bin, and our system selects it for the jury in the vast majority of cases. GPT-OSS-20B also shows strong alignment in the medium character count bin, where it achieves its highest performance and is also the most selected judge.

The framework also demonstrates robustness beyond simply picking the single best-performing judge. For instance, in the low character count bin, GPT-OSS-20B has the highest kendall’s Tau. While it is selected frequently, Gemini 2.5 Flash, which also performs exceptionally well, is selected more often. This illustrates that the system does not rely on a single judge; rather, it identifies a pool of highly reliable judges for a given context and assembles an optimal jury from that pool.

These findings reinforce the importance of constructing dynamic juries that adapt to specific data characteristics and demonstrate the potential of our framework to predict judge reliability based on interpretable data properties.

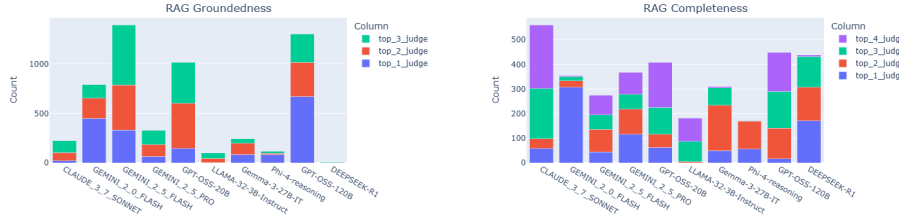


Figure 3: Selection frequency of the judge in the jury. Top k judge means that the judge has the k-th highest reliability score in the jury. Claude 3.7 Sonnet and DeepSeek R1 are favored in completeness, while Gemini 2.5 Flash is more often selected for groundedness.



Figure 4: RAG Groundedness analysis by response character count. (Left) Kendall’s Tau correlation for three single judges across low, medium, and high response character count bins. (Right) The selection percentage of these judges in the final dynamic jury for data points within each bin. The analysis shows that judge performance degrades with longer responses, particularly for Gemini 2.0 Flash. Our system’s jury selection adapts to this, heavily favoring the more reliable Gemini 2.5 Flash in the high-count bin.

Table 2 presents the top five most important features for the judge reliability XGBoost model, as determined by permutation feature importance (Fisher et al., 2019), with the summarization groundedness task as an example. For illustration, we focus on Gemini 2.5 Pro and GPT-OSS-120B, while the complete

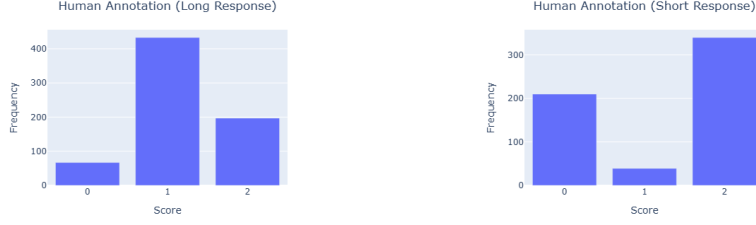


Figure 5: RAG Groundedness: Distribution of annotation scores across response lengths. Longer responses tend to receive more score 1s (moderately ungrounded), while shorter responses are more often assigned scores of 0 (severely ungrounded) or 2 (fully grounded).

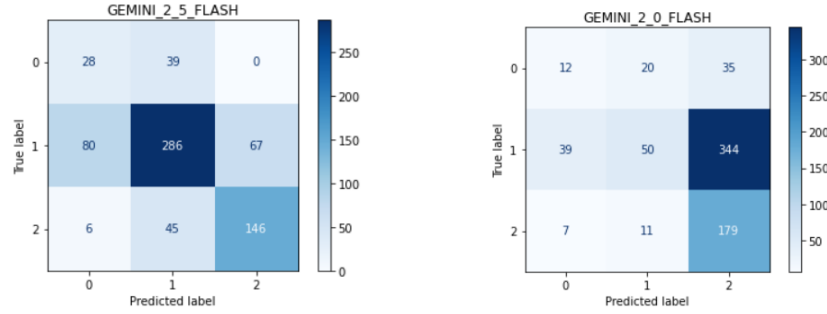


Figure 6: RAG Groundedness: Confusion matrix for long responses. Gemini 2.0 Flash frequently assigns score 2 to ungrounded content, leading to reduced performance in this region.

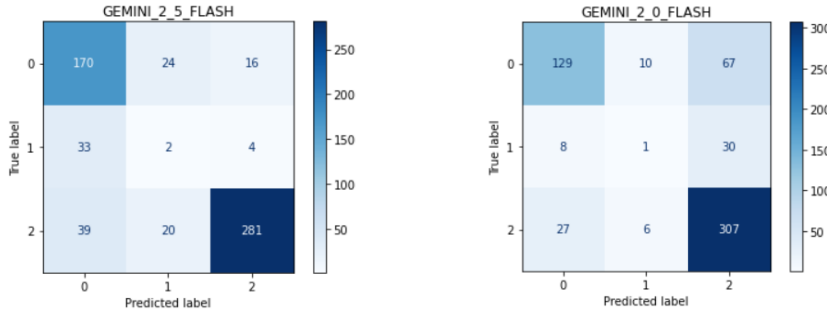


Figure 7: RAG Groundedness: Confusion matrix for short responses. Gemini 2.0 Flash performs better in this region, as human annotation scores are predominantly 0 (severely ungrounded) or 2 (fully grounded), making them easier to distinguish.

results for all judges are in Appendix E.2. The analysis reveals the variation in feature importance across judges. For instance, embedding-related features are more influential for GPT-OSS-120B, suggesting that different judges rely on distinct data properties when assessing reliability. We further aggregate the top five features that frequently appear across tasks, with results provided in Appendix E.2. The analysis reveals clear task-specific trends: text size-related features, such as word count and compression ratio, along with token entropy, are more prominent in RAG tasks. In contrast, embedding-based features, including PCA components and embedding similarity, play a more significant role in summarization tasks. These findings align with the ablation analysis in Appendix H.1, which shows that removing embedding features leads to a greater performance drop in the summarization task compared to RAG. These observations imply that evaluation reliability is task-dependent and further demonstrate that our approach effectively links data characteristics to judge reliability, enabling more informed and adaptive jury construction across diverse evaluation scenarios.

Table 2: Top 5 important features for the summarization groundedness task. Gemini 2.5 Pro relies more on embedding-based features.

Judge	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5
Gemini 2.5 Pro	output compression	input pca7	input reading index	output embedding science	output pca10
GPT-OSS-120B	output pca1	input pca1	output compression	output embedding business	output embedding legal

### 5.3 Analysis of Jury Failure in Evaluation Tasks

In this section, we investigate the conditions under which the dynamic jury fails to accurately assess model outputs. Given that jury performance is inherently dependent on the individual scoring behaviors of its constituent judges, our analysis aims to identify common regions and conditions where judges exhibit unreliable evaluations.

To uncover data attributes associated with jury failure, we train XGBoost models using a set of text attributes we measured as predictors and a binary response variable indicating jury success or failure. Feature importance is subsequently assessed via permutation-based methods to identify the most influential attributes affecting jury outcomes. For each of the top-ranked features, we conduct a binning analysis to evaluate jury performance across discrete intervals. Using the RAG groundedness task as a case study, we find that the most predictive features are related to the size and complexity of the generated text such as token entropy and character count. Figure 8 illustrates jury performance across bins of these two attributes, revealing a clear trend: as the length and entropy of generated text increase, the jury’s ability to reliably assess ground-

edness diminishes. The complete results for top features identified in this task are presented in Appendix E.5.

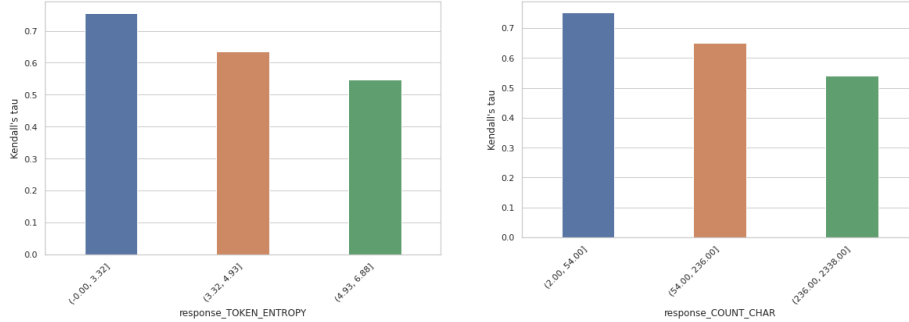


Figure 8: Jury Performance Across Binned Segments of Token Entropy and Character Count in Generated Responses. Higher jury accuracy is observed when outputs are shorter or exhibit lower textual complexity.

We further analyze the predicted reliability scores of the top-performing judges selected to form the jury, comparing cases of jury success versus failure. Using the RAG groundedness task as a representative example, we focus on the optimal jury configuration of three judges selected from a pool of ten for each evaluation instance. Figure 9 presents a box plot summarizing the predicted reliability scores of the top three judges across both successful and failed jury outcomes. The results exhibit a consistent pattern: judges involved in successful jury decisions tend to have higher predicted reliability scores than those in failure cases. This trend is corroborated by similar findings in the summarization completeness task, detailed in Figure 9 as well. In that task, the optimal jury size is seven, we highlight the top three judges for clarity.

These observations validate the effectiveness of our jury construction strategy, which prioritizes the selection of judges based on their predicted reliability. Moreover, the results suggest that the XGBoost models used to estimate individual judge reliability accurately capture behavioral patterns. When the models predict high reliability, the corresponding judges typically perform well, contributing to successful jury outcomes. Conversely, when even the top-ranked judges struggle, the overall jury performance deteriorates.

## 6 Conclusion

In this work, we addressed the critical challenge of creating scalable and reliable evaluation systems for Large Language Models. We introduced LLM Jury-on-Demand, a novel framework that moves beyond the static aggregation methods of prior jury-based systems by learning to predict judge reliability and dynamically assembling an expert jury for each data point. Our experimental results demonstrated that this adaptive approach consistently outperforms both single-



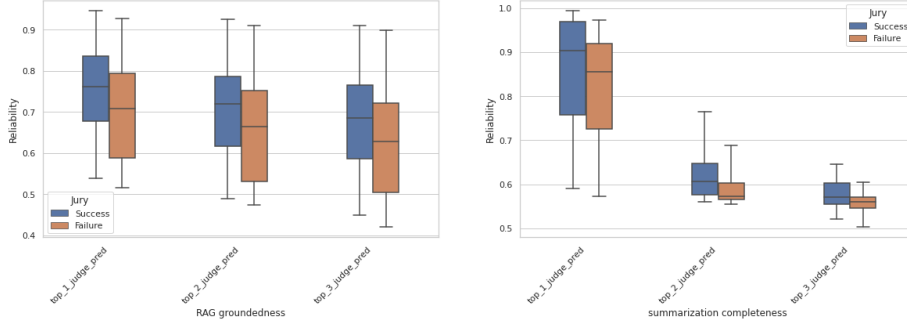


Figure 9: Predicted Reliability Scores of Top Judges for Successful and Failed Jury Outcomes in RAG Groundedness and Summarization Completeness Tasks. Higher reliability scores are consistently observed in successful jury cases, validating the effectiveness of judge selection based on model-predicted performance.

judge and static-jury baselines in aligning with human expert judgement. This confirms our central hypothesis that for automated evaluation to be trustworthy, it must be context-aware and adaptive, rather than static.

While our results are promising, this work has several limitations that open clear paths for future research. Our current framework relies on a human-annotated dataset to train the reliability predictors; future work could explore semi-supervised or self-supervised techniques to reduce this dependency and enhance scalability. Furthermore, we conduct experiments to assess the framework’s ability to generalize beyond its training domains by training on a subset of domains and applying it to held-out domains. The results indicate partial generalization: certain learned patterns transfer effectively to new domains, while others do not (see Appendix J for details). These findings suggest that the framework’s generalizability is contingent on both the diversity of the training data and the characteristics of the unseen domains. As additional annotated data becomes available and incorporated into training, we anticipate that the framework’s capacity to generalize will improve, enabling broader applicability across diverse areas.

Another promising direction for future work is mitigating bias in judge scores. For example, certain judges may consistently assign higher or lower scores compared to human annotations. We explore score calibration in Appendix K. Our experiments show that calibration can sometimes improve alignment between judge scores and human annotations, but in other cases, it may amplify the bias. Addressing this challenge remains an open problem, and we plan to investigate more robust approaches in future work.

## References

- Claude 3.7 sonnet system card. URL <https://api.semanticscholar.org/CorpusID:276612236>.
- Marah Abdin, Sahaj Agarwal, Ahmed Awadallah, Vidhisha Balachandran, Harkirat Behl, Lingjiao Chen, Gustavo de Rosa, Suriya Gunasekar, Mojan Javaheripi, Neel Joshi, et al. Phi-4-reasoning technical report. arXiv preprint arXiv:2504.21318, 2025.
- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. arXiv preprint arXiv:2508.10925, 2025.
- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. Ms marco: A human generated machine reading comprehension dataset. arXiv preprint arXiv:1611.09268, 2016.
- Joachim Baumann, Paul Röttger, Aleksandra Urman, Albert Wendsjö, Flor Miriam Plaza del Arco, Johannes B. Gruber, and Dirk Hovy. Large language model hacking: Quantifying the hidden risks of using llms for text annotation, 2025. URL <https://arxiv.org/abs/2509.08825>.
- Steven Bird, Ewan Klein, and Edward Loper. Natural Language Processing with Python. O'Reilly, 2009. ISBN 978-0-596-51649-9.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. Journal of machine Learning research, 3(Jan):993–1022, 2003.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877—1901, 2020.
- Nitay Calderon, Roi Reichart, and Rotem Dror. The alternative annotator test for llm-as-a-judge: How to statistically justify replacing human annotators with llms. arXiv preprint arXiv:2501.10970, 2025.
- Yixin Cao, Shibo Hong, Xinze Li, Jiahao Ying, Yubo Ma, Haiyuan Liang, Yantao Liu, Zijun Yao, Xiaozhi Wang, Dan Huang, et al. Toward generalizable evaluation in the llm era: A survey beyond benchmarks. arXiv preprint arXiv:2504.18838, 2025.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, pp. 785—794, 2016.

- Zhumin Chu, Qingyao Ai, Yiteng Tu, Haitao Li, and Yiqun Liu. Pre: A peer review based large language model evaluator. arXiv preprint arXiv:2401.15641, 2024.
- Norman Cliff. Dominance statistics: Ordinal analyses to answer ordinal questions. Psychological Bulletin, 114:494–509, 11 1993. doi: 10.1037/0033-2909.114.3.494.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. arXiv preprint arXiv:2507.06261, 2025.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. A dataset of information-seeking questions and answers anchored in research papers. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 4599–4610, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers), pp. 4171–4186, 2019.
- Alexander R Fabbri, Wojciech Kry\`sci\`nski, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. Summeval: Re-evaluating summarization evaluation. Transactions of the Association for Computational Linguistics, 9:391–409, 2021.
- Shangbin Feng, Wenxuan Ding, Alisa Liu, Zifeng Wang, Weijia Shi, Yike Wang, Zejiang Shen, Xiaochuang Han, Hunter Lang, Chen-Yu Lee, et al. When one llm drools, multi-llm collaboration rules. arXiv preprint arXiv:2502.04506, 2025.
- Aaron Fisher, Cynthia Rudin, and Francesca Dominici. All models are wrong, but many are useful: Learning a variable’s importance by studying an entire class of prediction models simultaneously. Journal of Machine Learning Research, 20(177):1–81, 2019.
- Mingqi Gao and Xiaojun Wan. DialSummEval: Revisiting summarization evaluation for dialogues. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 5693–5709, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.418. URL <https://aclanthology.org/2022.naacl-main.418/>.

- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to generate text with citations. In Empirical Methods in Natural Language Processing (EMNLP), 2023.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. arXiv preprint arXiv:2411.15594, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025.
- Nan Hu, Jiaoyan Chen, Yike Wu, Guilin Qi, Hongru Wang, Sheng Bi, Yongrui Chen, Tongtong Wu, and Jeff Z. Pan. Can LLMs evaluate complex attribution in QA? automatic benchmarking using knowledge graphs. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 17096—17118, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.837. URL <https://aclanthology.org/2025.acl-long.837/>.
- Ian Jolliffe. Principal component analysis. In International encyclopedia of statistical science, pp. 1094—1096. Springer, 2011.
- M. G. KENDALL. A new measure of rank correlation. Biometrika, 30(1-2):81–93, 06 1938. ISSN 0006-3444. doi: 10.1093/biomet/30.1-2.81. URL <https://doi.org/10.1093/biomet/30.1-2.81>.
- Peter Kincaid, Robert P. Fishburne, Richard L. Rogers, and Brad S. Chissom. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. 1975. URL <https://api.semanticscholar.org/CorpusID:61131325>.
- Yuhoo Lee, Taewon Yun, Jason Cai, Hang Su, and Hwanjun Song. UniSumEval: Towards unified, fine-grained, multi-dimensional summarization evaluation for LLMs. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Findings of the Association for Computational Linguistics: EMNLP 2024, pp. 3941–3960, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.227. URL <https://aclanthology.org/2024.findings-emnlp.227/>.

- Haitao Li, Junjie Chen, Qingyao Ai, Zhumin Chu, Yujia Zhou, Qian Dong, and Yiqun Liu. Calibraeval: Calibrating prediction distribution to mitigate selection bias in llms-as-judges. arXiv preprint arXiv:2410.15393, 2024a.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. Llms-as-judges: a comprehensive survey on llm-based evaluation methods. arXiv preprint arXiv:2412.05579, 2024b.
- Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Halueval: A large-scale hallucination evaluation benchmark for large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 6449–6464, 2023.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In Text summarization branches out, pp. 74–81, 2004.
- Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In International Conference on Learning Representations, 2013. URL <https://api.semanticscholar.org/CorpusID:5959482>.
- Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learning. In Proceedings of the 22nd international conference on Machine learning, pp. 625–632, 2005.
- Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, KaShun Shum, Randy Zhong, Juntong Song, and Tong Zhang. RAGTruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 10862–10878, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.585. URL <https://aclanthology.org/2024.acl-long.585/>.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp. 311–318, 2002.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162. URL <https://aclanthology.org/D14-1162/>.
- Hossein A Rahmani, Emine Yilmaz, Nick Craswell, and Bhaskar Mitra. Judgeblender: Ensembling judgments for automatic relevance assessment. arXiv preprint arXiv:2412.13268, 2024.

- Wout Schellaert, Fernando Martínez-Plumed, and José Hernández-Orallo. Analysing the predictability of language model performance. ACM Transactions on Intelligent Systems and Technology, 16(2):1–26, 2025.
- Kayla Schroeder and Zach Wood-Doughty. Can you trust llm judgments? reliability of llm-as-a-judge. arXiv preprint arXiv:2412.12509, 2024.
- Yuchen Shen and Xiaojun Wan. Opinsummeval: Revisiting automated evaluation for opinion summarization. arXiv preprint arXiv:2310.18122, 2023.
- Anil Kumar Shukla. Large language model evaluation in 2025: Smarter metrics that separate hype from trust. 2025.
- Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. ASQA: Factoid questions meet long-form answers. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pp. 8273–8288, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.566. URL <https://aclanthology.org/2022.emnlp-main.566/>.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. Advances in neural information processing systems, 33:3008–3021, 2020.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. arXiv preprint arXiv:2503.19786, 2025.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
- Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. Replacing judges with juries: Evaluating llm generations with a panel of diverse models. arXiv preprint arXiv:2404.18796, 2024.
- Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. In NeurIPS, 2023.

- Frank Wilcoxon. Individual comparisons by ranking methods. Biometrics bulletin, 1(6):80–83, 1945.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pp. 2369–2380, 2018.
- Qinyuan Ye, Harvey Yiyun Fu, Xiang Ren, and Robin Jia. How predictable are large language model capabilities? a case study on big-bench. arXiv preprint arXiv:2305.14947, 2023.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675, 2019.
- Ruochen Zhao, Wenxuan Zhang, Yew Ken Chia, Weiwen Xu, Deli Zhao, and Lidong Bing. Auto-arena: Automating llm evaluations with agent peer battles and committee discussions. arXiv preprint arXiv:2405.20267, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in neural information processing systems, 36:46595—46623, 2023.

## A List of Data Features

Table 3: Data features in judge reliability model.

Feature Name	Explanation	Category
COUNT_WORD	Number of words in the context.	Text size
COUNT_CHAR	Number of characters in the context.	Text size
COUNT_SENTENCE	Number of sentences in the context.	Text size
COUNT_PARAGRAPH	Number of paragraphs in the context.	Text size
CHAR_COMPRESSION	The ratio of the number of characters in the output context (summary or answer) to those in the input context (article or cited context).	Text size

Feature Name	Explanation	Category
WORD_COMPRESSION	The ratio of the number of words in the output context (summary or answer) to those in the input context (article or cited context).	Text size
NUM_WORD_SENTENCE	Average number of words per sentence.	Text size
NUM_CHAR_WORD	Average number of characters per word.	Text size
DIFFICULT_WORD	Number of difficult words in the context, defined as words with more than two syllables.	Special words
STOP_WORDS	Number of stop words in the context, such as “I”, “to”, “and”, “of”, etc.	Special words
MODALITY	Number of modality verbs in the context, such as “can”, “could”, “should”, etc.	Special words
NUMBER_COUNT	Count of numbers in the context, such as date, time, percent etc.	Special words
NAMED_ENTITY	Count of named entities in the context, including person, organization, date, time, Geo-Political entity, location and money.	Special words
FACTUAL_DENSITY	Number of entities divided by context length.	Special words
NGRAM_COUNT	Count of n(3)-grams in the context.	Special words
NEGATION_SENTENCE	Count of sentences with negation words such as “no”, “not”, “never”, etc.	Special words
COUNT_QUESTION	Number of questions in the context.	Special words
TOKEN_ENTROPY	Shannon entropy on token distribution.	Text complexity
LEXICAL_DIVERSITY	Number of unique words divided by the total number of words in the context.	Text complexity
READING_INDEX	Flesch reading ease index, measuring the difficulty of reading the context.	Text complexity
NGRAM_REPETITION	N(3)-gram repetition ratio in the context.	Text complexity
SENTENCE_SIMILARITY	Average cosine similarity between each pair of sentence embeddings.	Text complexity



Feature Name	Explanation	Category
SYNTACTIC_ AMBIGUITY	The average number of syntactically ambiguous POS tags (IN, TO) across sentences. These tags indicate structural complexity or multiple possible parses.	Text complexity
SEMANTIC_ AMBIGUITY	The average number of WordNet senses per word across all sentences in the text. A higher average suggests more potential meanings and interpretive complexity.	Text complexity
COREFERENCE_CHAIN	The average number of pronouns per sentence.	Text complexity
COREFERENCE_ AMBIGUOUS	Number of pronoun-ambiguous sentences in the context. Sentences with more than one pronoun are classified as ambiguous.	Text complexity
SYNTACTIC_ANOMALY	Number of syntactic anomaly sentences in the context. A sentence is syntactic anomaly if either subject or verb is missing, or both are missing.	Text complexity
RHETORICAL_STRUCTURE	Number of sentences with discourse markers (however, therefore) and rhetorical structure (moreover, in contrast, thus, instead, etc.).	Text complexity
POLARITY	The polarity score of the context, measuring the emotional tone of the text.	Text complexity
SUBJECTIVITY	The subjectivity score of the context, measuring the degree of personal opinion or factuality.	Text complexity
PCA	Text embeddings are computed via mean pooling and reduced in dimensionality using PCA. The top 10 principal components are used as features.	Embedding
Topic similarity	Cosine similarity between each text's embedding and a set of predefined topic embeddings - market, bank, business, tech, education, politics, legal, sports, media, science.	Embedding

## B Evaluation Metric Definitions

For each task, we focus on evaluating three core metrics: groundedness, relevance, and completeness.

**Groundedness:** Assesses how well the output is supported by the context of the input. A grounded output accurately reflects the source information without introducing unsupported or fabricated content. This dimension is closely related to the concept of hallucination in language models.

**Relevance:** Measures the degree to which the output includes only essential and contextually appropriate information, avoiding extraneous or off-topic content. However, for RAG, annotated data which assesses output (answer) relevance is not readily available, so instead we check retrieval relevance. Specifically, how closely and thoroughly the retrieved context addresses the posed question. A context is considered relevant if it is clearly focused on the question and provides sufficient information to support a complete and accurate answer. Similarly, we can assess how relevant the context is with respect to the reference answer.

**Completeness:** Captures whether the output includes all critical information from the input context, ensuring comprehensive coverage.

## C List of Datasets

The datasets used for different evaluation metrics are listed in Table 4 (summarization) and Table 5 (RAG). We prioritized datasets with annotated scores for completeness, groundedness, or relevance. However, annotated data for the completeness metric is relatively scarce. To address this, we simulate incomplete outputs by removing sentences from multi-sentence references and assigning scores accordingly. This approach is applied to the SummEval dataset for summarization and to all three datasets used in the RAG task.

Table 4: List of Datasets for Summarization

Metric	Data	Size	Annotation (Ann.)	Ann. scale	Domain
Completeness	TL;DR (Stiennon et al., 2020)	1680	Coverage score measures how much important information from the original post is covered	1 to 7	Reddit Discussion

Table 4: List of Datasets for Summarization

Metric	Data	Size	Annotation (Ann.)	Ann. scale	Domain
	UniSumEval (Lee et al., 2024)	1623	Completeness ratio measures the proportion of key facts inferable from the summary	0 to 1 (fraction)	Nine different domains including Wikihow, CNN/DM, GovReport, PubMed, etc.
	SummEval (Fabbri et al., 2021)	793	Original data does not have completeness score. To create annotation for completeness, we assign score 5 to the reference summary and remove different proportions of sentences in the reference summary to create incomplete summaries with scores 1 to 4	1 to 5	CNN/DM
Groundedness	SummEval (Fabbri et al., 2021)	1600	Consistency score measures the factual alignment between the summary and the summarized source. A factually consistent summary contains only statements that are entailed by the source document.	1 to 5	CNN/DM

Table 4: List of Datasets for Summarization

Metric	Data TL;DR (Stiennon et al., 2020)	Size 1597	Annotation (Ann.) Accuracy score, measures to what degree the statements in the summary are stated in the post	Ann. scale 1 to 7	Domain Reddit Discussion
	DialSummEval (Gao & Wan, 2022)	1400	Consistency score measures how well the summary aligns with the dialogue in fact. It focuses on whether the summary contains factual errors.	1 to 5	SAMSum (message dialogues)
	UniSumEval (Lee et al., 2024)	1624	Faithfulness score measures the proportion of factually correct summary sentences.	0 to 1 (frac- tion)	9 different domains including Wikihow, CNN/DM, GovRe- port, PubMed, etc.
Relevance	OpinSummEval (Shen & Wan, 2023)	1400	Aspect rele- vance, measures whether the mainly discussed aspects in the reviews are covered exactly by the summary. It focuses on whether summary correctly reflects the main aspects in the reviews.	1 to 5	Yelp reviews

Table 4: List of Datasets for Summarization

Metric	Data SummEval (Fabbri et al., 2021)	Size 1600	Annotation (Ann.) Relevance score, measures selection of important content from the source. The summary should include only important information from the source document.	Ann. scale 1 to 5	Domain CNN/DM
	DialSummEval (Gao & Wan, 2022)	1400	Relevance score, measures how well the summary captures the key points of the dialogue. It focuses on whether all and only the important aspects are contained in the summary.	1 to 5	SAMSum (message dialogues)
	UniSumEval (Lee et al., 2024)	1623	Conciseness score, measures the proportion of summary sentences aligned with the key-facts.	0 to 1 (frac- tion)	Nine different domains including Wikihow, CNN/DM, GovRe- port, PubMed, etc.

Table 5: List of Datasets for RAG

Metric	Data	Size	Annotation (Ann.)	Ann. scale	Domain
Completeness	ASQA (Stelmakh et al., 2022)	3231	Original data does not have completeness score. To create annotations for completeness, we assign score 2 to the reference answer and remove different proportions of sentences in the answer to create incomplete answers with scores 0 and 1	0, 1, 2	Wikipedia
	ALCE (Gao et al., 2023)	593	Original data does not have completeness score. To create annotations for completeness, we assign score 2 to the reference answer and remove different proportions of sentences in the answer to create incomplete answers with scores 0 and 2	0, 1, 2	Wikipedia and Reddit

Table 5: List of Datasets for RAG

Metric	Data QASPER (Dasigi et al., 2021)	Size 561	Annotation (Ann.) Original data does not have completeness score. To create annotations for completeness, we assign score 2 to the reference answer and remove different proportions of sentences in the answer to create incomplete answers with scores 0 and 3	Ann. scale 0, 1, 2	Domain NLP research papers
Groundedness	RagTruth (Niu et al., 2024)	3206	We assigned scores 0 to 2 based on count of hallucination spans in the output	0, 1, 2	MS Marco
	HaluEval (Li et al., 2023)	3000	Whether output contained hallucinated content	0, 1	HotPot-QA
	CAQA (Hu et al., 2025)	3000	Whether the cited text supports the answer. There are 4 labels, supportive, partially supportive, contradict and irrelevant	0,1,2	Knowledge Graph generated questions
Relevance	MS MARCO (Bajaj et al., 2016)	3200	Whether the cited text is relevant to the question	0, 1, 2, 3	Bing queries
	HotpotQA (Yang et al., 2018)	3000	Whether the cited text is relevant to the question	0,1,2	Wikipedia- based QA

Table 5: List of Datasets for RAG

Metric	Data	Size	Annotation (Ann.)	Ann. scale	Domain
	ALCE (Gao et al., 2023)	2896	Whether the cited text is relevant to the answer	0,1,2	Wikipedia and Reddit

## D Prompt Template

We use the same prompt template across tasks. Below we list the prompt for summarization (Table 6) and RAG (Table 7), with dimension groundedness.

Table 6: Prompt for summarization.

System prompt:
You are a helpful, respectful and honest assistant. Follow instructions attentively. Do not add any extraneous information.
User prompt:
You will be given an input context followed by a corresponding summary. Your task is to evaluate the quality of the summary with respect to groundedness.
Definition of Groundedness: Groundedness refers to how well the summary is supported by the content of the input context. A grounded summary should accurately reflect the information presented in the input context without introducing unsupported claims.
Provide a score for this summary on a scale from 1 (worst) to 5 (best). Valid scores are: 1, 2, 3, 4, or 5.
Output format:
[score number](on its own line, only one number here, no brackets or letters or 'score:')
[explanation](starting on the next line)
Conversation: source
Summary: output



Table 7: Prompt for RAG.

---

System prompt:
You are a helpful, respectful and honest assistant. Follow instructions attentively. Do not add any extraneous information.

---

User prompt:
You will be given a question ('Question' below) followed by a response ('Response' below) for the question. After that, cited background information is provided ('Context' below). The response was generated by a LLM based on the cited background information. Your task is to evaluate the quality of the response with respect to groundedness.
Definition of Groundedness: Groundedness refers to how well the response is supported by the content of the cited background information. A grounded response should accurately reflect the cited background information without introducing unsupported claims.
Provide a score for the response on a scale of 0 (bad), 1 (fair), or 2 (good). Valid scores are: 0, 1, or 2.
Output format:
[score number](on its own line, only one number here, no brackets or letters or 'score:')
[explanation](starting on the next line)
Question: question
Response: response
Context: cited text

---

## E Full Experimental Results

We first show the jury performance, then analyze the interactions between judges, tasks and data properties.

### E.1 Full Results by Task and Dataset

Tables 8 - 11 present the complete Kendall’s Tau correlation results of our experiments. Due to the table size limit, we split the results for clarity. Tables 8 and 10 compare our Jury-on-Demand system against the four Static Jury baselines (Average-All, Average-TopK, Weighted-Regression, and Weighted-Tau). Tables 9 and 11 compare Jury-on-Demand against the 10 Single-Judge baselines. Each row corresponds to a specific evaluation set, either an “Overall” aggregation or an individual source dataset. All values represent the mean Kendall’s Tau rank correlation coefficient ( $\pm$  standard deviation) calculated across 10 independent runs. Higher values indicate better performance and stronger alignment with human judgment.

We also report statistical significance and effect sizes. Specifically, we perform one-sided Wilcoxon signed-rank tests (Wilcoxon, 1945) to compare the Tau differences between Jury-on-Demand and either static juries or single judges. The corresponding p-values are presented in Table 12 through Table 15. Among these p-values, 77% are statistically significant ( $p < 0.05$ ). Values in parentheses represent Cliff’s delta (Cliff, 1993), a non-parametric effect size metric that quantifies the difference between two groups, with Jury-on-Demand serving as the baseline. According to conventional thresholds, an effect size is considered large if it exceeds 0.47 and medium if it falls between 0.33 and 0.47. Across all Cliff’s delta values, 80% are classified as either large (70%) or medium (10%). This high proportion of significant p-values and substantial effect sizes indicates that, in most cases, Jury-on-Demand outperforms static baselines and single judges.

Table 8: Summarization Results: Jury-on-Demand vs Static Jury baselines. Numbers in parentheses are standard deviation. **Bold** indicates the highest mean and its std in the row, underline indicates the second highest mean and its std.

Data	Jury-on-Demand	Static Jury (Average-All)	Static Jury (Average-TopK)	Static Jury (Weighted-Regression)	Static Jury (Weighted-Tau)
<b>Completeness</b>					
Overall	<b>0.48</b> (0.03)	0.44 (0.02)	<u>0.47</u> (0.02)	0.45 (0.02)	0.45 (0.02)
SummEval	<b>0.72</b> (0.05)	0.60 (0.06)	0.67 (0.08)	<u>0.69</u> (0.04)	0.61 (0.05)
TL;DR	0.38 (0.08)	0.40 (0.05)	<b>0.44</b> (0.05)	<u>0.41</u> (0.02)	0.41 (0.04)
UniSumEval	<b>0.66</b> (0.04)	0.59 (0.03)	<u>0.63</u> (0.04)	0.59 (0.04)	0.61 (0.05)
<b>Groundedness</b>					
Overall	0.54 (0.04)	0.55 (0.02)	<b>0.56</b> (0.03)	0.56 (0.02)	<u>0.56</u> (0.02)
DialSummEval	0.67 (0.05)	0.66 (0.02)	<b>0.69</b> (0.02)	<u>0.68</u> (0.02)	0.68 (0.02)
SummEval	0.61 (0.08)	<u>0.65</u> (0.04)	0.62 (0.05)	0.64 (0.03)	<b>0.65</b> (0.04)
TL;DR	0.43 (0.05)	0.45 (0.06)	0.46 (0.05)	<b>0.46</b> (0.05)	<u>0.46</u> (0.06)
UniSumEval	0.62 (0.07)	0.63 (0.07)	0.64 (0.07)	<u>0.64</u> (0.07)	<b>0.65</b> (0.06)
<b>Relevance</b>					
Overall	<b>0.64</b> (0.02)	0.62 (0.02)	0.63 (0.02)	0.61 (0.03)	<u>0.64</u> (0.01)
DialSummEval	<b>0.69</b>	0.65	0.63	0.65	<u>0.66</u>

Data	Jury-on-Demand	Static Jury (Average-All)	Static Jury (Average-TopK)	Static Jury (Weighted-Regression)	Static Jury (Weighted-Tau)
	<b>(0.04)</b>	(0.03)	(0.04)	(0.03)	<u>(0.03)</u>
OpinSummEval	<b>0.46</b> <b>(0.06)</b>	0.42 (0.04)	0.40 (0.07)	0.42 (0.04)	<u>0.44</u> <u>(0.04)</u>
SummEval	<u>0.71</u> <u>(0.07)</u>	0.70 (0.04)	<b>0.72</b> <b>(0.05)</b>	0.69 (0.08)	0.70 (0.04)
UniSumEval	<b>0.42</b> <b>(0.09)</b>	0.41 (0.09)	0.41 (0.11)	0.40 (0.10)	<u>0.42</u> <u>(0.09)</u>

Table 9: Summarization Results: Jury-on-Demand vs Single Judge baselines. Numbers in parentheses are standard deviation. Here Gemn. is Gemini and LL is LLAMA. **Bold** indicates the highest mean and its std in the row, underline indicates the second highest mean and its std.

Data	Jury-on-Demand	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	Gem ma3	Phi4	Deep Seek R1
<b>Completeness</b>											
Overall	<b>0.48</b> <b>(0.03)</b>	0.42 (0.03)	0.45 (0.03)	0.46 (0.02)	0.45 (0.02)	0.35 (0.04)	<u>0.46</u> <u>(0.02)</u>	0.12 (0.04)	0.26 (0.04)	0.18 (0.04)	0.17 (0.05)
Summ Eval	<b>0.72</b> <b>(0.05)</b>	0.57 (0.09)	0.63 (0.08)	<u>0.68</u> <u>(0.07)</u>	0.58 (0.11)	-0.10 (0.15)	0.67 (0.08)	0.12 (0.14)	0.33 (0.16)	0.27 (0.06)	0.68 (0.04)
TL;DR	0.38 (0.08)	0.37 (0.07)	0.39 (0.05)	<u>0.40</u> <u>(0.06)</u>	<b>0.43</b> <b>(0.06)</b>	0.34 (0.05)	0.40 (0.04)	0.09 (0.07)	0.14 (0.06)	0.01 (0.05)	0.05 (0.05)
Uni Sum Eval	<b>0.66</b> <b>(0.04)</b>	0.52 (0.04)	0.54 (0.07)	0.59 (0.04)	0.52 (0.05)	0.62 (0.04)	<u>0.66</u> <u>(0.04)</u>	0.29 (0.10)	0.55 (0.06)	0.29 (0.07)	0.41 (0.08)
<b>Groundedness</b>											
Overall	<u>0.54</u> <u>(0.04)</u>	<b>0.56</b> <b>(0.02)</b>	0.49 (0.02)	0.49 (0.03)	0.51 (0.03)	0.45 (0.03)	0.52 (0.03)	0.23 (0.03)	0.51 (0.02)	0.24 (0.03)	0.20 (0.04)
Dial Summ Eval	0.67 (0.05)	0.59 (0.04)	0.63 (0.03)	<b>0.70</b> <b>(0.03)</b>	<u>0.69</u> <u>(0.03)</u>	0.65 (0.04)	0.66 (0.03)	0.33 (0.06)	0.64 (0.04)	0.24 (0.06)	0.33 (0.08)
Summ Eval	0.61 (0.08)	<u>0.63</u> <u>(0.05)</u>	0.56 (0.06)	0.61 (0.04)	<b>0.64</b> <b>(0.05)</b>	0.56 (0.06)	0.60 (0.04)	0.21 (0.09)	0.59 (0.06)	0.17 (0.09)	0.12 (0.10)
TL;DR	<b>0.43</b> <b>(0.05)</b>	0.39 (0.04)	0.41 (0.05)	0.40 (0.05)	<u>0.42</u> <u>(0.04)</u>	0.29 (0.06)	0.39 (0.06)	0.10 (0.05)	0.42 (0.04)	0.11 (0.07)	0.13 (0.05)
Uni Sum Eval	<u>0.62</u> <u>(0.07)</u>	0.61 (0.07)	0.57 (0.08)	0.61 (0.08)	<b>0.64</b> <b>(0.07)</b>	0.61 (0.09)	0.59 (0.10)	0.10 (0.18)	0.55 (0.09)	0.31 (0.09)	0.24 (0.18)
<b>Relevance</b>											

Data	Jury-on-Demand	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	Gemma3	Phi4	Deep Seek R1
Overall	<b>0.64</b> (0.02)	0.58 (0.02)	0.59 (0.02)	0.58 (0.01)	0.58 (0.01)	0.56 (0.02)	<u>0.61</u> (0.02)	0.19 (0.08)	0.59 (0.03)	0.25 (0.02)	0.34 (0.05)
Dial Summ Eval	<b>0.69</b> (0.04)	0.52 (0.06)	0.56 (0.05)	0.59 (0.03)	0.51 (0.04)	0.66 (0.04)	<u>0.68</u> (0.03)	0.21 (0.09)	0.59 (0.03)	0.44 (0.04)	0.42 (0.08)
Opin Summ Eval	<b>0.46</b> (0.06)	0.31 (0.08)	0.36 (0.05)	0.30 (0.06)	0.33 (0.05)	0.37 (0.07)	<u>0.36</u> (0.08)	0.22 (0.09)	<u>0.39</u> (0.08)	0.18 (0.09)	0.25 (0.07)
Summ Eval	<b>0.71</b> (0.07)	0.67 (0.05)	0.69 (0.06)	0.64 (0.09)	0.66 (0.05)	0.66 (0.05)	0.68 (0.09)	0.40 (0.17)	<u>0.69</u> (0.05)	0.25 (0.09)	0.44 (0.06)
Uni Sum Eval	<u>0.42</u> (0.09)	<b>0.43</b> (0.11)	0.37 (0.09)	0.38 (0.10)	<b>0.43</b> (0.07)	0.32 (0.07)	0.37 (0.09)	0.14 (0.12)	0.40 (0.10)	0.11 (0.08)	0.28 (0.12)

Table 10: RAG Results: Jury-on-Demand vs Static Jury baselines. Numbers in parentheses are standard deviation. **Bold** indicates the highest mean and its std in the row, underline indicates the second highest mean and its std.

Data	Jury-on-Demand	Static Jury (Average-All)	Static Jury (Average-TopK)	Static Jury (Weighted-Regression)	Static Jury (Weighted-Tau)
<b>Completeness</b>					
Overall	<b>0.50</b> (0.03)	<u>0.38</u> (0.03)	0.36 (0.03)	0.34 (0.03)	0.31 (0.03)
ALCE	<b>0.47</b> (0.07)	<u>0.38</u> (0.09)	0.28 (0.08)	0.34 (0.11)	0.23 (0.10)
ASQA	<b>0.54</b> (0.05)	<u>0.38</u> (0.05)	0.38 (0.04)	0.36 (0.04)	0.34 (0.03)
QASPER	<b>0.44</b> (0.08)	<u>0.41</u> (0.08)	0.27 (0.08)	0.35 (0.07)	0.24 (0.11)
<b>Groundedness</b>					
Overall	<b>0.68</b> (0.02)	0.58 (0.02)	0.59 (0.02)	<u>0.65</u> (0.02)	0.47 (0.02)
CAQA	<b>0.68</b> (0.03)	0.56 (0.03)	0.60 (0.03)	<u>0.62</u> (0.03)	0.50 (0.05)
HaluEval	<b>0.77</b> (0.02)	0.73 (0.02)	<u>0.74</u> (0.04)	0.74 (0.04)	0.53 (0.05)
RagTruth	<b>0.57</b> (0.03)	0.53 (0.05)	0.34 (0.07)	<u>0.55</u> (0.02)	0.15 (0.08)
<b>Relevance</b>					
Overall	<b>0.61</b>	0.57	0.54	<u>0.59</u>	0.53

Data	Jury-on-Demand	Static Jury (Average-All)	Static Jury (Average-TopK)	Static Jury (Weighted-Regression)	Static Jury (Weighted-Tau)
	<b>(0.01)</b>	(0.01)	(0.01)	<u>(0.01)</u>	(0.01)
ALCE	<b>0.61</b> <b>(0.03)</b>	0.60 (0.03)	0.44 (0.04)	<u>0.60</u> <u>(0.03)</u>	0.38 (0.03)
HotpotQA	<b>0.90</b> <b>(0.02)</b>	0.86 (0.02)	<u>0.87</u> <u>(0.03)</u>	0.85 (0.02)	0.86 (0.03)
MS MARCO	<b>0.46</b> <b>(0.04)</b>	0.39 (0.04)	<u>0.45</u> <u>(0.04)</u>	0.44 (0.03)	0.44 (0.03)

Table 11: RAG Results: Jury-on-Demand vs Single Judge baselines. Numbers in parentheses are standard deviation. Here Gemn. is Gemini and LL is LLAMA. **Bold** indicates the highest mean and its std in the row, underline indicates the second highest mean and its std.

Data	Jury-on-Demand	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	Gem ma3	Phi4	Deep Seek R1
<b>Completeness</b>											
Overall	<b>0.50</b> <b>(0.03)</b>	0.36 (0.02)	0.34 (0.02)	0.35 (0.02)	<u>0.37</u> <u>(0.02)</u>	0.33 (0.02)	0.35 (0.02)	0.19 (0.04)	0.33 (0.02)	0.16 (0.04)	0.21 (0.03)
ALCE	<b>0.47</b> <b>(0.07)</b>	0.38 (0.09)	0.34 (0.10)	0.30 (0.10)	0.34 (0.08)	<u>0.40</u> <u>(0.07)</u>	0.36 (0.07)	0.16 (0.10)	0.33 (0.07)	0.08 (0.10)	0.32 (0.09)
ASQA	<b>0.54</b> <b>(0.05)</b>	<u>0.42</u> <u>(0.03)</u>	0.38 (0.04)	0.33 (0.05)	0.34 (0.05)	0.33 (0.04)	0.38 (0.05)	0.22 (0.05)	0.36 (0.05)	0.13 (0.04)	0.20 (0.06)
QASPER	<b>0.44</b> <b>(0.08)</b>	0.31 (0.07)	0.23 (0.08)	0.41 (0.08)	0.35 (0.08)	0.40 (0.08)	<u>0.43</u> <u>(0.07)</u>	0.05 (0.10)	0.32 (0.10)	-0.02 (0.07)	0.15 (0.14)
<b>Groundedness</b>											
Overall	<b>0.68</b> <b>(0.02)</b>	0.56 (0.02)	0.51 (0.02)	0.62 (0.02)	0.58 (0.02)	0.62 (0.01)	<u>0.63</u> <u>(0.02)</u>	0.11 (0.03)	0.06 (0.02)	0.08 (0.03)	0.28 (0.03)
CAQA	<b>0.68</b> <b>(0.03)</b>	0.56 (0.03)	0.59 (0.03)	0.59 (0.03)	0.56 (0.02)	0.60 (0.02)	<u>0.60</u> <u>(0.03)</u>	0.08 (0.06)	0.02 (0.02)	0.01 (0.03)	0.26 (0.04)
Halu Eval	<u>0.77</u> <u>(0.02)</u>	0.67 (0.04)	0.52 (0.03)	0.74 (0.04)	0.73 (0.03)	0.76 (0.03)	<b>0.78</b> <b>(0.03)</b>	0.20 (0.05)	0.17 (0.04)	0.14 (0.05)	0.40 (0.03)
Rag Truth	<b>0.57</b> <b>(0.03)</b>	0.40 (0.06)	0.30 (0.04)	<u>0.56</u> <u>(0.03)</u>	0.49 (0.05)	0.51 (0.05)	0.52 (0.05)	0.12 (0.05)	0.14 (0.05)	0.07 (0.09)	0.14 (0.07)
<b>Relevance</b>											
Overall	<b>0.61</b> <b>(0.01)</b>	0.55 (0.01)	0.51 (0.02)	0.58 (0.01)	0.57 (0.01)	0.57 (0.01)	<u>0.58</u> <u>(0.01)</u>	0.19 (0.02)	0.41 (0.02)	0.23 (0.02)	0.30 (0.02)
ALCE	<b>0.61</b> <b>(0.03)</b>	0.51 (0.03)	0.53 (0.04)	0.57 (0.04)	0.58 (0.05)	<u>0.59</u> <u>(0.03)</u>	0.58 (0.03)	0.11 (0.04)	0.43 (0.03)	0.13 (0.05)	0.28 (0.05)
Hotpot QA	<b>0.90</b>	0.87	0.85	0.89	0.89	0.89	<u>0.90</u>	0.40	0.78	0.45	0.60

Data	Jury-on-Demand	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	Gemma3	Phi4	Deep Seek R1
	<b>(0.02)</b>	(0.03)	(0.02)	(0.02)	(0.02)	(0.02)	<u>(0.01)</u>	(0.04)	(0.05)	(0.04)	(0.05)
MS MARCO	<b>0.46</b>	0.40	0.40	0.43	<u>0.44</u>	0.39	0.43	0.14	0.17	0.10	0.19
	<b>(0.04)</b>	(0.04)	(0.05)	(0.03)	<u>(0.03)</u>	(0.04)	(0.03)	(0.06)	(0.05)	(0.05)	(0.06)

Table 12: Summarization: p-value of Wilcoxon test for the Tau difference between Jury-On-Demand and Static Jury baselines. Numbers in paranthesis () are effect size Cliff’s delta.

Data	Static Jury (Average-All)	Static Jury (Average-TopK)	Static Jury (Weighted-Regression)	Static Jury (Weighted-Tau)
<b>Completeness</b>				
Overall	0.001 (0.76)	0.161 (0.26)	0.024 (0.56)	0.019 (0.64)
Summ - Eval	0.001 (0.88)	0.116 (0.34)	0.116 (0.44)	0.001 (0.86)
TL;DR	0.722 (-0.08)	0.981 (-0.50)	0.722 (-0.06)	0.919 (-0.24)
UniSum - Eval	0.002 (0.82)	0.014 (0.46)	0.001 (0.86)	0.005 (0.54)
<b>Groundedness</b>				
Overall	0.958 (-0.22)	0.981 (-0.42)	0.995 (-0.36)	0.997 (-0.38)
Summ - Eval	0.958 (-0.36)	0.652 (-0.14)	0.968 (-0.34)	0.976 (-0.42)
TL;DR	0.903 (-0.14)	0.958 (-0.30)	0.986 (-0.36)	0.919 (-0.16)
UniSum - Eval	0.500 (-0.10)	0.813 (-0.10)	0.688 (-0.18)	0.903 (-0.32)
Dial - Summ - Eval	0.313 (0.08)	0.958 (-0.40)	0.784 (-0.22)	0.688 (-0.16)
<b>Relevance</b>				
Overall	0.002 (0.60)	0.080 (0.14)	0.002 (0.64)	0.116 (0.24)
Dial - Summ - Eval	0.001 (0.62)	0.002 (0.74)	0.007 (0.54)	0.032 (0.46)
Opin - Summ - Eval	0.001 (0.34)	0.001 (0.46)	0.010 (0.42)	0.007 (0.26)
Summ - Eval	0.216	0.903	0.032	0.313

Data	Static Jury (Average-All)	Static Jury (Average- TopK)	Static Jury (Weighted- Regression)	Static Jury (Weighted- Tau)
	(0.16)	(-0.08)	(0.18)	(0.14)
UniSum - Eval	0.461 (0.08)	0.461 (0.04)	0.312 (0.08)	0.500 (0.04)

Table 13: Summarization: p-value of Wilcoxon test for the Tau difference between Jury-On-Demand and single judge. Here Gemn. is Gemini and LL is LLAMA. Numbers in paranthesis () are effect size Cliff’s delta.

Data	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LL- 3.2	Gem ma3	Phi4	Deep Seek R1
<b>Completeness</b>										
Overall	0.002 (0.92)	0.024 (0.62)	0.032 (0.50)	0.024 (0.50)	0.001 (1.00)	0.042 (0.40)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
Summ - Eval	0.005 (0.92)	0.014 (0.66)	0.096 (0.32)	0.005 (0.76)	0.001 (1.00)	0.096 (0.46)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.042 (0.56)
TL;DR	0.312 (0.04)	0.500 (0.02)	0.652 (- 0.12)	0.920 (- 0.34)	0.065 (0.44)	0.615 (0.04)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)
UniSum - Eval	0.001 (1.00)	0.001 (0.90)	0.002 (0.78)	0.001 (1.00)	0.003 (0.62)	0.042 (0.12)	0.001 (1.00)	0.002 (0.94)	0.001 (1.00)	0.001 (1.00)
<b>Groundedness</b>										
Overall	0.967 (- 0.38)	0.001 (0.68)	0.001 (0.76)	0.005 (0.40)	0.001 (1.00)	0.080 (0.24)	0.001 (1.00)	0.005 (0.48)	0.001 (1.00)	0.001 (1.00)
Summ - Eval	0.652 (- 0.20)	0.065 (0.36)	0.652 (- 0.08)	0.813 (- 0.24)	0.065 (0.36)	0.500 (0.08)	0.001 (1.00)	0.313 (0.02)	0.001 (1.00)	0.001 (1.00)
TL;DR	0.019 (0.52)	0.116 (0.10)	0.053 (0.38)	0.313 (0.14)	0.001 (0.94)	0.019 (0.38)	0.001 (1.00)	0.116 (0.24)	0.001 (1.00)	0.001 (1.00)
UniSum - Eval	0.461 (0.02)	0.053 (0.32)	0.461 (0.04)	0.784 (- 0.20)	0.461 (0.08)	0.188 (0.22)	0.001 (1.00)	0.007 (0.24)	0.001 (1.00)	0.001 (1.00)
Dial - Summ - Eval	0.001 (0.86)	0.024 (0.50)	0.935 (- 0.44)	0.839 (- 0.38)	0.188 (0.20)	0.348 (0.06)	0.001 (1.00)	0.014 (0.36)	0.001 (1.00)	0.001 (1.00)
<b>Relevance</b>										
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.002 (0.72)	0.001 (1.00)	0.002 (0.82)	0.001 (1.00)	0.001 (1.00)

<b>Data Claude 3.7</b>	<b>Gemn. 2.0 Flash</b>	<b>Gemn. 2.5 Flash</b>	<b>Gemn. 2.5 Pro</b>	<b>GPT- OSS- 20B</b>	<b>GPT- OSS- 120B</b>	<b>LL- 3.2</b>	<b>Gem ma3</b>	<b>Phi4</b>	<b>Deep Seek R1</b>	
Dial - Summ - Eval	0.001 (1.00)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.065 (0.40)	0.216 (0.06)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)
Opin - Summ - Eval	0.001 (0.94)	0.001 (0.84)	0.001 (0.94)	0.001 (0.94)	0.001 (0.78)	0.001 (0.68)	0.001 (1.00)	0.001 (0.42)	0.001 (1.00)	0.001 (1.00)
Summ - Eval	0.024 (0.40)	0.053 (0.26)	0.001 (0.48)	0.053 (0.46)	0.007 (0.50)	0.065 (0.18)	0.001 (1.00)	0.188 (0.18)	0.001 (1.00)	0.001 (1.00)
UniSum - Eval	0.577 (- 0.04)	0.042 (0.24)	0.065 (0.18)	0.615 (0.00)	0.001 (0.64)	0.003 (0.26)	0.001 (0.94)	0.313 (0.10)	0.001 (0.98)	0.001 (0.64)

Table 14: RAG: p-value of Wilcoxon test for the Tau difference between Jury-On-Demand and Static Jury baselines. Numbers in paranthesis () are effect size Cliff’s delta.

<b>Data</b>	<b>Static Jury (Average-All)</b>	<b>Static Jury (Average- TopK)</b>	<b>Static Jury (Weighted- Regression)</b>	<b>Static Jury (Weighted- Tau)</b>
<b>Completeness</b>				
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
ALCE	0.001 (0.58)	0.001 (0.96)	0.003 (0.78)	0.001 (0.94)
ASQA	0.001 (0.96)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
QASPER	0.188 (0.24)	0.001 (0.92)	0.007 (0.66)	0.005 (0.84)
<b>Groundedness</b>				
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (0.82)	0.001 (1.00)
CAQA	0.001 (1.00)	0.001 (0.98)	0.001 (0.78)	0.001 (1.00)
Halu - Eval	0.002 (0.74)	0.001 (0.48)	0.005 (0.52)	0.001 (1.00)
Rag - Truth	0.001 (0.44)	0.001 (1.00)	0.010 (0.32)	0.001 (1.00)
<b>Relevance</b>				
Overall	0.001 (0.80)	0.001 (1.00)	0.003 (0.90)	0.001 (1.00)
ALCE	0.001	0.001	0.066	0.001



Data	Static Jury (Average-All)	Static Jury (Average- TopK)	Static Jury (Weighted- Regression)	Static Jury (Weighted- Tau)
	(0.36)	(1.00)	(0.30)	(1.00)
Hotpot - QA	0.001 (0.92)	0.003 (0.74)	0.001 (1.00)	0.001 (0.86)
MS MARCO	0.001 (0.84)	0.138 (0.20)	0.003 (0.28)	0.001 (0.36)

Table 15: RAG: p-value of Wilcoxon test for the Tau difference between Jury-On-Demand and single judge. Numbers in paranthesis () are effect size Cliff’s delta.

Data	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LL- 3.2	Gem ma3	Phi4	Deep Seek R1
<b>Completeness</b>										
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
ALCE	0.001 (0.56)	0.001 (0.70)	0.001 (0.80)	0.001 (0.76)	0.001 (0.54)	0.001 (0.70)	0.001 (1.00)	0.001 (0.82)	0.001 (1.00)	0.001 (0.88)
ASQA	0.001 (0.98)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (0.96)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)
QASPER	0.005 (0.78)	0.001 (0.96)	0.246 (0.16)	0.005 (0.64)	0.053 (0.30)	0.216 (0.10)	0.001 (1.00)	0.002 (0.64)	0.001 (1.00)	0.001 (0.92)
<b>Groundedness</b>										
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (0.96)	0.001 (1.00)	0.001 (0.96)	0.001 (0.92)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
CAQA	0.001 (1.00)	0.001 (0.96)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (0.92)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
Halu - Eval	0.001 (1.00)	0.001 (1.00)	0.014 (0.42)	0.001 (0.64)	0.024 (0.28)	0.461 (- 0.06)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
Rag - Truth	0.001 (1.00)	0.001 (1.00)	0.001 (0.24)	0.001 (0.80)	0.001 (0.72)	0.001 (0.62)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
<b>Relevance</b>										
Overall	0.001 (1.00)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)	0.001 (0.98)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
ALCE	0.001 (1.00)	0.001 (0.96)	0.001 (0.62)	0.001 (0.56)	0.001 (0.44)	0.001 (0.48)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
Hotpot - QA	0.001 (0.82)	0.001 (1.00)	0.001 (0.34)	0.002 (0.56)	0.001 (0.64)	0.001 (0.22)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)	0.001 (1.00)
MS MARCO	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001

Data	Claude 3.7	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LL- 3.2	Gem ma3	Phi4	Deep Seek R1
	(0.72)	(0.66)	(0.52)	(0.34)	(0.80)	(0.52)	(1.00)	(1.00)	(1.00)	(1.00)

## E.2 Analyzing the Interaction Between Judges, Tasks, and Data Attributes

Fig. 10 summarizes the selection frequency of each judge, revealing that summarization juries tend to include a broader set of judges compared to those used in RAG tasks. Within RAG, Claude 3.7 Sonnet and DeepSeek R1 are frequently selected for completeness evaluation but are rarely chosen for groundedness. In contrast, Gemini 2.5 Flash is commonly selected for groundedness but appears less frequently in completeness evaluations. GPT OSS 20B and GPT OSS 120B are consistently selected across both metrics.

To further explore how data properties influence model performance and judge selection, we examine Kendall’s tau and selection patterns across different bins of key attributes. Following Section 5.2. We focus on summarization completeness task and property compression ratio. Fig. 11 compares judge performance across low, medium, and high compression ratio bins. Performance improves as the compression ratio increases, likely because judges find it more difficult to identify incomplete summaries, which tend to have lower compression ratios. Fig. 11 shows that Gemma performs particularly poorly when the compression ratio is low. Upon reviewing the data, Gemma’s scores, and its explanations, we found that it sometimes struggles to distinguish between the source context and the summary. As a result, it incorrectly assigns a score of 2, interpreting the context as part of the summary. Examples illustrating this issue are provided in the Section E.3. Stronger models—such as the Gemini series, GPT models, and Claude 3.7 Sonnet—do not exhibit this issue and are able to correctly identify missing content in the summaries.

Jury selection, shown in the right plot, is generally consistent with judge performance across the bins. In the low compression bin where Gemma fails, it is selected least often. In contrast, Gemini 2.0 Flash has the highest performance in this bin and is selected in nearly all juries, showing strong alignment. However, the selection mechanism again proves more sophisticated than just selecting the top-ranked judge. In the high compression ratio bin, Gemini 2.5 Flash (red) achieves the highest Kendall’s Tau. Yet, Gemini 2.0 Flash (blue), which also performs well, is selected more frequently. This demonstrates that the reliability model identifies multiple judges as reliable in this context and dynamically constructs juries based on this broader reliability assessment rather than overfitting to a single “best” judge for the bin.

These findings reinforce the importance of constructing dynamic juries that adapt to specific data characteristics, and demonstrate the potential of predicting judge reliability based on interpretable data properties.

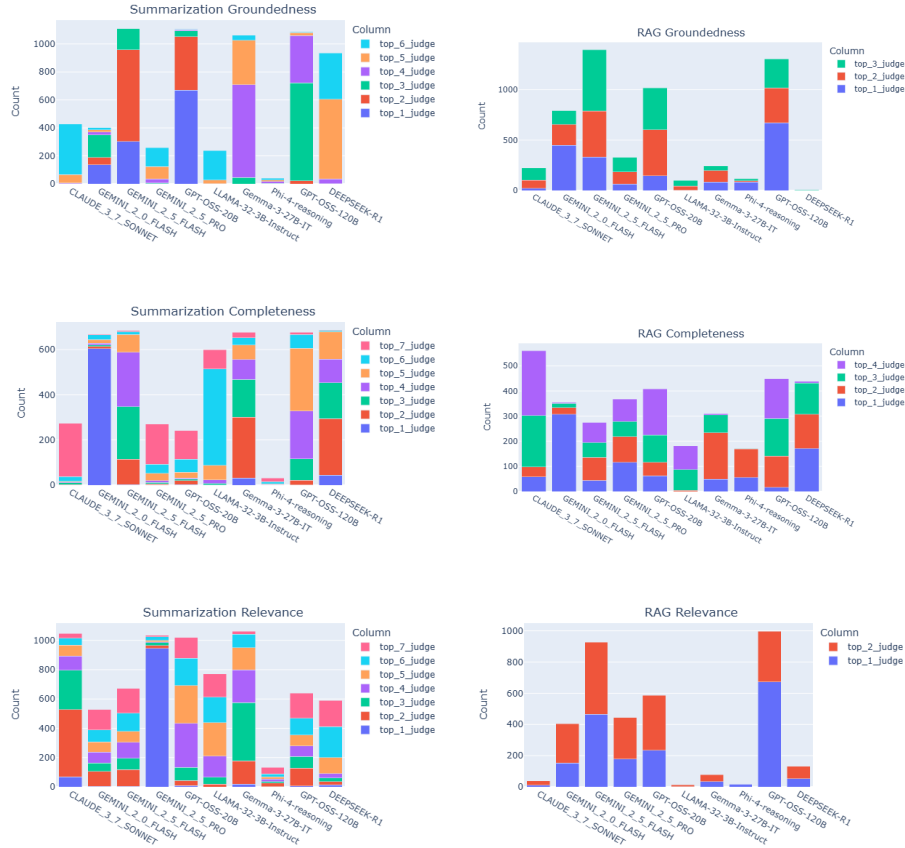


Figure 10: Selection frequency of the judge in the jury. Top k judge means that the judge has the k-th highest reliability score in the jury. Summarization juries tend to incorporate a more diverse set of judges compared to those used in RAG tasks. For RAG Claude 3.7 Sonnet and DeepSeek R1 are frequently selected for completeness evaluation but are rarely chosen for groundedness. In contrast, Gemini 2.5 Flash is commonly selected for groundedness but appears less frequently in completeness evaluations. GPT OSS 20B and GPT OSS 120B are consistently selected across both metrics.

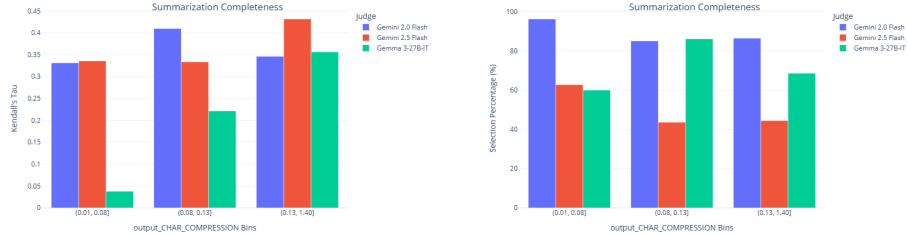


Figure 11: Summarization Completeness analysis by character compression ratio. (Left) Kendall’s Tau correlation for three single judges across low, medium, and high compression ratio bins. (Right) The selection percentage of three judges in the final dynamic jury for data points within each bin. Performance generally degrades at lower compression ratio, with Gemma 3-27B-IT failing significantly in this region. The jury selection percentage correctly mirrors this, heavily selecting Gemini 2.0 Flash which is the most reliable in that bin.

Tables 16-17 present the top five most important features for each judge’s XGBoost model, as determined by permutation feature importance (Fisher et al., 2019), using the summarization groundedness task as an illustrative example. The results show substantial variation in the top-ranked features across different judges, suggesting that each judge’s reliability is influenced by distinct data properties.

Fig. 12 aggregates the top five features that frequently appear across tasks, revealing clear task-specific patterns. For instance, character count is more prominent in RAG tasks, while compression ratio and embedding-related features such as PCA components and embedding similarity are more influential in summarization tasks. These findings align with the ablation analysis in Appendix H.1, which shows that removing embedding features leads to a greater performance drop in the summarization task compared to RAG. For RAG, removing text size-related features results in a larger decline than removing embedding features.

These observations imply that evaluation reliability is task-dependent, and further demonstrate that our approach effectively links data characteristics to judge reliability, enabling more informed and adaptive jury construction across diverse evaluation scenarios.

### E.3 Human-Reviewed Examples

In this section, we present examples where specific judges fail to evaluate correctly. For illustration, we select one example and one judge from each of the two tasks: RAG groundedness and summarization completeness. Fig. 13 shows how Gemini 2.0 Flash evaluates groundedness in the RAG task. It fails to identify ungrounded content in the response—for instance, Okavango Delta, which is not mentioned in the cited context. Fig. 14 illustrates how Gemma 3 27B IT

Table 16: The top 3 most important features for each judge’s XGBoost model from summarization groundedness. The results show substantial variation in the top-ranked features across different judges.

Judge	Feature 1	Feature 2	Feature 3
CLAUDE 3 .7 SONNET	input __ embedding __ similarity __ politics	input __ COUNT __ WORD	input __ SUBJECTIVITY
GEMINI 2.0 FLASH	input_pca_9	output_pca_9	output_pca_7
GEMINI 2.5 FLASH	input_pca_9	output __ embedding __ similarity __ legal	output_pca_1
GEMINI 2.5 PRO	input __ embedding __ similarity __ financemarket	output __ embedding __ similarity __ financemarket	output __ CHAR __ COMPRESSION
GPT-OSS-20B	output_pca_1	output __ WORD __ COMPRESSION	output_pca_2
LLAMA-32-3B-Instruct	output __ READING __ INDEX	input __ embedding __ similarity __ media	input __ READING __ INDEX
Gemma-3-27B-IT	input __ embedding __ similarity __ sports	output __ SENTENCE __ SIMILARITY	output __ embedding __ similarity __ media
Phi-4-reasoning	input __ DIFFICULT __ WORD	output __ NUM __ WORD __ SENTENCE	output __ LEXICAL __ DIVERSITY
GPT-OSS-120B	input __ embedding __ similarity __ financemarket	output __ embedding __ similarity __ financemarket	output __ CHAR __ COMPRESSION
DEEPSEEK-R1	input_pca_6	output __ embedding __ similarity __ legal	output __ embedding __ similarity __ financemarket

Table 17: The 4th and 5th most important features for each judge’s XGBoost model from summarization groundedness. The results show substantial variation in the top-ranked features across different judges.

Judge	Feature 4	Feature 5
CLAUDE 3 .7 SONNET	output __ DIFFICULT __ WORD	output_pca_8
GEMINI 2.0 FLASH	input_pca_7	output_pca_3
GEMINI 2.5 FLASH	input_pca_1	output __ WORD __ COMPRESSION
GEMINI 2.5 PRO	output_pca_3	output_pca_1
GPT-OSS-20B	output_pca_9	output __ embedding __ similarity __ business
LLAMA-32-3B-Instruct	output __ COUNT __ WORD	output_pca_1
Gemma-3-27B-IT	output __ COREFERENCE __ CHAIN	input __ SEMANTIC __ AMBIGUITY
Phi-4-reasoning	output_pca_5	output __ READING __ INDEX
GPT-OSS-120B	output_pca_3	output_pca_1
DEEPSEEK-R1	output __ SENTENCE __ SIMILARITY	input __ NAMED __ ENTITIE

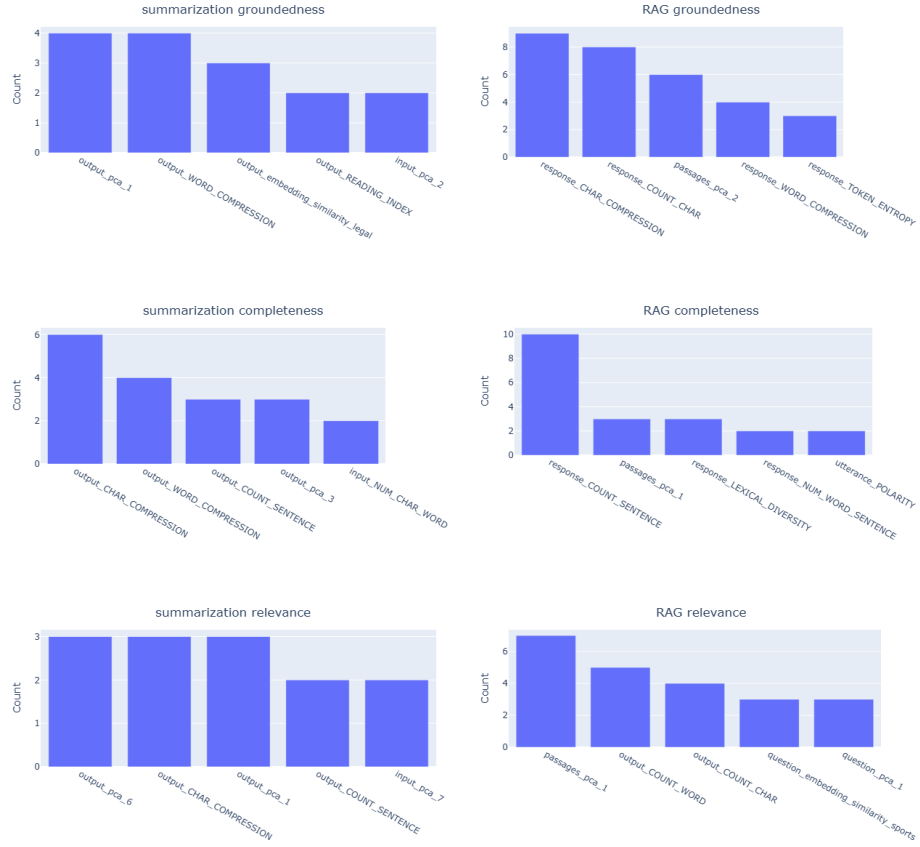


Figure 12: Aggregating the top five features that frequently appear across tasks. Character count is more prominent in RAG tasks, while compression ratio and embedding-related features such as PCA components and embedding similarity are more influential in summarization tasks.

fails to assess the completeness of a summary. The summary merely repeats the first sentence of the post and omits key details, yet Gemma incorrectly considers it complete and explains that it covers all key events.

```
df['question']

'botswana flag color meanings'

df['response']

"Sure! Here's my answer based on the given passages:\n\nThe Botswana flag has three horizontal s
tripes - blue, white, and black. The blue color represents water, specifically the Okavango Delt
a, which is a vital source of life and livelihood for many Batswana (the people of Botswana). Th
e white stripe symbolizes peace and unity among the different ethnic groups in Botswana, and is
also said to represent the salt pans found in the country. The black stripe stands for the ferti
le soil and the natural resources of the country.\n\nTherefore, the blue, white, and black color
s on the Botswana flag represent water, peace and unity, and natural resources, respectively."
```

```
df['context']

"passage 1:Botswana Flag: The Botswana flag was officially adopted on September 30, 1966. The co
lors on the flag correspond to those on the national coat of arms. The blue represents water, th
e white-black-white bands depict the racial harmony of the people as well as the pluralist natur
e of the society. They are inspired by the coat of the zebra, the national animal. Larger Botswa
na flag Botswana Coat of Arms: Botswana's coat of arms was adopted on January 25, 1966, and cont
ains a shield supported by two zebras. The zebras hold onto a tusk of ivory (symbolic of the for
mer ivory trade) and an ear of sorghum (an important crop to the nation.\n\npassage 2:Botswana f
lag colors meaning, significance and symbolism - Presidential Standard of Botswana - The meaning
of Botswana flag and What is the symbolism of the flag of Botswana Botswana flag colors symbolis
m\n\npassage 3:19. Botswana. Botswana used the UK flag as a protectorate till their independence
in September 30, 1966. The flag is credited for being one of the African flags that did not take
its meaning from the leading political party. Instead the flag designs and colors give cultural,
political, and regional interpretations.\n\n"
```

```
df['GEMINI_2_0_FLASH_groundedness_explanation']

'The response accurately reflects the information provided in the context. The colors of the fla
g (blue, white, and black) and their meanings (water, peace and unity, and natural resources) ar
e all mentioned in the context.\n'
```

Figure 13: Gemini 2.0 Flash fails to identify ungrounded content in the response—for instance, Okavango Delta, which is not mentioned in the cited context.

## E.4 Explanations on Embedding PCA-Related Features

In this section, we provide insights into the embedding PCA-related features. For illustration, we use the summarization completeness task as an example. Fig. 15 (left) shows the data distribution after applying K-means clustering based on the first two principal components (PCA 1 and PCA 2) of the embedding features. Fig. 15 (right) displays the same data distribution, colored by source dataset. It is evident that under PCA 1 and PCA 2, the TLDR dataset is clearly separated from the other two datasets, indicating that these embedding features capture meaningful dataset-level information. A likely explanation for this separation lies in differences in topic and writing style. TLDR samples originate from Reddit posts, which typically focus on personal or emotional dilemmas and are written in a less formal style. In contrast, source contexts from SummEval and UniSumm consist of more formally written news articles and reports. This observation is further supported by topic modeling using La-

```
df['context']

{'post': "Hi Reddit. \nMe (M 23)\nHer (F 19)\nBeen together for about 6 months.\n\nI'll cut st
raight to the chase. Yesterday I dropped my girlfriend at the train station so she could go to Melbourn
e to do some shopping. I received a text message from her a few hours later that she was on the train h
ome and that her aunty would be picking her up and she would be staying at her cousins house as she had
had a fight with her mum and didn't want to go home. She said she was in bed and was going to sleep. I
said that was fine. This morning, I found that my phone deleted most of my contacts during the night, i
ncluding my girlfriends number. I messaged her cousin this morning and told her to tell my girlfriend t
o message me when she woke up. I then found out that my girlfriend wasn't there, and didn't stay ther
e the night at all. \n\nI messaged my girlfriend asking what the fuck was going on and where she real
ly stayed. She then told me that she stayed at another cousins house in Melbourne since her mum had tol
d her to find somewhere else to live and that she didn't want me to worry about it. (I should add that
her ex lives in Melbourne, who she still talks to and has a kid with).\n\nI don't trust her, and I h
ave no way of knowing whether she's telling the truth or not. She's lied to me before. Should I get o
ut now? If she really has been kicked out of home, I hate the thought of adding to her troubles by leav
ing her. I've never broken up with someone, and I hate the thought of hurting someone.", 'title':
'Girlfriend lied to me.. not sure what to do', 'subreddit': 'relationships', 'site': None, 'ar
ticle': None}

df['summary']

' Hi Reddit. Me (M 23)\nHer (F 19)\nBeen together for about 6 months.'

df['Gemma-3-27B-IT_explanation']

"Explanation: The summary accurately captures all key events and the core dilemma presented in the pos
t, including the initial lie, the subsequent explanation, the reason for distrust, and the poster's int
ernal conflict."
```

Figure 14: Gemma 3 27B IT fails to assess the completeness of a summary. The summary merely repeats the first sentence of the post and omits key details, yet Gemma incorrectly considers it complete and explains that it covers all key events.

tent Dirichlet Allocation (LDA) (Blei et al., 2003). Below, we summarize the top words and inferred topics for each dataset:

- SummEval: Top words include club, hull, liverpool, time, and claim, suggesting topics related to sports and crime.
- UniSumm: Top words include text, rate, city, and officer, indicating topics such as technical instructions and city-related events.
- TLDR: Top words include want, boyfriend, girlfriend, love, and problem, reflecting themes of romantic desires and relationship issues.



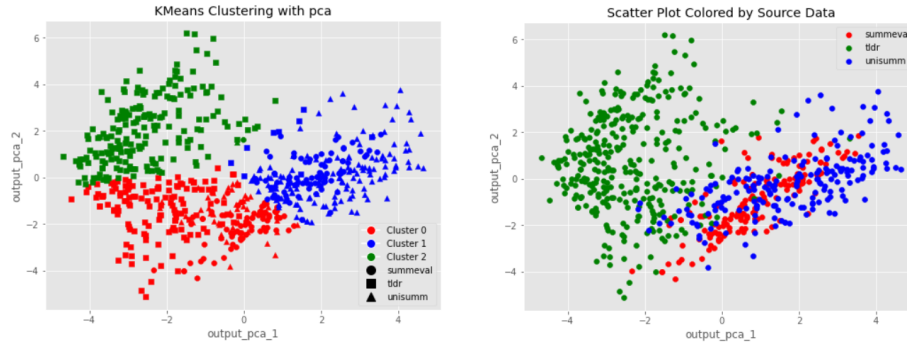


Figure 15: The left figure shows the data distribution after applying K-means clustering based on the first two principal components (PCA 1 and PCA 2) of the embedding features. The right figure displays the same data distribution, colored by source dataset. It is evident that under PCA 1 and PCA 2, the TLDR dataset is clearly separated from the other two datasets.

## E.5 Analysis of Jury Failure in Evaluation Tasks

Figure 16 presents the complete results of binning analyses of top features for the RAG groundedness task. The results reveal a clear trend: as the length or complexity of generated text increase, the jury’s ability to reliably assess groundedness decreases. It is hard to interpret the meaning of PCA features of generated text but a clear pattern is also shown.

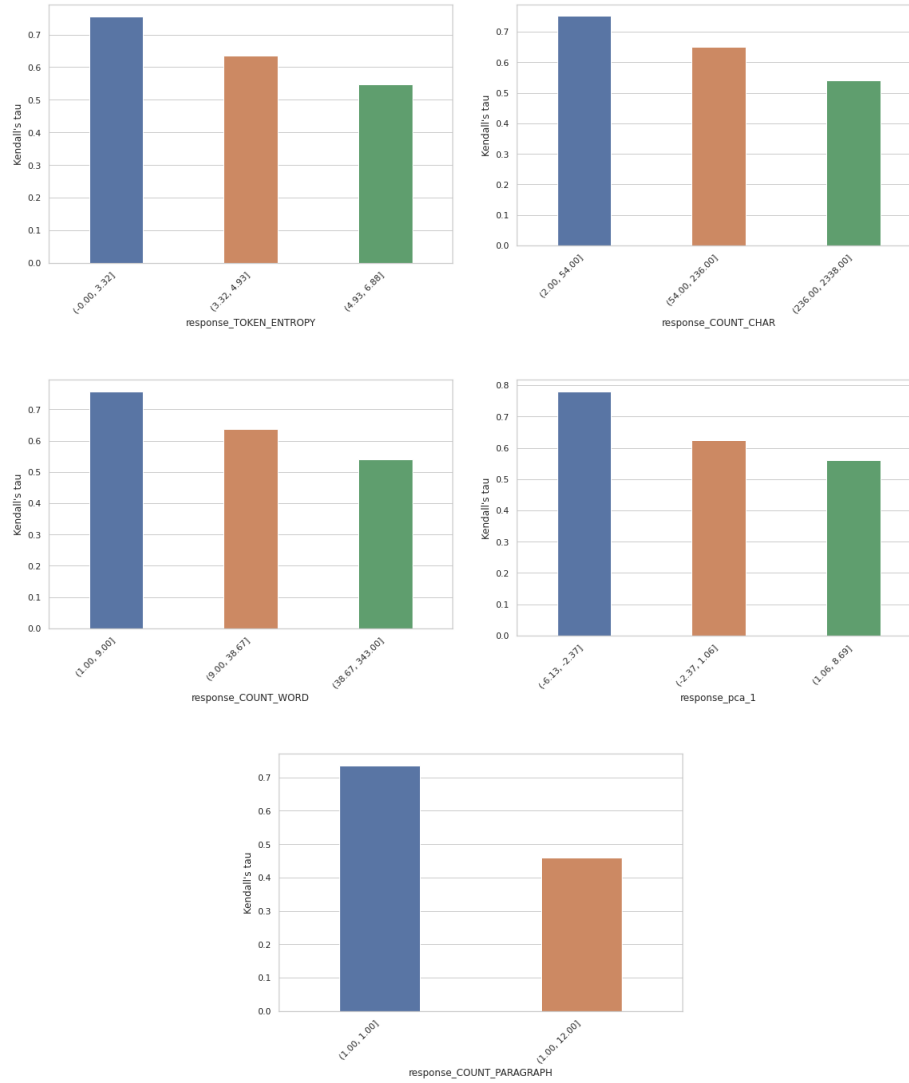


Figure 16: Jury performance across segments of top features indentified in the RAG groundedness task. Jury has better outcomes when the generated text is shorter or less complex.

## F Human Evaluation Consistency

We examine the consistency of human evaluations and compare them with the performance of our jury model. Both the SummEval and DialSumm datasets include three human annotators per evaluation dimension. Tables 18 to 21 report Kendall’s tau correlations between each individual annotation and the average of the three scores, which serves as the reference in our analysis. The top-left cell in each table presents the Kendall’s tau of the jury model. Our findings reveal that inter-annotator agreement is generally low, highlighting notable discrepancies among human judgments. While Kendall’s tau increases when comparing individual annotator scores to the average, substantial variation remains. For instance, in the SummEval Relevance dimension, the highest tau between a human annotator and the average is 0.875, whereas the lowest is only 0.411. Although the jury model rarely surpasses the best-performing annotator, it consistently outperforms the weaker ones. This suggests that the jury model offers a reliable and robust alternative to individual human evaluations.

Table 18: Kendall’s tau for human annotations – SummEval Groundedness

0.576	<b>Average score</b>	<b>Annotator 1</b>	<b>Annotator 2</b>	<b>Annotator 3</b>
<b>Average score</b>	1	0.893	0.893	0.735
<b>Annotator 1</b>	0.893	1	0.748	0.793
<b>Annotator 2</b>	0.893	0.748	1	0.807
<b>Annotator 3</b>	0.735	0.793	0.807	1

Table 19: Kendall’s tau for human annotations – SummEval Relevance.

0.696	<b>Average score</b>	<b>Annotator 1</b>	<b>Annotator 2</b>	<b>Annotator 3</b>
<b>Average score</b>	1	0.667	0.875	0.411
<b>Annotator 1</b>	0.667	1	0.455	0.388
<b>Annotator 2</b>	0.875	0.455	1	0.394
<b>Annotator 3</b>	0.411	0.388	0.394	1

## G Judge Reliability Prediction Model Performance

Table 22 presents the AUC scores of ROC curves for each judge reliability model trained using XGBoost on the testing set. The AUC values are relatively consistent across tasks; we include results for summarization completeness and RAG groundedness as representative examples. Result varies across judges and

Table 20: Kendall’s tau for human annotations – DialSumm Groundedness.

0.699	Average score	Annotator 1	Annotator 2	Annotator 3
<b>Average score</b>	1	0.647	0.712	0.679
<b>Annotator 1</b>	0.647	1	0.462	0.379
<b>Annotator 2</b>	0.712	0.463	1	0.351
<b>Annotator 3</b>	0.679	0.379	0.351	1

Table 21: Kendall’s tau for human annotations – DialSumm Relevance.

0.639	Average score	Annotator 1	Annotator 2	Annotator 3
<b>Average score</b>	1	0.741	0.758	0.622
<b>Annotator 1</b>	0.741	1	0.564	0.365
<b>Annotator 2</b>	0.758	0.564	1	0.344
<b>Annotator 3</b>	0.622	0.365	0.344	1

tasks, with most AUCs ranging between 0.63 and 0.78, indicating that the models demonstrate adequate predictive capability for these evaluation tasks.

Table 22: AUC of ROC for the judge reliability model. Here Gemn. is Gemini, DS is DeepSeek.

Metric	Claude 3.7 SON- NET	Gemn. 2.0 Flash	Gemn. 2.5 Flash	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LLAMAGemma -3.2- 3B In- struct	Phi- 3- 27B- IT	4- reason - ing	DS- R1
Summ - completeness	0.63	0.67	0.72	0.66	0.76	0.75	0.63	0.63	0.62	0.68
RAG - groundedness	0.70	0.78	0.67	0.68	0.76	0.73	0.61	0.77	0.75	0.73

## H Ablation Study and Model Weakness Analysis

### H.1 Ablation Study (Data Properties)

To investigate the influence of data property features on jury performance, we conduct ablation studies by selectively removing feature sets during model construction. As illustrative examples, we focus on two tasks: summarization completeness and RAG groundedness. The full categorization of features is provided in Appendix A. Specifically, we remove three groups of features in separate ex-

periments: (1) text size-related features and special word count features (jointly removed due to their high correlation), (2) text complexity features, and (3) embedding-based features. The results are presented in Tables 23 and 24 for summarization completeness and RAG groundedness, respectively. We observe that the jury model achieves its best performance when all feature sets are included, underscoring the importance of comprehensive feature representation. Although the performance differences are modest, this can be attributed to the internal correlations within each feature category. Additionally, different tasks exhibit varying sensitivity to feature sets. For instance, removing embedding features leads to a greater performance drop in the summarization completeness task than in RAG groundedness, where text size-related features have a more pronounced impact.

Table 23: Kendall’s tau for jury performance on summarization completeness under feature ablation. The number of judges (k) is indicated for each configuration. The jury using all features achieves the highest performance, particularly on SummEval.

Data	Text size + special words removed (k=6)	Text complex- ity removed (k=5)	Embedding removed (k=6)	All features (k=7)
Overall	0.487	0.478	0.471	0.488
SummEval	0.639	0.658	0.688	0.721
TL;DR	0.352	0.459	0.318	0.427
UniSumEval	0.619	0.595	0.589	0.612

Table 24: Kendall’s tau for jury performance on RAG grounded under feature ablation. The number of judges (k) is indicated for each configuration. The jury using all features achieves the highest performance, particularly on RagTruth.

Data	Text size + special words removed (k=4)	Text complex- ity removed (k=3)	Embedding removed (k=3)	All features (k=3)
Overall	0.667	0.672	0.667	0.678
CAQA	0.659	0.659	0.676	0.651
HaluEval	0.761	0.759	0.766	0.798
RagTruth	0.558	0.596	0.575	0.576

## H.2 Ablation Study (Jury Size - $K$ )

We conduct experiments to test the effectiveness of varying jury size compared to keeping a fixed value. We focus on the tasks Summarization-Completeness and RAG-Groundedness for this study. The tolerance level for all trained XGBoost models is set to 0 which means only the exactly matching scores are considered as correct. Performance is measured across 10 runs and average is taken for each jury size. Fig. 17 shows that the performance varies with jury size. In the Summarization-Completeness task, performance steadily improves as jury size increases, reaching its peak around a jury size of 7–8 before slightly declining. This indicates that adding more judges generally strengthens the overall decision by reducing individual biases, but very large juries introduce diminishing returns and slight degradation. Similarly, in the RAG-Groundedness task, performance starts low at smaller jury sizes and improves significantly with larger juries, peaking at 5 and 8 and then tapering off. In both cases, increasing jury size enhances robustness against noisy predictions, but there is an optimal range beyond which gains flatten or reverse. This shows that tuning jury size per task still provides significant improvements compared to fixing it, though the optimal size tends to be moderately large for both tasks.

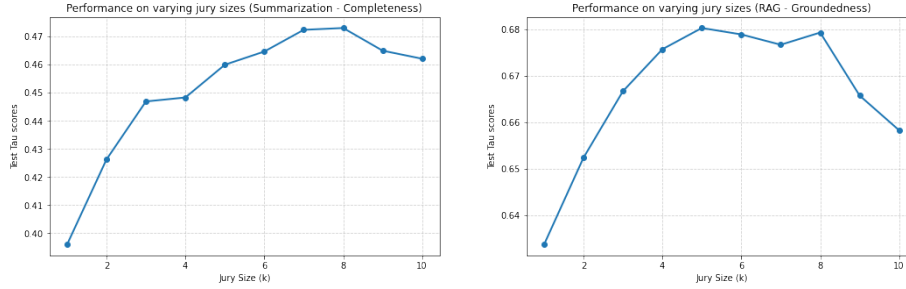


Figure 17: Test performance with varying jury size.

## H.3 Ablation Study (Tolerance levels - $\tau$ )

We consider two tolerance levels for summarization tasks (0 and 1) and a single tolerance level (0) for RAG tasks. For RAG tasks valid scores are 0-2. Thus, a tolerance level of 1 or more would mean that a score of 1 is always considered as correct. For the same reason, we do not experiment with a tolerance level of 2 or more for summarization tasks where the valid scores are 1-5.

For this study, we focus on Summarization-Completeness. We experiment with two tolerance levels: 0 and 1 in the original scale (1-5). After min-max normalization to  $[0, 1]$ , these correspond to 0 and 0.25. The performance of the jury is observed with all the XGBoost models trained either on tolerance of 0 or 0.25. Table 25 summarizes the means and standard deviations of the Kendall’s Tau across the 10 runs and shows the comparison with tuned tolerance models.

We observe that allowing different tolerance levels across different XGBoost models gives slightly better performance than a fixed tolerance level across all models.

Additionally, Table 25 further illustrates the importance of tolerance tuning. While Jury-on-Demand with variable tolerance achieves the best overall performance, the optimal fixed tolerance differs across datasets: TL;DR performs better with a tolerance of 0, whereas UniSumEval favors 0.25. This variability underscores that no single fixed tolerance can fit all datasets. In practical scenarios, especially for unseen datasets without human annotations, it is impossible to know the ideal tolerance beforehand. Therefore, adaptive approaches that allow tolerance to vary across models or instances are crucial for robust generalization.

Finally, the results with fixed tolerance levels are better than the static jury as we have chosen the best jury size ( $K = 7$ ) overall across the runs.

Table 25: Kendall’s tau for jury performance on Summarization-Completeness under tolerance ablation. The number of judges ( $K$ ) is fixed to 7 (overall best) for the fixed tolerance configurations.

Data	Fixed tolerance (0)	Fixed tolerance (0.25)	Jury on Demand (variable tolerance)	Static Jury
Overall	0.47 (0.02)	0.47 (0.00)	0.48 (0.03)	0.44 (0.02)
SummEval	0.69 (0.05)	0.75 (0.04)	0.72 (0.05)	0.60 (0.06)
TL;DR	0.38 (0.06)	0.37 (0.04)	0.38 (0.08)	0.40 (0.05)
UniSumEval	0.63 (0.04)	0.66 (0.05)	0.66 (0.04)	0.59 (0.03)

## H.4 Ablation Study (Prompt Variation)

We analyze the effect of using different prompts with the judges on the overall jury performance. We focus on the task Summarization-Completeness for our analysis. The performance with the prompt as described in Table 6 is compared against the performance with a slightly different prompt which omits the list of valid scores.

We observe that changing the prompt while maintaining its meaningfulness (the prompt should clearly explain the input and the expected output), changes the evaluation scores for roughly 30% of the data points by 1 score point (on average across datasets, judges and metrics). Fig. 18 shows the distribution of score differences for judging the Completeness of summaries using GEMINI 2.5 PRO as the judge.

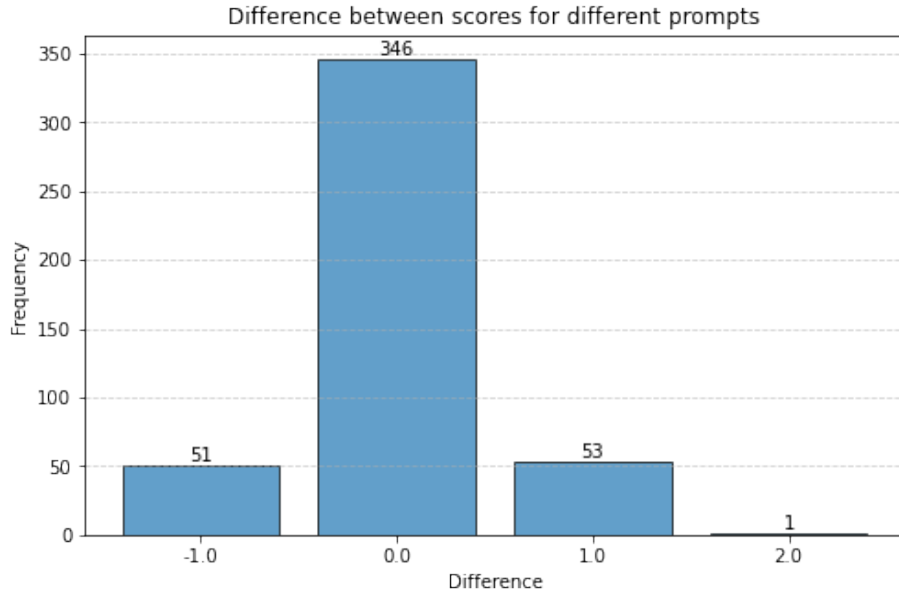


Figure 18: Score differences in the results of GEMINI 2.5 PRO judge with two different prompts

The XGBoost judge reliability models are then trained using the scores from the second prompt. We observe changes in the reliability scores across judges with the perturbed prompt. The distribution of the reliability score differences is shown in Fig. 19. It is seen that on average roughly 50% of the datapoints have less than 0.1 difference in the reliability scores of the judges.



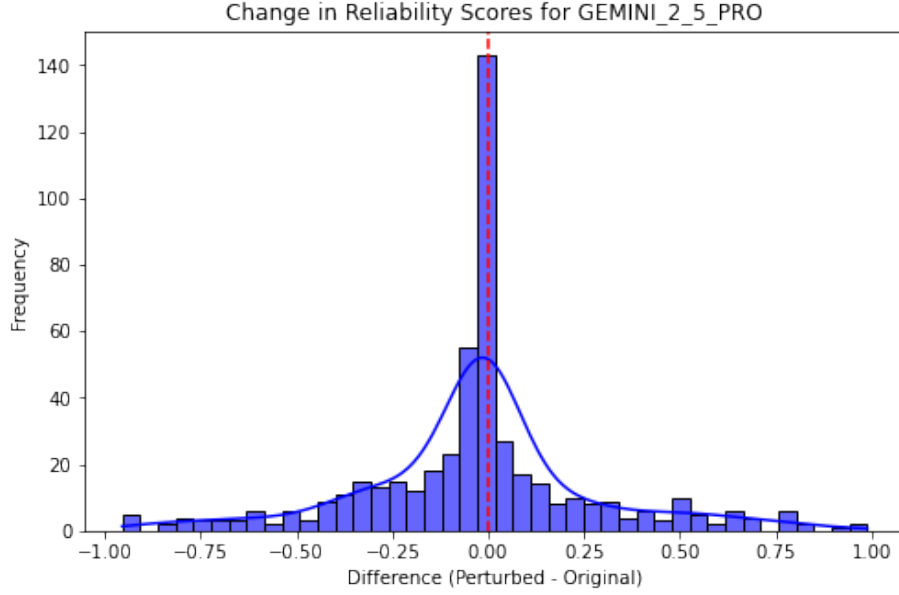


Figure 19: Predicted reliability score differences in the results of GEMINI 2.5 PRO judge with two different prompts

It is also observed that even though the scores might change by one score point, their agreement or disagreement with the human given scores generally remains consistent even with slightly different prompts. For example, in the case of Gemini 2.5 PRO as the judge, for 439/495 datapoints, the judge responses either both match the human score or both disagree. Table 26 summarizes the agreement changes with the human scores.

Table 26: Agreement changes with the human scores for two prompts on GEMINI 2.5 PRO judge

Measurement Metric	Both or Neither Match	Only matches one
Number of datapoints	439	56
<i>%tage</i> age of datapoints	88.7	11.3

We also observe how the number of appearances of the models in the selected jury panel changes with slight change in the prompts. Fig. 20 shows that there is limited change in the inclusion of individual judges in a given jury with the changed prompt.

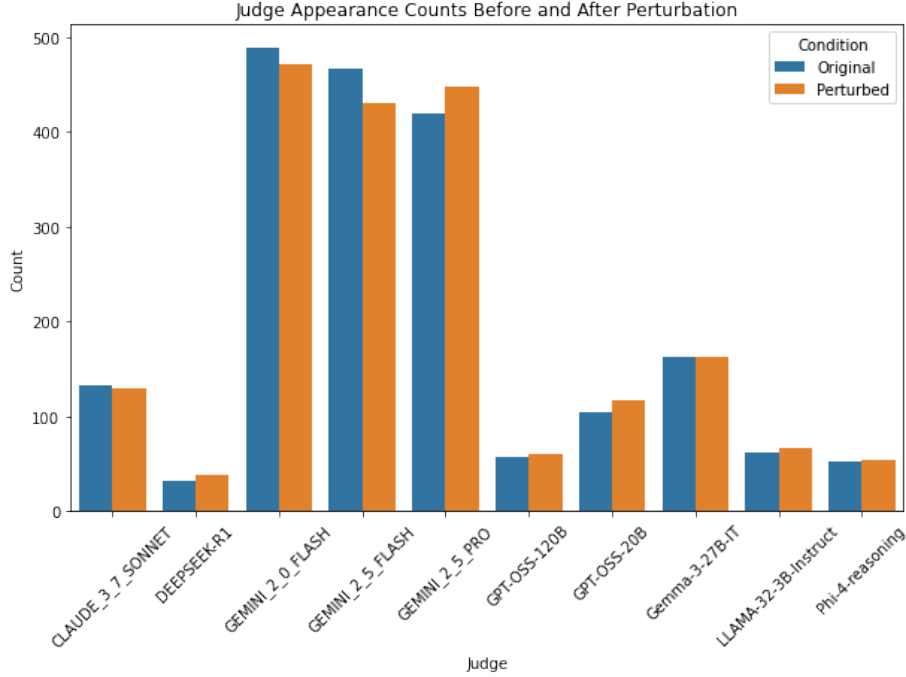


Figure 20: Number of appearances of various judges in the selected jury panel

We observe that on average, 3.5 out of 4 final judges in the jury remained the same on using different prompts.

Table 27: Distribution of the number of common judges between the two runs with different prompts

# of common judges	0	1	2	3	4
Counts	0	3	21	182	289

Table 27 shows that the jury compositions do not change much when using slightly different prompts.

The perturbation analysis demonstrates that the evaluation framework is highly robust to prompt variations. While minor prompt changes introduce some variability (about 30% of data points shift by one score point), the overall agreement with human evaluations remains largely stable. For instance, in the Gemini 2.5 PRO case, over 88% of instances (439/495) either consistently match or consistently diverge from human scores across prompt variations.

Reliability score changes exhibit minimal changes, with 50% of data points showing less than 0.1 difference, indicating that prompt-induced randomness is minimal and does not significantly affect judge reliability. Furthermore, jury composition stability is striking. On average, 3.5 out of 4 jury members remain

unchanged, and the aggregated jury scores show negligible deviation. This suggests that the methodology maintains integrity even under prompt perturbations.

The system’s resilience to prompt changes implies that the evaluation pipeline is not overly sensitive to linguistic nuances, making it suitable for real-world deployment where prompt variability is inevitable. This robustness ensures consistent decision-making and fairness in model evaluation, reinforcing confidence in the methodology for broader applications.

## H.5 Model Weakness Analysis

To identify conditions where the jury model underperforms, we analyzed feature distributions between instances with high and low prediction discrepancies - excluding embedding features for interpretability. We focus on the RAG groundedness task for illustration. The dataset was split into two groups:

1. High Difference: Instances with large prediction errors.
2. Low/No Difference: Instances with minimal errors.

For each feature, we computed mean and median values across both groups and ranked features by their median absolute differences in Table 28. This revealed two failure modes:

1. **Systematic Bias:** Features like factual density and named entities show consistent shifts in both mean and median, suggesting bias toward certain content structures.
2. **Distributional Fragility:** Features such as syntactic anomaly and subjectivity show high median shifts but low mean differences, indicating sensitivity to rare or irregular linguistic patterns.

Table 28: Top 5 features ranked by median differences between high-difference and low/no-difference groups.

Feature	Median Difference	Mean Difference
utterance_SYNTACTIC_ANOMALY	1.97	0.11
utterance_NAMED_ENTITY	0.77	0.37
utterance_SYNTACTIC_AMBIGUITY	0.74	0.22
context_FACTUAL_DENSISTY	0.63	0.36
context_SUBJECTIVITY	0.53	0.17

## I Judge Scoring Runtimes

Here we provide some notes the runtime for judges scoring select data. For this work, each judge was implemented in one of three distinct environments: Google Cloud API, Nvidia H200 with 140 GB RAM, or Nvidia TESLA V100 with 32 GB RAM. These choices were made based on the nature of the models (closed vs. open) and resource needs of the the open-weight models.

- GCP was selected for the closed models such as Claude 3.7 Sonnet and Gemini 2.5.
- The Nvidia H200 was used for larger open-weight models like Phi-4 Reasoning, Llama 3.2-3B, and Gemma 3 27B-IT, leveraging its high memory capacity for large-scale inference.
- V100 served as a baseline GPU environment for evaluating performance under constrained resources.

The runtime for each dataset varied depending on several factors including the size of the dataset, the computational environment used, and the specific model being executed. Table 29 summarizes the runtimes observed for select datasets during our experiments. Note that this is not intended as to represent a rigorous comparative study of the models but to provide a general sense of the time necessary for executing judge scoring.

LLM	Environment	TLDR	unisumeval	MS MARCO	QASPER
Claude 3.7 SONNET	GCP	3 mins	4 mins	7 mins	5 mins
Gemini 2.5 Pro		6 mins	6 mins	26 mins	2 mins
Gemini 2.5 Flash		2 mins	3 mins	13 mins	2 mins
Phi-4 Reasoning	H200	81 mins	116 mins	579 mins	36 mins
Llama 3.2-3B		9 mins	68 mins	136 mins	17 mins
Gemma 3 27B-IT		17 mins	66 mins	477 mins	17 mins
DeepSeek-R1	V100	76 mins	401 mins	721 mins	18 mins

Table 29: Observed model runtimes (in minutes) across different environments and datasets used in the analysis

## J Framework Generalizability to Unseen Domains

To assess the framework’s generalizability to unseen domains, we employ a leave-one-out procedure. Specifically, for each experiment, one data source is excluded from the training of XGBoost reliability models and jury construction, and the trained framework is then evaluated on the held-out source. This approach tests

whether the Jury-on-Demand mechanism consistently outperforms both static juries and individual judges in previously unseen domains.

Using the RAG-Relevance task as an illustrative example, the dataset comprises three sources: ALCE, Hotpot-QA, and MS MARCO. In the first iteration, ALCE is held out while the framework is trained on Hotpot-QA and MS MARCO; performance is then assessed on ALCE. The process is repeated for each remaining source. The complete results are presented in Table 30. Across all three cases, Jury-on-Demand achieves the highest performance, indicating that the learned patterns generalize effectively to held-out domains.

Table 30: Kendall’s tau on held-out data source - RAG Relevance

Held-out	Jury-on-Demand	Static-Jury	Claude 3.7	Gemn. 2.0	Gemn. 2.5	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	GemmaPhi-3	Phi-4	DeepSeek-R1
ALCE	0.62	0.59	0.58	0.55	0.57	0.57	0.58	0.6	0.22	0.42	0.13	0.28
Hotpot-QA	0.92	0.86	0.87	0.81	0.88	0.88	0.89	0.89	0.39	0.78	0.47	0.57
MS-MARCO	0.46	0.38	0.44	0.38	0.4	0.43	0.39	0.4	0.2	0.12	0.12	0.29

We extend our evaluation to the Summarization-Groundedness task, which includes four data sources: Summ-Eval, TL;DR, UniSum-Eval, and Dial-Summ-Eval. Table 31 reports the performance of the Jury-on-Demand framework under the leave-one-out setting for each source. The results indicate strong generalization for three sources, while performance on Summ-Eval is weaker.

These findings reinforce that the framework’s generalizability is influenced by both the diversity of training data and the characteristics of unseen domains. As additional annotated datasets become available and incorporated into training, we expect the framework’s ability to generalize to new domains to improve substantially.

## K Judge Score Calibration

Some judges may consistently assign lower or higher scores compared to human annotations. To address this, we apply score calibration. Specifically, we perform isotonic calibration (Niculescu-Mizil & Caruana, 2005) for each judge’s score within each dataset, then retrain the XGBoost model and construct the jury using the calibrated scores. For illustration, we focus on summarization completeness and RAG groundedness. Prior to calibration, we examine the difference between each judge’s mean score and the mean human annotation score (annotation score minus judge score). Fig. 21 show these differences for completeness in Unisumm and groundedness in RagTruth. We observe that weaker models—particularly LLAMA 3.2 3B Instruct, Gemma 3.2 7B IT, Phi

Table 31: Kendall’s tau on held-out data source - Summarization Groundedness

Held-out	Jury-on-Demand	Static-Jury	Claude 3.7	Gemn. 2.0	Gemn. 2.5	Gemn. 2.5 Pro	GPT-OSS-20B	GPT-OSS-120B	LL-3.2	GemmaPhi-3	Phi-4	DeepSeek-R1
Dial-Summ-Eval	0.7	0.64	0.63	0.63	0.69	0.65	0.66	0.68	0.3	0.69	0.26	0.28
Summ-Eval	0.63	0.68	0.63	0.54	0.68	0.66	0.64	0.68	0.33	0.63	0.23	0.02
TL;DR	0.46	0.43	0.38	0.43	0.32	0.33	0.27	0.41	0.06	0.42	0.12	0.14
Uni-Summ-Eval	0.63	0.53	0.52	0.56	0.59	0.58	0.6	0.59	0.15	0.58	0.17	-0.1

4 Reasoning and DeepSeek R1—tend to exhibit larger discrepancies.

Tables 33 and 32 compare Kendall’s tau before and after calibration. Overall performance changes are minimal, with calibrated results slightly improving on RAG groundedness (0.69 vs. 0.67) but slightly worsening on summarization completeness (0.46 vs. 0.49). Differences become more pronounced within certain datasets and for specific judges. For example, as shown in Fig. 22 (left), tau for Gemini 2.5 Flash on Unisumm completeness drops from 0.62 to 0.54 because many original score 5s are calibrated to 4, leading to underestimation of human annotation scores. Conversely, for RagTruth groundedness (Fig. 22 (right)), Gemini 2.5 Flash’s tau increases from 0.57 to 0.58 after calibration because scores of 0 are calibrated to 1, reducing bias in judge scores. In summary, calibration can be beneficial for certain tasks and judges, but it may also introduce under- or overestimation of human annotations, reducing alignment. Future work will explore strategies to mitigate judge score bias more effectively.

## L Comparison with Baselines

In this section, we examine how performance varies with the number of judges  $K$ , we focused on two tasks for illustration: Summarization-Completeness and RAG-Groundedness. We compared Jury-on-Demand (JOD) with two static jury baselines: Simple Average and Weighted Regression (see section 4.2 for definitions). For these static jury baselines, we first determine which judges to include. Specifically, we rank judges by their Kendall’s tau performance for each task and select the top  $K$  judges for the static jury. Table 34 shows the judge ranked by performance for each task. Tables 35 and 35 illustrate how Kendall’s tau changes with jury size  $K$  for Summarization-Completeness and RAG-Groundedness, respectively. Results are averaged over 10 runs. Our main observations are:

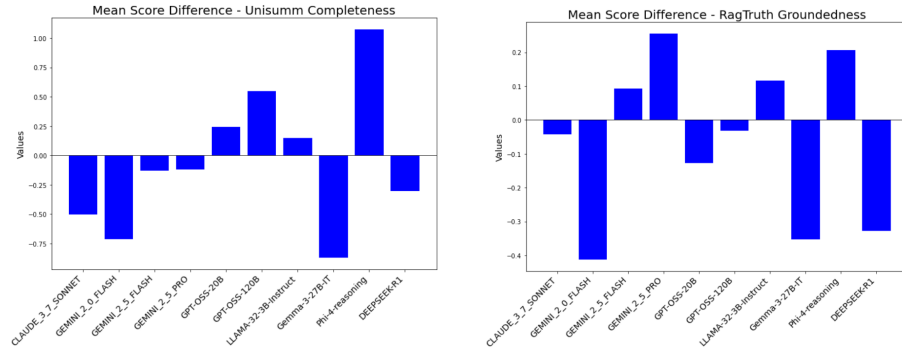


Figure 21: The difference between each judge’s mean score and the mean human annotation score (annotation score minus judge score) for completeness and groundedness in Unisumm.

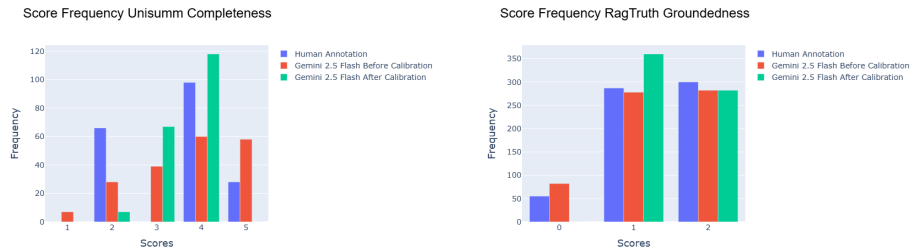


Figure 22: The difference between each judge’s mean score and the mean human annotation score (annotation score minus judge score) for completeness and groundedness in Unisumm.

Table 32: Summarization completeness: Kendall’s tau after calibration, with original values shown in parentheses. NA indicates cases where all scores were calibrated to a single value, making Kendall’s tau computation infeasible.

Model	Jury- on- Demand	Claude 3.7	Gemn. 2.0	Gemn. 2.5	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LL- 3.2	Gma. 3	Phi- 4	DS- R1
Overall	0.46 (0.49)	0.38 (0.39)	0.47 (0.46)	0.51 (0.46)	0.44 (0.43)	0.38 (0.39)	0.46 (0.47)	0.15 (0.13)	0.22 (0.25)	0.16 (0.22)	0.23 (0.27)
Summ- Eval	0.63 (0.72)	0.64 (0.59)	0.61 (0.62)	0.58 (0.60)	0.47 (0.46)	NA (- 0.02)	0.62 (0.63)	- (- 0.04)	0.07 (0.24)	0.31 (0.21)	0.69 (0.69)
TL;DR	0.35 (0.43)	0.31 (0.32)	0.41 (0.40)	0.36 (0.37)	0.35 (0.37)	0.29 (0.34)	0.39 (0.39)	0.06 (0.12)	0.23 (0.17)	NA (0.02)	0.12 (0.12)
Uni- Summ- Eval	0.58 (0.61)	0.41 (0.57)	0.46 (0.58)	0.54 (0.62)	0.46 (0.58)	0.66 (0.62)	0.61 (0.71)	0.26 (0.26)	0.54 (0.64)	0.33 (0.34)	0.35 (0.38)

Table 33: RAG groundedness: Kendall’s tau after calibration, with original values shown in parentheses. NA indicates cases where all scores were calibrated to a single value, making Kendall’s tau computation infeasible.

Model	Jury- on- Demand	Claude 3.7	Gemn. 2.0	Gemn. 2.5	Gemn. 2.5 Pro	GPT- OSS- 20B	GPT- OSS- 120B	LL- 3.2	Gma. 3	Phi- 4	DS- R1
Overall	0.69 (0.67)	0.56 (0.53)	0.50 (0.48)	0.64 (0.61)	0.60 (0.56)	0.63 (0.61)	0.63 (0.61)	NA (0.09)	0.10 (0.02)	0.03 (0.13)	0.28 (0.25)
CAQA	0.65 (0.65)	0.57 (0.57)	0.57 (0.58)	0.58 (0.57)	0.57 (0.57)	0.57 (0.57)	0.59 (0.58)	NA (0.1)	NA (0.01)	NA (0.03)	0.19 (0.23)
Halu- Eval	0.79 (0.79)	0.72 (0.70)	0.52 (0.51)	0.79 (0.77)	0.76 (0.76)	0.78 (0.78)	0.79 (0.79)	NA (0.14)	0.14 (0.14)	- (0.11) 0.04	0.39 (0.39)
Rag- Truth	0.55 (0.57)	0.48 (0.44)	NA (0.32)	0.58 (0.57)	0.50 (0.47)	0.55 (0.52)	0.53 (0.51)	NA (0.13)	NA (0.16)	NA (0.21)	NA (0.24)

- JOD consistently outperforms static juries: Performance varies with jury



size  $K$ , but JOD consistently outperforms both static baselines across all  $K$  values, except for a few cases in Summarization-Completeness. The performance margin of JOD over static juries is larger for RAG-Groundedness than for Summarization-Completeness.

- Effect of jury size differs by task: For Summarization-Completeness, there is no clear pattern on performance do not for different juries. JOD’s performance do not vary too much as  $K$  changes, Simple Average jury show performance decrease as  $K$  increases. For RAG-Groundedness, performance generally improves with larger  $K$ , except that Simple Average shows a drop at  $K = 10$ . These observations can be explained by the larger score bias in weaker judges. Large biases degrade performance when these judges are included, especially for Simple Average, which assigns equal weight regardless of bias.
- Weighted Regression struggles with small  $K$ : Weighted Regression performs worse than other jury methods when  $K$  is small, particularly for RAG-Groundedness. One contributing factor is that the regression coefficients for RAG-Groundedness are very small, which results in weighted regression scores that fall below human annotation scores. To ensure realistic outputs, we round up any weighted regression score that is lower than the minimum human annotation score. However, because the coefficients are small, this adjustment often causes many weighted regression scores to equal the lowest human annotation score, further degrading performance. Another potential reason for poor performance at small  $K$  is the high correlation among strong judges (e.g., Gemini models, GPT-OSS models, Claude 3.7 Sonnet), with correlation values around 0.7–0.8. Using such highly correlated features in regression can negatively impact model performance.

In summary, how performance changes with jury size  $K$  depends on multiple factors, including task characteristics, score distributions, and correlations among judge scores. Although JOD requires more training data and a more complex training process compared to static juries, its dynamic jury selection and adaptive weight assignment enable it to choose the best judges for each instance and achieve superior overall performance.

Table 34: Judges ranked by Kendall’s tau (high to low).

<b>K</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>
Summ- complete- ness	GPT- OSS- 120B	Gemn- 2.5	Gemn- 2.0	Gemn- 2.5- Pro	Claude	GPT- OSS- 20B	DS- R1	Gemma- 3	Phi 4	LL 3.2
RAG- grounded- ness	GPT- OSS- 120B	Gemn- 2.5	GPT- OSS- 20B	Gemn- 2.5- Pro	Claude	Gemn- 2.0	DS- R1	Phi 4	LL 3.2	Gemma- 3

Table 35: Summarization Completeness: Kendall’s tau across varying jury sizes for different baselines.

<b>K</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>
Jury-on-Demand	0.43	0.45	0.45	0.46	0.46	0.47	0.47	0.46	0.46
Static (Average K)	0.46	0.47	0.48	0.48	0.46	0.45	0.45	0.44	0.43
Static (Average K)	0.46	0.45	0.45	0.46	0.46	0.46	0.46	0.44	0.45

Table 36: RAG Groundedness: Kendall’s tau across varying jury sizes for different baselines.

<b>K</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>
Jury-on-Demand	0.65	0.67	0.67	0.68	0.68	0.67	0.68	0.67	0.66
Static (Average K)	0.64	0.65	0.65	0.65	0.65	0.64	0.64	0.63	0.58
Static (Average K)	0.58	0.58	0.56	0.54	0.52	0.60	0.65	0.65	0.65

## M Parameter Tuning in XGBoost

We use random search to tune the hyperparameters of the XGBoost models for judge reliability scores. The search space for the tuned parameters is provided in Table 37. Parameters not listed in the table are set to their default values.

Table 37: Parameter search space in XGBoost.

<b>Parameter</b>	<b>Search Space</b>
max_depth	2, 3, 4, 5, 6, 7, 8, 9
learning_rate	0.01, 0.03, 0.05, 0.07, 0.1, 0.2
n_estimators	500, 600, 800, 1000, 1200