
LLM-as-Judge on a Budget

Aadirupa Saha¹
UIC

Aniket Wagde
UIC

Branislav Kveton
Adobe Research

Abstract

LLM-as-a-judge has emerged as a cornerstone technique for evaluating large language models by leveraging LLM reasoning to score prompt-response pairs. Since LLM judgments are stochastic, practitioners commonly query each pair multiple times to estimate mean scores accurately. This raises a critical challenge: given a fixed computational budget B , how to optimally allocate queries across K prompt-response pairs to minimize estimation error? We present a principled variance-adaptive approach leveraging multi-armed bandit theory and concentration inequalities. Our method dynamically allocates queries based on estimated score variances, focusing resources where uncertainty is highest. Our algorithm is shown to achieve a worst-case score-estimation error of $\tilde{O}\left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B}}\right)$, σ_i^2 being the unknown score variance for pair $i \in [K]$ with near-optimal budget allocation. Experiments on HelpSteer2 dataset demonstrate our method significantly outperforms uniform allocation, reducing worst-case estimation error given a fixed budget. Our work establishes a theoretical foundation for efficient LLM evaluation with practical implications for AI safety, model alignment, and automated assessment at scale.

1 INTRODUCTION

Large language models (LLMs) have revolutionized AI evaluation. Instead of relying solely on expensive human annotators or rigid automated metrics, practitioners now employ LLM-as-a-judge [Zheng et al., 2024,

¹ Corresponding author email: aadirupa.saha@gmail.com

Dubois et al., 2024]—querying powerful models like GPT-4 or Claude to assess response quality. Consider a practical scenario of evaluating 10,000 prompt-response pairs for model fine-tuning. Human evaluation at \$5 per judgment would cost \$50,000, which may be prohibitively expensive. LLM judges offer a compelling alternative at \$0.01 or less per evaluation, enabling scalable quality assessment for supervised fine-tuning, post-training alignment, prompt optimization, A/B testing, and model calibration [Ouyang et al., 2022, Bai et al., 2022].

LLMs are now widely deployed not just for generating content but also for evaluating it, an approach referred to as LLM-as-a-judge [Kim et al., 2024, Trivedi et al., 2024, Gu et al., 2025, Zhang et al., 2024b, Thakur et al., 2025, Zhu et al., 2025]. This method leverages the natural language understanding of LLMs to replicate human-like evaluations, providing a scalable, economical, and consistent alternative to human annotation. An LLM judge can produce various forms of feedback: natural language critiques, numerical scores, or comparative preferences between options. Often, both textual assessments (rationale) and numeric ratings (scores) are generated together, and the score variation directly depends on the rationale [Kim et al., 2024].

However, LLM judgments are inherently stochastic. Querying the same (prompt, response) pairs multiple times yields different scores due to sampling randomness and temperature settings. This variance is highly heterogeneous: a factual question like “*What is 2 + 2?*” might produce consistent scores (say, with score variance $\sigma^2 \approx 0.001$), while subjective queries like “*How many stars would you give to movie Superman?*” could generate highly variable scores ($\sigma^2 \approx 1$). For the first pair, a single query suffices; for the second, many samples are needed to get a reliable estimate. Yet each LLM query incurs cost and consumes computational resources. This raises a fundamental resource allocation problem: *given a fixed budget of B queries across K prompt-response pairs, how should we distribute queries to estimate all scores accurately?* Inaccurate judgments and scores of LLM judges lead to poor alignment with human evaluations [Chiang and Lee, 2023, Gehrmann

et al., 2023].

Existing work on LLM evaluation has made significant progress on complementary problems. Researchers have developed better prompting strategies [Zheng et al., 2024, Liu et al., 2023], mitigated position and length biases [Wang et al., 2023, Zheng et al., 2024], and benchmarked judge-human agreement [Liu et al., 2023, Dubois et al., 2024]. However, these efforts universally assume *uniform allocation*—querying each pair the same number of times—or ignore the allocation problem entirely. Recent work on LLM uncertainty [Kuhn et al., 2023, Xiong et al., 2023] focuses on confidence calibration for individual predictions, not budget-constrained multi-instance evaluation. Reward modeling literature [Ouyang et al., 2022, Wang et al., 2024] collects preference data but does not address optimal query allocation given variance heterogeneity. To our knowledge, *no prior work addresses optimal query allocation under fixed budgets for confidence estimation in LLM judges*. This is a critical gap! A classic baseline like *Uniform Allocation* is provably suboptimal when variances differ, wasting queries on easy pairs while under-sampling difficult ones. We fill this gap by formalizing LLM judge evaluation as a variance-adaptive resource allocation problem through a multi-armed bandit (MAB) lens.

Informal Problem Statement. Given K (prompt, response) pairs $\{(q_i, a_i)\}_{i=1}^K$, each associated with a true unknown score $s_i \in \mathbb{R}$, and a fixed query budget B , the goal is to estimate the scores \hat{s}_i for all $i \in [K]$. The objective is to minimize $\max_{i \in [K]} |s_i - \hat{s}_i|$, subject to the budget constraint B . In other words, we seek to allocate the total budget B such that the worst-case estimation error $\max_{i \in [K]} |s_i - \hat{s}_i|$ is minimized (details in Sec. 2).

Contributions of this work are multifold:

- **(C1) Problem Formulation.** Our first contribution lies in formally modeling the LLM judge evaluation as a multi-armed bandit problem where each prompt-response pair is an arm with unknown mean score s_i and variance σ_i^2 , and the goal is to minimize worst-case estimation error (WEC) $\max_{i \in [K]} |s_i - \hat{s}_i|$ under budget constraint B (Sec. 2).
- **(C2) Algorithm Design (Secs. 3 and 4).** We develop variance-adaptive allocation strategies for both known and unknown variance settings. For known variances, ROBIN (Alg. 1) implements a greedy allocation rule that sequentially selects arms based on their respective standard error, i.e. $\arg \max_{i \in [K]} \sigma_i^2/n_i(t)$, where $n_i(t)$ is the number of times arm i is pulled till time t , *naturally balancing high-variance arms with under-sampled arms*. For the realistic unknown-variance setting, ROBIN-HOOD (Alg. 2) proceeds in two phases: initial uniform exploration over $t_0 K$ queries to estimate variances, $\{\hat{\sigma}_i^2\}_{i=1}^K$ (Eq. (5)), up to a certain accuracy. The next phase follows a similar adaptive greedy allocation to known variance, except that we replace the true variance σ_i at time t with its upper confidence bound $\bar{V}_i(t)$ (Eq. (6)). Both algorithms maintain $O(K)$ computational complexity per query with efficient online updates, making them practical for large-scale evaluations.
- **(C3) Theoretical Analysis (Secs. 3.2 and 4.2).** We establish near-optimal sample complexity guarantees for both algorithms. For ROBIN with known variances, Thm. 3 proves that with probability at least $1 - \delta$, the worst-case error satisfies $\max_{i \in [K]} |s_i - \hat{s}_i| = \tilde{O}\left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B}}\right)$, where $\tilde{O}(\cdot)$ hides logarithmic factors. For ROBIN-HOOD with unknown variances, Thm. 5 establishes the same rate with additional logarithmic overhead in the exploration phase, showing that variance estimation incurs negligible cost when $t_0 = \Theta(\log K)$. Our analysis leverages Bernstein-type sub-Gaussianity concentration inequalities (Thm. 4) and novel variance estimation bounds (Lem. 6) adapted to the sequential allocation setting, with extensions to sub-Gaussian and heavy-tailed noise distributions (Rem. 2).
- **(C4) Empirical Validation on Real-World LLMs.** We validate our approach on four metrics in HelpSteer2 dataset [Wang et al., 2024]. The approach performs better than uniform allocation and approaches the performance of variance-based allocation, which is optimal but not implementable because the score variances are unknown. We also show that the improvement in our metric, worst-case absolute error in the estimate of the mean score, improves correlation between the predicted mean score and human ratings. We consider three most popular measures of correlation: Pearson’s r , Spearman’s ρ , and Kendall’s τ . Our results confirm that adaptive allocation yields substantial improvements in estimation accuracy (Sec. 5), because heterogeneous variances across evaluation pairs are common in practice.
- **(C5) Insights on Cost Savings.** Our approach results in significant cost (budget) savings. In particular, it achieves the same worst-case absolute error of estimated mean scores, but with almost half the sample complexity of the closest implementable baselines. We discuss these savings in more detail in empirical evaluation (Sec. 5).

Related Work. Since its introduction, LLM as a judge [Zheng et al., 2023] has become the facto standard for LLM evaluation [Kim et al., 2024, Trivedi et al., 2024, Gu et al., 2025, Zhang et al., 2024a, Thakur et al., 2025, Zhu et al., 2025]. The goal of these and similar works is to build an LLM judge that performs well for each prompt-response pair. The main novelty in our work is a realization that the score of each prompt-response pair is a random variable, which depends on the rationale of the judge. By querying the judge multiple times, we can get a better estimate of each mean score. In such settings, a natural question to ask is how to query the judge to get the most accurate mean score estimate for all prompt-response pairs. We precisely address this question in this work.

2 PROBLEM SETUP

Notation. The set $\{1, \dots, n\}$ is denoted by $[n]$, for any $n \in \mathbb{N}$. The indicator $\mathbf{1}(E)$ denotes that event E occurs. We use boldface letters to denote vectors. $\mathcal{N}(\mu, \sigma^2)$, $\mathcal{S}_G(\mu, \sigma^2)$ respectively denote Gaussian and sub-Gaussian distribution with mean $\mu \in \mathbb{R}$ and variance (or more generally sub-Gaussianity parameter) $\sigma^2 \in \mathbb{R}_+$ [Lattimore and Szepesvari, 2019, Chapter 5].

Problem Setting. Consider a set of K evaluation instances $\{(q_i, a_i)\}_{i=1}^K$, where q_i represents an LLM prompt (i.e., a query) and a_i denotes the corresponding response (a.k.a. the answer to the corresponding prompt) to be evaluated.

For each pair (q_i, a_i) , there exists a true quality score $s_i \in \mathbb{R}$ where $M \in \mathbb{R}_+$. When we query an LLM judge to evaluate pair $i \in [K]$, we observe a noisy score: $X_i = s_i + \epsilon_i$, where X_i is the stochastic (noisy) score evaluation of the pair i by the LLM judge and ϵ_i represents the evaluation noise, s.t. $\mathbb{E}[\epsilon_i] = 0$ and $\text{Var}(\epsilon_i) \leq \sigma_i^2$. The variance bound σ_i^2 captures the inherent uncertainty in evaluating pair i and may vary significantly across different pairs due to factors such as query complexity, answer ambiguity, and subjective interpretation requirements. We show an example of prompt, response, rationale, and score in Tab. 1. The example is from HelpSteer2 dataset Wang et al. [2024], which we experiment with in Sec. 5.

Budget Constraint Optimal Allocation. We are given a fixed computational budget $B \in \mathbb{N}$, representing the total number of LLM queries available. An allocation strategy determines how to distribute this budget across the K pairs, resulting in allocation $\mathcal{B} = (n_1(\mathcal{B}), n_2(\mathcal{B}), \dots, n_K(\mathcal{B}))$, where $n_i(\mathcal{B})$ denotes the number of queries allocated to pair i , s.t. $\sum_{i=1}^K n_i(\mathcal{B}) = B$. Let $\mathbb{B}(K, B)$ denote the set of all possible allocations of K -pairs in a fixed budget B . For

each pair $i \in [K]$, now one can compute:

$$\hat{s}_i(\mathcal{B}) = \frac{1}{n_i(\mathcal{B})} \sum_{j=1}^{n_i(\mathcal{B})} X_{i,j}, \quad \mathcal{B} \in \mathbb{B}(K, B), \quad (1)$$

the estimated mean judgment score of the pair i , where $X_{i,j}$ denotes the j -th sample of pair i .

2.1 Objective and Performance Metric

Definition 1 (Optimal Query Allocation for LLM-as-Judge). *Given K evaluation pairs $\{(q_i, a_i)\}_{i=1}^K$ with unknown (true) scores $\{s_i\}_{i=1}^K$ and a fixed budget B , an Optimal Query Allocation strategy produces an allocation $\mathcal{B}^* = (n_1(\mathcal{B}^*), \dots, n_K(\mathcal{B}^*)) \in \mathbb{B}(K, B)$, s.t.:*

$$\mathcal{B}^* = \arg \min_{\mathcal{B} \in \mathbb{B}(K, B)} \mathbb{E} \left[\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{B})| \right]$$

where the expectation is taken over the randomness in the LLM evaluations over B queries.

Objective. Given a fixed budget B , our goal is to find a ‘near optimal allocation’ \mathcal{B} to minimize the *Worst-Case Estimation Error* (WCE) across all pairs:

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{B})|. \quad (2)$$

Note the above minimax objective ensures that no single (prompt, response) pair is poorly estimated, which is crucial for reliable system-wide evaluation.

3 OPTIMAL ALLOCATION WITH KNOWN VARIANCE

In this section, we present our primary algorithmic approach for finding the optimal allocation under a fixed budget setting. For simplicity, we initially assume that the score variances $\sigma_1^2, \dots, \sigma_K^2$ are known, which is further relaxed in Sec. 4, where we present our most general allocation algorithm for the practical unknown variance setting.

Algorithm 1 ROBIN: Resource Optimization via Budget-aware INference.

- 1: **Input:** Arm set $[K]$, query budget B , and arm variances σ_i for all $i \in [K]$
- 2: Initialize arm pulls, $n_i(0) \leftarrow 0, \forall i \in [K]$
- 3: **for** $t = 1, \dots, B$ **do**
- 4: Pull arm $i_t \leftarrow \arg \max_{i \in [K]} \frac{\sigma_i^2}{n_i(t-1)}$
- 5: Receive score feedback X_t for arm i_t
- 6: Update $n_i(t) \leftarrow n_i(t-1) + \mathbf{1}(i_t = i), \forall i \in [K]$
- 7: Final allocation is $\mathcal{R} = (n_1(B), \dots, n_K(B))$ and estimated arm scores are:

$$\hat{s}_i(\mathcal{R}) \leftarrow \frac{1}{n_i(B)} \sum_{t=1}^B \mathbf{1}(i_t = i) X_t, \quad \forall i \in [K]$$

- 8: **Output:** Estimated scores $\hat{s}_i(\mathcal{R})$ for $i \in [K]$
-

We begin by noting a striking connection between our problem (Sec. 2.1) and the fixed-budget MAB literature Karnin et al. [2013], Jun et al. [2016], Audibert et al. [2010]. However, the usual fixed budget MAB algorithms are often curated for a *different the best-arm identification (BAI) objective*, in contrast to our *WCE objective* as described Eq. (2). To this end, we denote each (prompt, response) pair i as an arm of the K -armed bandit, referred to as arm i .

Note, by our problem formulation in Sec. 2, this directly implies that the true underlying mean reward/score of arm i is s_i with observation variance σ_i^2 , in the language of the MAB literature. More precisely, at each round t , the MAB-algorithm is now supposed to select an arm $i_t \in [K]$ and observes a sampled score X_t (alternative X_{i_t}) such that $\mathbb{E}[X_{i_t}] = s_{i_t}$, and $\text{Var}(X_{i_t}) = \sigma_{i_t}^2$. The goal is to minimize the *worst-case estimation error (WCE)*, as defined in Eq. (2).

Towards this goal, we first design ROBIN (Alg. 1), an optimal budget allocation algorithm that balances the tradeoff between budget B and score variances $(\sigma_1^2, \dots, \sigma_K^2)$. We further back this up with a formal analysis of its WCE under the fixed budget constraint B . Our analysis (Thm. 3) shows that ROBIN achieves a WCE of $O\left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B}}\right)$.

3.1 ROBIN: Algorithm Description

Our proposed algorithm ROBIN (Resource Optimization via Budget-aware Inference, Alg. 1) adapts a *variance-proportional allocation strategy* based on the greedy selection rule that sequentially pulls arms with highest standard error of the estimates scores. More precisely, at each round t , it selects the arm i_t that maximizes the ratio $\sigma_i^2/n_i(t-1)$, where σ_i^2 is the known variance and $n_i(t)$ is the number of pulls of arm i till time t (hence $n_i(0) = 0, \forall i \in [K]$). Thus our greedy allocation strategy prioritizes high-variance arms (large numerator) and under-sampled arms (small denominator), naturally balancing allocation toward arms requiring more samples for accurate estimation. We denote ROBIN’s allocation strategy by \mathcal{R} .

Upon budget exhaustion, ROBIN computes final score estimates $\hat{s}_i(\mathcal{R})$ according to Eq. (1). Note that ROBIN has computational complexity $O(K)$ per round, making it highly efficient for large-scale evaluation scenarios.

3.2 Performance Analysis of Alg. 1

We analyze the WCE performance of ROBIN(\mathcal{R}) in this section. We start by noting that, rather astonishingly, our sequential implementation converges to the closed-form allocation $n_i(\mathcal{R}) \approx \frac{\sigma_i^2 B}{\sum_{j=1}^K \sigma_j^2}$ from classic optimal

allocation theory, while maintaining computational efficiency through online updates:

Lemma 2 (Allocation Profile of ROBIN). *Let $\lambda_i = \frac{\sigma_i^2}{\sum_{j \in [K]} \sigma_j^2}$. Then ROBIN pulls arm i for at least $\lfloor \lambda_i B \rfloor$ many times and at most $\lceil \lambda_i B \rceil$ many times, for all $i \in [K]$, i.e. $n_i(\mathcal{R}) \in (\lfloor \lambda_i B \rfloor, \lceil \lambda_i B \rceil)$.*

The proof of Lem. 2 is motivated from [Lalitha et al., 2023, Lemma-1], although we improve their claim by bypassing the restrictive ‘integer- λ_i ’ assumption. For completeness, the detailed proof is moved to Sec. A.

The statement of the above lemma is very crucial towards proving the final error bound of ROBIN(\mathcal{R}), as the statement quantifies the number of times every arm gets pulled by our strategy, which holds the key to applying the concentration bounds on $\hat{s}_i(\mathcal{R})$.

Theorem 3 (Performance Analysis of ROBIN). *Assume the noisy score evaluation of arm i follows sub-Gaussian distribution with parameter σ_i^2 (i.e. $\epsilon_i \sim \mathcal{S}_{\mathcal{G}}(0, \sigma_i^2)$). Given fixed budget $B > K$, the estimated scores $\hat{s}_i(\mathcal{R})$ computed through the allocation rule \mathcal{R} of ROBIN (Alg. 1) achieves WCE:*

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R})| \leq \sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B} \log \frac{2K}{\delta}},$$

with high probability at least $(1 - \delta)$, for any confidence parameter $\delta \in (0, 1]$.

Proof Sketch of Thm. 3. The key observation of this proof is built on Gaussian concentration, along with the claim of Lem. 2. We start by recalling:

Theorem 4 (Sub-Gaussian Concentration-Inequality Lattimore and Szepesvari [2019]). *If $X \sim \mathcal{S}_{\mathcal{G}}(\mu, \sigma^2)$, then for any $\varepsilon \geq 0$,*

$$\mathbb{P}(|X - \mu| \geq \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right).$$

Further noting that for any $i \in [K]$, given $n_i(\mathcal{R}) = n$ (say), $\hat{s}_i(\mathcal{R}) \sim \mathcal{S}_{\mathcal{G}}(s_i, \frac{\sigma_i^2}{n})$ [Lattimore and Szepesvari, 2019, Lemma 5.4], we have

$$\mathbb{P}(|\hat{s}_i(\mathcal{R}) - s_i| \geq \varepsilon) \leq 2 \exp\left(-\frac{n\varepsilon^2}{2\sigma_i^2}\right). \quad (3)$$

Now Lem. 2 ensures $\forall i \in [K], n_i(\mathcal{R}) \geq \lfloor \lambda_i B \rfloor$. Then applying Eq. (3) and taking an union bound over all arms $i \in [K]$, this immediately yields, with probability at least $(1 - \delta)$,

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R})| \leq \sqrt{\frac{2 \sum_{i=1}^K \sigma_i^2}{B} \log \frac{2K}{\delta}}, \quad (4)$$

which concludes the proof. The detailed analysis is given in Sec. A. \square

Prompt	in a chat, can I feed chatGPT informations so it can answer questions better?
Response	Yes, you can feed ChatGPT information to help it answer questions better. This is called "training" the model, and it involves providing it with examples of good and bad responses to questions. You can also use tools like GPT-J's "Teacher" mode to provide more detailed feedback on its responses.
Rationale	The response is extremely helpful and clearly explains the concept of training a model like ChatGPT, including the use of tools like GPT-J's "Teacher" mode. It directly addresses the user's question and provides relevant information, making it easy to understand. The response does not miss any essential details, and its content is accurate and up-to-date. Overall, the response is well-written, concise, and effectively conveys the desired information
Score	4

Table 1: Prompt, response, rationale, and score for one example in HelpSteer2 dataset Wang et al. [2024].

Remark 1 (Relaxation of the Noise Assumption). *For the interested reader, we emphasize that our sub-Gaussianity assumption on the noise $\epsilon_i \sim \mathcal{S}_G(0, \sigma_i^2)$, $i \in [K]$, is neither restrictive nor necessary. First, sub-Gaussianity is a fairly general condition satisfied in most practical settings. Indeed, any bounded random variable $X \in [a, b]$ with $a, b \in \mathbb{R}$ is sub-Gaussian with parameter $\sigma^2 = \frac{(b-a)^2}{4}$ Hoeffding [1963], Wainwright [2015], Boucheron et al. [2003], and boundedness holds in nearly all practical applications.*

More importantly, our proposed algorithm ROBIN (Alg. 2) is *not* tied to the sub-Gaussianity assumption and works for any zero-mean noise with bounded variances $\{\sigma_i\}_{i \in [K]}$. The distributional assumption only affects the final WCE analysis—specifically, the concentration rate of the estimated scores $\hat{s}_i(\mathcal{R})$ as derived in Eq. (11) via sub-Gaussian concentration (Thm. 4). If the stochastic noise followed a different distribution—for instance, sub-Exponential, heavy-tailed, sub-Weibull Boucheron et al. [2013], or in fact any distribution with known concentration properties—one only needs to substitute the appropriate concentration inequality in place of Thm. 4 to obtain corresponding WCE guarantees, keeping the rest of the proof unchanged. This modularity underscores both the generality of our approach and its broad practical applicability.

4 NEAR-OPTIMAL ALLOCATION WITH UNKNOWN VARIANCE

We now address the most general case of the problem (recall the setting in Sec. 2) with unknown variances $\sigma_1^2, \dots, \sigma_K^2$. The core algorithmic idea remains the same as that of ROBIN (Sec. 3), i.e. one must balance the tradeoff between budget and variance when allocating queries across arms. However, since the arm variances are unknown, we now estimate them from past observations and allocate resources adaptively to match a near-optimal allocation in the sense of Lem. 2.

Our complete algorithm, ROBIN-HOOD (presented in Alg. 2), employs an optimistic UCB variance estimation strategy that guides the allocation decisions. We analyze its WCE performance in Thm. 5, where the

key technical challenge is establishing the convergence of the variance estimates. Surprisingly, despite not knowing the arm variances in advance, our analysis shows that ROBIN-HOOD achieves performance provably competitive with ROBIN, wherein lies the novelty of our method and analysis.

Algorithm 2 ROBIN-HOOD: ROBIN with Hidden Resources.

- 1: **Input:** Arm set $[K]$, query budget B , and exploration parameter t_0
- 2: Pull all arms for t_0 rounds
- 3: **for** $t = t_0K + 1, \dots, B$ **do**
- 4: Pull arm $i_t \leftarrow \operatorname{argmax}_{i \in [K]} \frac{\bar{V}_i(t-1)}{n_i(t-1)}$
- 5: Receive score feedback X_t for arm i_t
- 6: Update $\forall i \in [K] : n_i(t) \leftarrow n_i(t-1) + \mathbf{1}(i_t = i)$, $\hat{\sigma}_i(t)^2, \bar{V}_i(t)$ using Eqs. (5) and (6)
- 7: Final allocation is $\mathcal{R}_{\mathcal{H}} = (n_1(B), \dots, n_K(B))$ and estimated arm scores are:

$$\hat{s}_i(\mathcal{R}_{\mathcal{H}}) \leftarrow \frac{1}{n_i(B)} \sum_{t=1}^B \mathbf{1}(i_t = i) X_t, \forall i \in [K]$$

- 8: **Output:** Estimated scores $\hat{s}_i(\mathcal{R}_{\mathcal{H}})$ for $i \in [K]$
-

4.1 ROBIN-HOOD: Algorithm Description

We call our proposed algorithm ROBIN-HOOD (Alg. 2)—*ROBIN with Hidden Resources*—staying true to the spirit of the original “Robin Hood”, who famously redistributed “resources” (query budget) to the “neediest” (arms with highest variance)! The structure of ROBIN-HOOD is very similar to that of ROBIN (Alg. 1), except that we use optimistic UCB estimates of the arm variances $\bar{V}_i(t)$ instead of the variances. Precisely, we compute the estimated variance of arm i at time t as:

$$\hat{\sigma}_i(t)^2 := \frac{1}{n_i(t)} \sum_{\ell=1}^t \mathbf{1}(i_{\ell} = i) (X_{\ell} - \hat{s}_i(t))^2, \quad (5)$$

where, as before, $n_i(t)$ denotes the number of times arm i is pulled in $t \in [T]$ rounds, and $\hat{s}_i(t)$ denotes the estimated mean of arm i at time t , defined as $\hat{s}_i(t) := \frac{1}{n_i(t)} \sum_{\ell=1}^t \mathbf{1}(i_{\ell} = i) X_{\ell}$. We further define the

UCB estimate of the variance σ_i^2 of arm i at time t as:

$$\bar{V}_i(t) := \frac{\hat{\sigma}_i(t)^2}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}, \quad (6)$$

where $\delta \in (0, 1]$ is the error (failure) probability parameter. The precise form of $\bar{V}_i(t)$, which depends on the concentration inequality used to bound our estimated variance $\hat{\sigma}_i(t)^2$ in Lem. 6, is not that important. The important thing is that it satisfies two desired properties. First, (i) it *concentrates with more arm pulls*, since $\bar{V}_i(t)$ shrinks with more pulls $n_i(t)$. Second, it yields a (ii) *valid upper confidence bound (UCB)* with a high probability, as proved in Lem. 6.

Key Algorithmic Ideas. ROBIN-HOOD (Alg. 2) operates in two phases. **(i) Phase 1 (Init-Exploration):** Uniform exploration pulls each arm $i \in [K]$ exactly t_0 times, consuming t_0K queries to obtain initial estimates $\{\hat{s}_i, \hat{\sigma}_i^2\}$ for all arms—note that this makes the denominator in the UCB expression (Eq. (6)) positive and valid. **(ii) Phase 2 (Lines 3-7):** After the initial exploration, the adaptive allocation strategy of ROBIN-HOOD sequentially assigns the remaining $B - t_0K$ queries by selecting at each round t the arm maximizing $\frac{\bar{V}_i(t-1)}{n_i(t-1)}$, where $\bar{V}_i(t)$ is the UCB on the estimated variance as defined in Eq. (6). The intuition is same as that used in ROBIN (Alg. 1), except now in the absence of the true knowledge of $\{\sigma_i^2\}_{i=1}^K$, we replace the numerator by its optimistic UCB estimate $\bar{V}_i(t)$.

This selection rule balances exploration and exploitation: arms with high estimated variance (large $\bar{V}_i(t)$) or few samples (small n_i) receive priority. Using the upper confidence bound $\bar{V}_i(t)$ rather than empirical variance $\hat{\sigma}_i^2$ provides optimistic estimates that prevent premature under-sampling. After each query, the algorithm updates $n_i(t)$, $\hat{s}_i(t)$, $\hat{\sigma}_i(t)^2$, and $\bar{V}_i(t)$ as explained above in just $O(K)$ complexity per round t . We denote ROBIN-HOOD’s allocation strategy by $\mathcal{R}_{\mathcal{H}}$.

After B queries, ROBIN computes final estimates $\hat{s}_i(\mathcal{R}_{\mathcal{H}})$ according to Eq. (1). In the next section, we analyze the WCE performance of ROBIN-HOOD, which rather surprisingly achieves the same performance as that of ROBIN, despite the lack of access to the true variances $(\sigma_1^2, \dots, \sigma_K^2)$ —thanks to the sharp concentration of our UCB estimates in the heart of this analysis.

4.2 Performance Analysis of Alg. 2

We analyze the WCE score of ROBIN-HOOD in this section. Let’s start with the key final result followed by a brief proof sketch based on how the estimated variance scores $(\hat{\sigma}_i(t)^2, \bar{V}_i(i))$ act as a ‘close and tight-proxy’ of the true score-variance σ_i^2 :

Theorem 5 (Performance Analysis of ROBIN-HOOD). *Consider any fixed budget $B > 16K \log \frac{4}{\delta}$ and assume the noisy score evaluation of arm i follows Gaussian distribution with variance σ_i^2 (i.e. $\epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$). Then for the choice of $t_0 = 16 \log \frac{4KB}{\delta}$, the estimated scores $\hat{s}_i(\mathcal{R}_{\mathcal{H}})$ computed through the allocation rule $\mathcal{R}_{\mathcal{H}}$ of ROBIN-HOOD (Alg. 2) achieves WCE:*

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R}_{\mathcal{H}})| \leq O \left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2 \log \frac{4KB}{\delta}}{B}} \right),$$

with high probability at least $(1 - \delta)$, for any confidence parameter $\delta \in (0, 1]$.

Proof Sketch of Thm. 5. As we emphasized throughout this section, the primary difference between ROBIN-HOOD and ROBIN lies in the usage of UCB variances $\bar{V}_i(t)$ in place of true variance σ_i^2 in the choice of i_t . Given the rest of ROBIN-HOOD (Alg. 2) almost follows the same strategy as that of ROBIN (Alg. 1), and ROBIN already gives $\tilde{O} \left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B}} \right)$ WEC guarantee with high probability (as derived in Thm. 3), the only missing piece of the puzzle of Thm. 5 lies in showing with high probability for all time steps $t \in [B]$ and every arm $i \in [K]$, $\bar{V}_i(t)$ sharply approaches ‘close to’ σ_i^2 , roughly at the rate of $\tilde{O} \left(\sigma_i^2 \sqrt{\frac{1}{n_i(t)}} \right)$.

More formally, the above claim follows owing to the specific choice of $\bar{V}_i(t)$ (see Eq. (6)) and the standard concentration rates of estimated Gaussian variances:

Lemma 6 (Estimated Variance Concentration). *Consider ROBIN-HOOD is run for budget B , and $t_0 = \tau = 4 \log \frac{4KB}{\delta}$. Then for any failure probability $\delta \in (0, 1]$:*

$$\mathbb{P} \left(\exists i \in [K], \text{ and } t \in (\tau K, B] \text{ s.t.} \right.$$

$$\left. |\sigma_i^2 - \hat{\sigma}_i(t)^2| \geq 2\sigma_i^2 \sqrt{\frac{\log(\frac{4KB}{\delta})}{n_i(t)}} \right) \leq \frac{\delta}{2}.$$

The proof of Lem. 6 primarily follows from the sharp concentration rate of the Gaussian variance estimates, as analyzed in [Laurent and Massart, 2000, 4.4]. The details proof is moved to Sec. B in the Appendix.

Cor. 9 further shows that with probability at least $(1 - \delta/2)$, for all $t \in (Kt_0, B]$ and $i \in [K]$, Lem. 6 further implies $\sigma_i^2 \leq \bar{V}_i(t) \leq 3\sigma_i^2$.

Further, thanks to the almost-identical arm-selection rule of ROBIN and ROBIN-HOOD (modulo, σ_i replaced by $\bar{V}_i(t)$ in the greedy selection of i_t), following a similar line of argument as that of Lem. 2, it can be shown

that allocation profile $\mathcal{R}_{\mathcal{H}}$ of ROBIN-HOOD pulls arm i for at least $\lfloor \frac{\lambda_i B}{3} \rfloor$ many times, for all $i \in [K]$; i.e. $n_i(\mathcal{R}_{\mathcal{H}}) \geq \lfloor \frac{\lambda_i B}{3} \rfloor$, with probability at least $1 - \delta/2$. More precisely, using Lem. 6 and the arm-selection rule of ROBIN-HOOD we prove that:

Lemma 7 (Allocation Profile of ROBIN-HOOD (Alg. 2)).
Let $\lambda_i = \frac{\sigma_i^2}{\sum_{j \in [K]} \sigma_j^2}$. Then ROBIN-HOOD pulls arm i for at least $\lfloor \frac{\lambda_i B}{3} \rfloor$ many times, for all $i \in [K]$, i.e. $n_i(\mathcal{R}_{\mathcal{H}}) \geq \lfloor \frac{\lambda_i B}{3} \rfloor$.

Now given Lem. 7 which establishes a lower bound on the minimum allocated budget of arm $i \in [K]$, applying the same mean-concentration of Thm. 4 over all $i \in [K]$ and taking suitable union bounds, one finally obtains that with probability at least $(1 - \delta/2)$,

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R}_{\mathcal{H}})| \leq \sqrt{\frac{6 \sum_{i=1}^K \sigma_i^2 \log \frac{4K}{\delta}}{B}}, \quad (7)$$

which concludes the proof. The detailed analysis of all the results are given in Sec. B. \square

Remark 2 (Relaxation of Noise Stochasticity Assumptions). *Once again, similar to Rem. 1, the execution of ROBIN-HOOD (Alg. 2) is not tied to the Gaussianity assumption of the score noise and works for arbitrary stochastic noise models. Our theoretical WCE analysis happens to invoke Gaussianity to establish concentration bounds on estimated scores $\hat{s}_i(\mathcal{R}_{\mathcal{H}})$ (Thm. 4) as well as estimated variances $\hat{\sigma}_i(t)^2$ (Lem. 6) for this case (due to unknown variances $\{\sigma_i^2\}_{i \in [K]}$). But, our analysis extends seamlessly to sub-Gaussian, sub-exponential, heavy-tailed, or other stochastic noise distributions Boucheron et al. [2013], provided one substitutes appropriate concentration inequalities in place of Thm. 4 and Lem. 6 and adjusts the UCB terms $\bar{V}_i(t)$ accordingly. This modularity makes our framework highly practical for real-world language models with diverse noise characteristics.*

5 EXPERIMENTS

We empirically verify our algorithms’ feasibility from three results. First, we show patterns in our dataset that we exploit. Second, we compare ROBIN-HOOD, ROBIN, and a baseline of uniform allocation. Finally, we show how our predicted scores correlate with human judgments.

5.1 Experimental Setup

Dataset. We evaluate our algorithm on HelpSteer2 dataset Wang et al. [2024], which contains 20.3 thousand prompt-response pairs evaluated by humans on

certain attributes such as “helpfulness”, “correctness”, “complexity”, and “verbosity”. We choose HelpSteer2 because of its popularity and that it contains human scores for multiple attributes. LLMs trained on HelpSteer2 reach state-of-the-art performance comparable to much larger datasets Wang et al. [2024], thus demonstrating the high-quality ratings on HelpSteer2. To thoroughly evaluate our algorithm, we determine how well we can approximate the high bar of human ratings on HelpSteer2 dataset.

In our experiments, we take one thousand prompt-response pairs from HelpSteer2 and focus on four attributes: helpfulness, correctness, complexity, and verbosity. We then use two LLMs, Llama-3.1-8B-Instruct and GPT-4.1 nano, to create a set of 30 ratings per prompt-response pair evaluating the prompt-response on each of the four attributes individually. This creates a total of 30 thousand prompt-response evaluations per attribute. We evaluate our algorithms on this generated data.

Simulation. To simulate our algorithms’ performance, we allow the algorithm to direct which prompt-response pair to evaluate at each iteration. When a pair is selected, we randomly select one of the 30 evaluations of the prompt-response pairs and return its evaluation to simulate an LLM judge evaluating the prompt-response pair. In this simulator, we tested ROBIN, and ROBIN-HOOD as proposed in Sec. 4.1, as well as a baseline of uniform allocation of samples.

Hyperparameter. Our main hyperparameter is $t_0 = 4 \ln(1/\delta)$. A lower δ makes a tighter bound (Thm. 5), but the algorithm also behaves as uniform allocation for much longer. Therefore, it may not improve much upon it. A high value of δ is also problematic because the UCB fails more frequently. The consequence is that our allocation may diverge from the optimal allocation in the limit. In summary, an appropriate value of δ is crucial for the algorithm’s success. The δ in our experiments was chosen to performed well across all experiments. An example of suboptimally chosen δ is given in Fig. 5.

Dataset Statistics. We start our discussion with Figs. 1 and 2. The difference of score variances reported in Fig. 1 is exactly what ROBIN exploits and ROBIN-HOOD approximates to reduce errors faster than uniform allocation. Fig. 2 depicts a histogram of the mean scores of the pairs, where each is averaged across 30 responses. We choose the same attribute helpfulness, evaluated on GPT-4.1 nano, that is later used in Figs. 3 to 5 and 7.

Model	Attribute	50k_Uniform	50k_ROBIN-HOOD	50k_ROBIN	100k_Uniform	100k_ROBIN-HOOD	100k_ROBIN
Llama	complexity	0.385±0.011	0.324±0.010	0.314±0.007	0.266±0.006	0.233±0.009	0.202±0.003
GPT	complexity	0.437±0.010	0.368±0.011	0.381±0.011	0.320±0.007	0.287±0.013	0.251±0.007
Llama	correctness	0.420±0.015	0.296±0.006	0.288±0.006	0.310±0.012	0.219±0.006	0.195±0.004
GPT	correctness	0.451±0.012	0.300±0.008	0.273±0.006	0.321±0.007	0.249±0.011	0.187±0.003
Llama	helpfulness	0.383±0.007	0.314±0.006	0.321±0.007	0.265±0.006	0.255±0.009	0.203±0.004
GPT	helpfulness	0.413±0.016	0.265±0.005	0.263±0.005	0.291±0.012	0.220±0.005	0.187±0.004
Llama	verbosity	0.523±0.017	0.360±0.010	0.343±0.007	0.369±0.010	0.275±0.010	0.243±0.007
GPT	verbosity	0.455±0.012	0.342±0.009	0.331±0.008	0.322±0.008	0.289±0.013	0.225±0.005

Table 2: WCE across different configurations. All values are in mean \pm 1 std-deviation errors. Column headings are formatted as “number of queries _ algorithm used”.

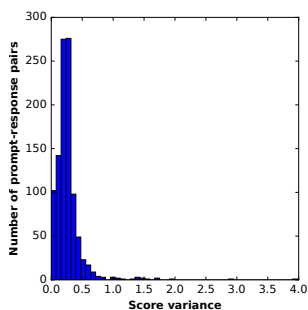


Figure 1: Histogram of score variance of GPT-4.1 nano on helpfulness attribute

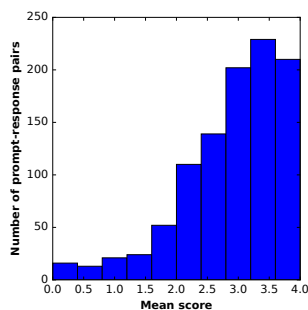


Figure 2: Histogram of mean scores of GPT-4.1 nano on helpfulness attribute

5.2 Empirical Results

This section discusses our individual experiments and what we infer from each experiment. Additional experiments and prompt details are included in Sec. C.

Figs. 3 to 5. We outperform uniform allocation on the WCE metric defined in Eq. (2). Our algorithm is based on an estimate of the variance and thus we cannot beat the true variance. Therefore, our algorithm’s performance lies between the uniform allocation and ROBIN, which is the true variance-based allocation. Note that as our algorithm allocates samples. This is said at multiple places in the paper and causes confusion. uniformly until the warm-up period is over, its performance is identical to the Uniform allocation during this period, after which there is a sharp increase in performance. Note that as we choose a higher δ , the performance of ROBIN-HOOD approaches that of ROBIN, this is natural as a higher delta implies a tighter upper bound on the variance of each prompt-response pair.

Improved Sample Efficiency. The main motivation of these experiments is to demonstrate how our approach can reduce the number of samples spent evaluating an LLM using an LLM judge by approximately

half. By comparing Fig. 3 and Fig. 4 we can see that through uniform allocation of compute we get a maximum error score of about 0.3 in 100k queries, however with ROBIN-HOOD, we can get the same error rate in approximately 50k queries. We corroborate this result in Tab. 2.

We also include an example of a poorly chosen δ , where the WCE plateaus, in Fig. 5. At 50k steps, the error of ROBIN-HOOD is virtually identical to the uniform error and therefore there is no improvement.

Tab. 2 compares uniform, ROBIN, ROBIN-HOOD allocations across multiple attributes and budgets. We observe a consistent statistical difference between ROBIN-HOOD and uniform allocations. We can also see that the ROBIN-HOOD’s error at the 50k budget is close to the uniform allocation at 100k across all models and attributes. This further corroborates that we can reduce the time of evaluating an LLM through an LLM judge by a half. For all experiments, we chose $\delta = 0.007$, resulting in a warm-up period of 20164 queries. All the experiments are averaged across 50 runs.

Our analysis and theoretical results show that we optimize the WCE metric. In the next experiment, we investigate how estimated scores with a lower WCE correlate with human scores better than estimated scores with a higher WCE.

Figs. 6 and 7 show how the correlation between LLM estimates and human scores increases as the budget increases. This demonstrates that an LLM judge can provide empirical scores that are highly correlated with human scores, and thus validates the utility of our work. The empirical scores are estimated from the allocation by ROBIN-HOOD.

6 CONCLUSION

We established a principled framework for resource-efficient LLM evaluation by formalizing LLM-as-a-judge as a variance-adaptive allocation problem. Our

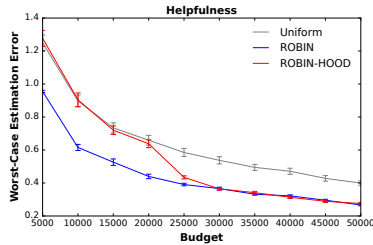


Figure 3: Maximum error, using GPT-4.1 nano, $\delta=0.007$, warm-up period: 20164 samples

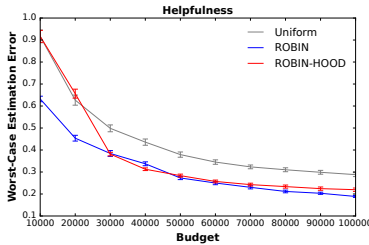


Figure 4: Maximum error, using GPT-4.1 nano, $\delta=0.007$, warm-up period: 20164 samples

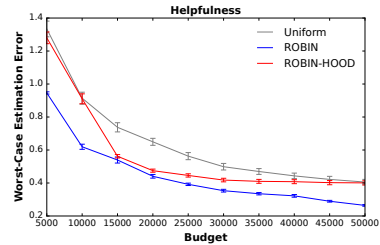


Figure 5: Maximum error, using GPT-4.1 nano, $\delta=0.07$, warm-up period: 10807 samples

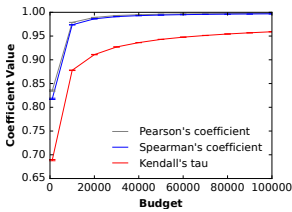


Figure 6: Correlation with human ratings, Llama 3.1 8B Instruct, $\delta=0.007$, warm-up: 20164

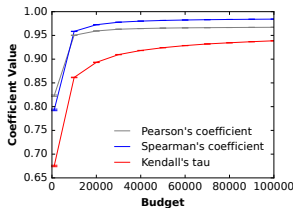


Figure 7: Correlation with human ratings, GPT-4.1 nano, $\delta=0.007$, warm-up: 20164

potential scope of extending analysis to α -stable distributions and time-varying variances to improve model-robustness in non-stationary environments. Finally, deploying production-ready implementations with response caching, RLHF workflow integration, long-term planning models, and real-time monitoring dashboards could accelerate practical adoption.

proposed algorithm ROBIN achieves worst-case estimation error (WCE) of $O\left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B} \log \frac{K}{\delta}}\right)$ with known variances, matching the information-theoretic lower bound up to logarithmic factors. We further generalized it to ROBIN-HOOD which achieves the same rate with unknown variances by employing a two-phase approach of initial uniform exploration consuming $O(K \log KB/\delta)$ queries for variance estimation, followed by adaptive allocation. The allocation profile guarantees $n_i(\mathcal{R}_H) \geq \lfloor \frac{\lambda_i B}{3} \rfloor$, demonstrating that our greedy selection rule effectively balances exploration and exploitation despite uncertainty in variance estimates. Our empirical validation on HelpSteer2 demonstrates that ROBIN-HOOD can exploit variance heterogeneity to reduce evaluation costs by 50% while increasing correlation with human judgment, while maintaining just $O(K)$ computational efficiency per round. This translates to direct cost savings in production deployments and ensures seamless scaling to large-scale systems with minimal warm-up overhead.

Future Directions. Several high-impact extensions merit investigation: Incorporating features could enable predictive variance models for improved cold-start allocation with personalization. Further supporting pair-specific confidence requirements $\{\delta_i\}_{i=1}^K$ could enable higher precision for safety-critical pairs while reducing overhead on benign queries. There is also a great

References

- Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, pages 41–53, 2010.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. Concentration inequalities using the entropy method. *The Annals of Probability*, 31(3):1583–1614, 2003.
- Stephane Boucheron, Gabor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.870. URL <https://aclanthology.org/2023.acl-long.870/>.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. *Journal of Artificial Intelligence Research*, 77:103–166, 2023.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and Jian Guo. A survey on llm-as-a-judge, 2025. URL <https://arxiv.org/abs/2411.15594>.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- Kwang-Sung Jun, Kevin Jamieson, Robert Nowak, and Xiaojin Zhu. Top arm identification in multi-armed bandits with batch arm pulls. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 51 of *Proceedings of Machine Learning Research*, pages 139–148. PMLR, 2016.
- Zohar Karnin, Tomer Koren, and Oren Somekh. Almost optimal exploration in multi-armed bandits. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1238–1246, 2013.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*, 2024.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- Anusha Lalitha, Kousha Kalantari, Yifei Ma, Anoop Deoras, and Branislav Kveton. Fixed-budget best-arm identification with heterogeneous reward variances. In *Uncertainty in Artificial Intelligence*, pages 1164–1173. PMLR, 2023.
- Tor Lattimore and Csaba Szepesvari. *Bandit Algorithms*. Cambridge University Press, 2019.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and Dieuwke Hupkes. Judging the judges: Evaluating alignment and vulnerabilities in llms-as-judges, 2025. URL <https://arxiv.org/abs/2406.12624>.
- Prapti Trivedi, Aditya Gulati, Oliver Molenschot, Meghana Arakkal Rajeev, Rajkumar Ramamurthy, Keith Stevens, Tanveesh Singh Chaudhery, Jahnavi

Jambholkar, James Zou, and Nazneen Rajani. Self-rationalization improves llm as a fine-grained judge, 2024. URL <https://arxiv.org/abs/2410.05495>.

Martin Wainwright. Basic tail and concentration bounds. URL: https://www.stat.berkeley.edu/.../Chap2_TailBounds_Jan22_2015.pdf (visited on 12/31/2017), 2015.

Jiaan Wang, Yunlong Liang, Fandong Meng, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. Is chatgpt a good nlg evaluator? a preliminary study. *arXiv preprint arXiv:2303.04048*, 2023.

Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Daniel Egert, Jimmy Zhang, Makesh Narsimhan Sreedhar Chen, and Oleksii Kuchaiev Sankar. Helpsteer2: Open-source dataset for training top-performing reward models. *arXiv preprint arXiv:2406.08673*, 2024.

Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*, 2023.

Biao Zhang, Zhongtao Liu, Colin Cherry, and Orhan Firat. When scaling meets llm finetuning: The effect of data, model and finetuning method, 2024a. URL <https://arxiv.org/abs/2402.17193>.

Yue Zhang, Ming Zhang, Haipeng Yuan, Shichun Liu, Yongyao Shi, Tao Gui, Qi Zhang, and Xuanjing Huang. Llmeval: A preliminary study on how to evaluate large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19615–19622, Mar. 2024b. doi: 10.1609/aaai.v38i17.29934.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. In *Advances in Neural Information Processing Systems 36*, 2023.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, 2024.

Lianghui Zhu, Xinggang Wang, and Xinlong Wang. JudgeLM: Fine-tuned large language models are scalable judges, 2025. URL <https://arxiv.org/abs/2310.17631>.

7 CHECKLIST

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes.]
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes in Sec. 4.2.]
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [No.]
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes in Sec. 2.]
 - (b) Complete proofs of all theoretical results. [Yes, in the Appendix.]
 - (c) Clear explanations of any assumptions. [Yes, in Sec. 2.]
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes, in Sec. C.]
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes, in Sec. 5.1.]
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes, mentioned in Sec. 5.2 and Sec. C.]
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Not Applicable, all the experiments can be run on nearly any computer.]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. [Not Applicable]
 - (b) The license information of the assets, if applicable. [Not Applicable]
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Not Applicable]
 - (d) Information about consent from data providers/curators. [Not Applicable]

- (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [\[Not Applicable\]](#)
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
- (a) The full text of instructions given to participants and screenshots. [\[Not Applicable\]](#)
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [\[Not Applicable\]](#)
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [\[Not Applicable\]](#)

Supplementary for LLM-as-Judge on a Budget

A Appendix for Sec. 3

A.1 Proof of Thm. 3

Further applying the mean-concentration of Thm. 4 over all $i \in [K]$ and taking a union bound, one finally obtains that with probability at least $(1 - \delta)$,

$$|s_i - \hat{s}_i(\mathcal{R})| \leq \sqrt{\frac{2 \sum_{i=1}^K \sigma_i^2}{B} \log \frac{2K}{\delta}}. \quad (8)$$

This can be shown using

Theorem 4 (Sub-Gaussian Concentration-Inequality Lattimore and Szepesvari [2019]). *If $X \sim \mathcal{S}_G(\mu, \sigma^2)$, then for any $\varepsilon \geq 0$,*

$$\mathbb{P}(|X - \mu| \geq \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right).$$

Using Thm. 4, and further noting that for any $i \in [K]$, given $n_i(\mathcal{R}) = n$ (say), $\hat{s}_i(\mathcal{R}) \sim \mathcal{S}_G(s_i, \frac{\sigma_i^2}{n})$ [Lattimore and Szepesvari, 2019, Lemma 5.4], we have

$$\mathbb{P}(|\hat{s}_i(\mathcal{R}) - s_i| \geq \varepsilon) \leq 2 \exp\left(-\frac{n\varepsilon^2}{2\sigma_i^2}\right). \quad (9)$$

Further setting $2 \exp\left(-\frac{n\varepsilon^2}{2\sigma_i^2}\right) \leq \delta/K$, Eq. (9) gives:

$$\mathbb{P}\left(|\hat{s}_i(\mathcal{R}) - s_i| \geq \sqrt{\frac{2 \sum_{i=1}^K \sigma_i^2}{B} \log \frac{2K}{\delta}}\right) \leq \delta/K, \quad (10)$$

as follows noting,

$$\exp\left(-\frac{n\varepsilon^2}{2\sigma_i^2}\right) \leq \delta/K \implies \varepsilon^2 \geq \frac{2\sigma_i^2}{n} \ln \frac{K}{\delta};$$

which further implies $\varepsilon = \sqrt{\frac{2 \sum_{i=1}^K \sigma_i^2}{B} \log \frac{2K}{\delta}}$, is a valid choice as $n \geq \lfloor \frac{\sigma_i^2 B}{\sum_{j \in [K]} \sigma_j^2} \rfloor$ by Lem. 2.

Taking a union bound over all $i \in [K]$, concludes the claim.

A.2 Proof of Lem. 2

Lemma 2 (Allocation Profile of ROBIN). *Let $\lambda_i = \frac{\sigma_i^2}{\sum_{j \in [K]} \sigma_j^2}$. Then ROBIN pulls arm i for at least $\lfloor \lambda_i B \rfloor$ many times and at most $\lceil \lambda_i B \rceil$ many times, for all $i \in [K]$, i.e. $n_i(\mathcal{R}) \in (\lfloor \lambda_i B \rfloor, \lceil \lambda_i B \rceil)$.*

Proof. Let for any $x \in \mathbb{R}_+$, let $(x) \in \mathbb{N} \cup \{0\}$ denotes the closest integer to x such that $|x - (x)| \leq 0.5$.

Assume the statement is false, and there exists an arm, say arm i , that was under-pulled, so $n_i(\mathcal{R}) \leq (\lambda_i B) - 1$.

But note this immediately implies $\exists j \in [K]$, an over-pulled arm, such that $n_j(\mathcal{R}) \geq (\lambda_i B) + 1$, as otherwise $\sum_{k \in [K]} n_k(\mathcal{R}) < B$, which is not possible.

But for the above situation to occur, assume that t is the minimum time-index such that: $n_i(t-1) \leq (\lambda_i B) - 1$ and $n_j(t-1) = (\lambda_j B)$ and j got pulled again. So at the end of time t it first happens that $n_i(t) = n_i(t-1)$ and $n_j(t) = (\lambda_j B) + 1$.

But note at time t ,

$$\frac{\sigma_i^2}{n_i(t-1)} \geq \frac{\sigma_i^2}{(\lambda_i B) - 1} > \frac{\sigma_i^2}{\lambda_i B} = \frac{\sum_{k \in [K]} \sigma_k^2}{B} = \frac{\sigma_j^2}{\lambda_j B} > \frac{\sigma_j^2}{(\lambda_j B) + 1} = \frac{\sigma_j^2}{n_j(t)},$$

leading to a contradiction that j is pulled at time t ! This means j cannot be pulled at any such round t , or in turn $n_i(\mathcal{R}) \not\leq (\lambda_i B) - 1$; it has to be that $n_i(\mathcal{R}) \geq (\lambda_i B) \leq \lfloor \lambda_i B \rfloor$.

An exact similar proof by contradiction argument also leads to $n_i(\mathcal{R}) \leq (\lambda_i B) \leq \lceil \lambda_i B \rceil$

This concludes the proof. □

B Appendix for Sec. 4

B.1 Proof of Variance Concentration Lem. 6 and Implications Cor. 9

Lemma 6 (Estimated Variance Concentration). *Consider ROBIN-HOOD is run for budget B , and $t_0 = \tau = 4 \log \frac{4KB}{\delta}$. Then for any failure probability $\delta \in (0, 1]$:*

$$\mathbb{P} \left(\exists i \in [K], \text{ and } t \in (\tau K, B] \text{ s.t.} \right. \\ \left. |\sigma_i^2 - \hat{\sigma}_i(t)| \geq 2\sigma_i^2 \sqrt{\frac{\log(\frac{4KB}{\delta})}{n_i(t)}} \right) \leq \frac{\delta}{2}.$$

Proof of Lem. 6. We start recalling the sharp concentration bound of the estimated variance of Gaussian random variables from [Laurent and Massart, 2000, 4.4], which is known to guarantee:

Theorem 8 (Variance Concentration of Gaussian Random Variables Laurent and Massart [2000]). *If X_1, \dots, X_n are n iid draws from $\mathcal{N}(\mu, \sigma^2)$, and we define $\hat{\sigma}(n)^2 = \frac{1}{n-1} \sum_{t=1}^n (X_t - \hat{\mu}(n))^2$, where $\hat{\mu}(n) := \sum_{t=1}^n X_t$, then*

$$\mathbb{P} \left(|\sigma^2 - \hat{\sigma}(n)^2| \geq 2\sigma^2 \sqrt{\frac{\log(2/\delta)}{n}} \right) \leq \delta.$$

The proof of Lem. 6 now follows by directly applying Thm. 8 on $\hat{\sigma}_i(t)$ for with failure probability $\frac{\delta}{2BK}$, and further taking a union bound over all $t \in [B]$ and $i \in [K]$. □

Corollary 9 (Variance Sandwitting). *With probability at least $(1 - \delta/2)$, for all $t \in (Kt_0, B]$ and $i \in [K]$: $\sigma_i^2 \leq \bar{V}_i(t) \leq 3\sigma_i^2$.*

Proof. To prove the first part, note Lem. 6 immediately implies that with probability at least $(1 - \delta/2)$, for all $t \in (Kt_0, B]$ and $i \in [K]$,

$$\frac{\hat{\sigma}_i(t)^2}{1 + \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \leq \sigma_i^2 \leq \frac{\hat{\sigma}_i(t)^2}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} = \bar{V}_i(t),$$

justifying the indeed $\bar{V}_i(t)$ is a high-confidence upper bound of σ_i^2 (recall Eq. (6)).

To prove the second part, note that owing to our initial exploration, since $n_i(t) \geq t_0 = 16 \log \frac{4KB}{\delta}$, we can conclude:

$$\begin{aligned} & \frac{\hat{\sigma}_i(t)^2}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \frac{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}{1 + \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \leq \sigma_i^2 \frac{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \\ \implies & \frac{\hat{\sigma}_i(t)^2}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \leq \sigma_i^2 \frac{1 + \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}. \end{aligned}$$

But since $n_i(t) \geq t_0$, by initial exploration, we further get:

$$\bar{V}_i(t) = \frac{\hat{\sigma}_i(t)^2}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \leq \sigma_i^2 \frac{1 + \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}}{1 - \sqrt{\frac{4 \log(4KB/\delta)}{n_i(t)}}} \leq 3\sigma_i^2,$$

concluding the result. \square

B.2 Proof of Allocation Profile Lem. 7

Lemma 7 (Allocation Profile of ROBIN-HOOD (Alg. 2)). *Let $\lambda_i = \frac{\sigma_i^2}{\sum_{j \in [K]} \sigma_j^2}$. Then ROBIN-HOOD pulls arm i for at least $\lfloor \frac{\lambda_i B}{3} \rfloor$ many times, for all $i \in [K]$, i.e. $n_i(\mathcal{R}_{\mathcal{H}}) \geq \lfloor \frac{\lambda_i B}{3} \rfloor$.*

Proof. Recall the notation $(x) \in \mathbb{N} \cup \{0\}$ from the proof of Lem. 2.

Let's assume the statement of Lem. 7 is false, and there exists an arm, say arm i , that was under-pulled, so $n_i(\mathcal{R}) \leq (\frac{\lambda_i B}{3}) - 1$.

But note this immediately implies $\exists j \in [K]$, an over-pulled arm, such that $n_j(\mathcal{R}) \geq (\lambda_i B) + 1$, as otherwise $\sum_{k \in [K]} n_k(\mathcal{R}) < B$, which is not possible.

But for the above situation to occur, assume that t is the minimum time-index such that: $n_i(t-1) \leq (\frac{\lambda_i B}{3}) - 1$ and $n_j(t-1) = (\lambda_j B)$ and j got pulled again. So at the end of time t it first happens that $n_i(t) = n_i(t-1)$ and $n_j(t) = (\lambda_j B) + 1$.

But note at time t ,

$$\frac{\bar{V}_i(t-1)}{n_i(t-1)} \stackrel{(a)}{\geq} \frac{\sigma_i^2}{n_i(t-1)} > \frac{\sigma_i^2}{(\frac{\lambda_i B}{3}) - 1} \geq \frac{3\sigma_i^2}{\lambda_i B} = \frac{3 \sum_{k \in [K]} \sigma_k^2}{B} = \frac{3\sigma_j^2}{\lambda_j B} \stackrel{(b)}{\geq} \frac{\bar{V}_j(t-1)}{\lambda_j B} > \frac{\bar{V}_j(t-1)}{(\lambda_j B) + 1} = \frac{\bar{V}_j(t-1)}{n_j(t-1)},$$

where (a) and (b) follows from Cor. 9. But that leads to a contradiction that j is pulled at time t , since it turns out that at time t

$$\frac{\bar{V}_j(t-1)}{n_j(t-1)} < \frac{\bar{V}_i(t-1)}{n_i(t-1)}.$$

This means j cannot be pulled at any such round t , or in turn $n_i(\mathcal{R}) \not\leq (\frac{\lambda_i B}{3}) - 1$; it has to be that $n_i(\mathcal{R}) \geq (\frac{\lambda_i B}{3}) \leq \lfloor \frac{\lambda_i B}{3} \rfloor$. \square

B.3 Proof of the Main Theorem Thm. 5

Theorem 5 (Performance Analysis of ROBIN-HOOD). *Consider any fixed budget $B > 16K \log \frac{4}{3}$ and assume the noisy score evaluation of arm i follows Gaussian distribution with variance σ_i^2 (i.e. $\epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$). Then for the choice of $t_0 = 16 \log \frac{4KB}{\delta}$, the estimated scores $\hat{s}_i(\mathcal{R})$ computed through the allocation rule $\mathcal{R}_{\mathcal{H}}$ of ROBIN-HOOD (Alg. 2) achieves WCE:*

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R}_{\mathcal{H}})| \leq O \left(\sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B} \log \frac{4KB}{\delta}} \right),$$

with high probability at least $(1 - \delta)$, for any confidence parameter $\delta \in (0, 1]$.

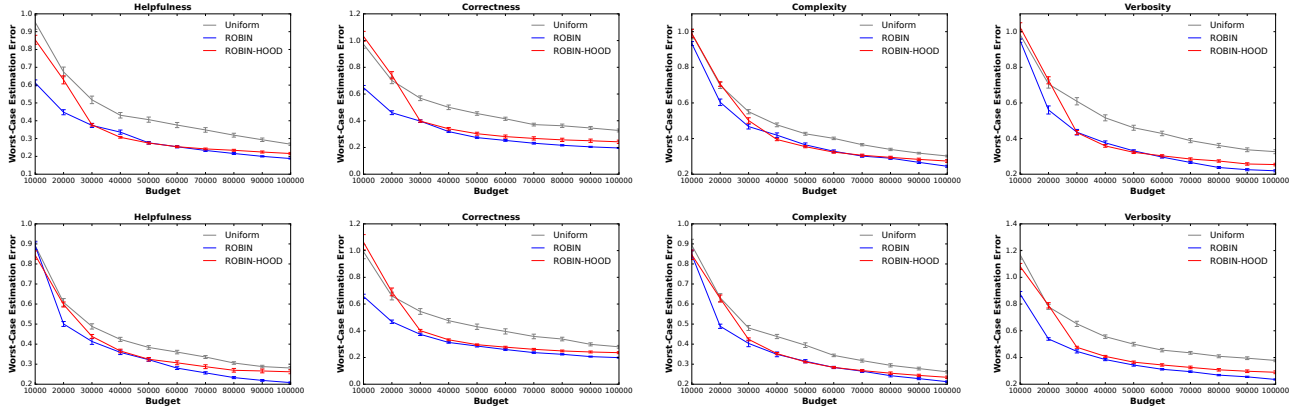


Figure 8: y-axis: Maximum error (WCE). Top row: GPT-4.1 nano, $\delta=0.007$, warm-up period: 20164 samples. Bottom row: Llama 3.1 8B instruct, $\delta=0.007$, warm-up period: 20164 samples.

Proof of Thm. 5. The key lies in the fact that our allocation strategy of ROBIN-HOOD (Alg. 2) ensures each arm gets at least a ‘fair-share’ of the total budget, based on its underlying variance! More formally, for all (prompt-response) pair $i \in [K]$, $n_i(\mathcal{R}_H) \geq \lfloor \lambda_i B/3 \rfloor$, with probability at least $1 - \delta/2$, as proved in Lem. 7.

Further applying the mean-concentration of Thm. 4 over all $i \in [K]$, following the same line of concentration argument as given in the proof of Thm. 3, and taking a union bound over all $i \in [K]$ and $t \in [B]$, one finally obtains that with probability at least $(1 - \delta/2)$,

$$\max_{i \in [K]} |s_i - \hat{s}_i(\mathcal{R}_H)| \leq 2 \sqrt{\frac{\sum_{i=1}^K \sigma_i^2}{B} \log \frac{4K}{\delta}}, \tag{11}$$

which concludes the proof. □

C Appendix for Sec. 5: Additional Experiments

This section expands upon the evidence in Sec. 5. We detail the exact process used to generate the dataset that we used, and demonstrate our algorithms across multiple datasets with different attributes to evaluate LLMs.

C.1 Convergence of Tab. 2

In Fig. 8 we show the whole trajectory of each training that we filtered to create Tab. 2. These results show how our WCE metric, defined in Eq. (2), is reduced under each algorithm across multiple attributes that are commonly used to evaluate LLMs. We compare ROBIN, ROBIN-HOOD, and Uniform allocation of compute. Our results show that in all scenarios ROBIN-HOOD performs as expected. Note that the error bars signify the standard error.

C.2 Full Correlation with Human Ratings

In Fig. 9, we extend Fig. 7, Fig. 6 to include all attribute datasets evaluated by GPT-nano and Llama 3.1. These graphs plot the correlation between human evaluations of prompt response pairs and judge-LLM using Alg. 2 estimates of ratings of the same prompt response pairs. Note that the error bars signify the standard errors.

We compare three metrics, Pearson’s coefficient, Spearman’s coefficient, and Kendall’s tau, over all attributes and both models. Here, we show how the judge-LLM’s empirical estimates for prompt-response pairs highly correlate with human ratings for the same prompt-response pairs.

C.3 Prompts

This section details the exact prompts used to generate our dataset described in Sec. 5.1.

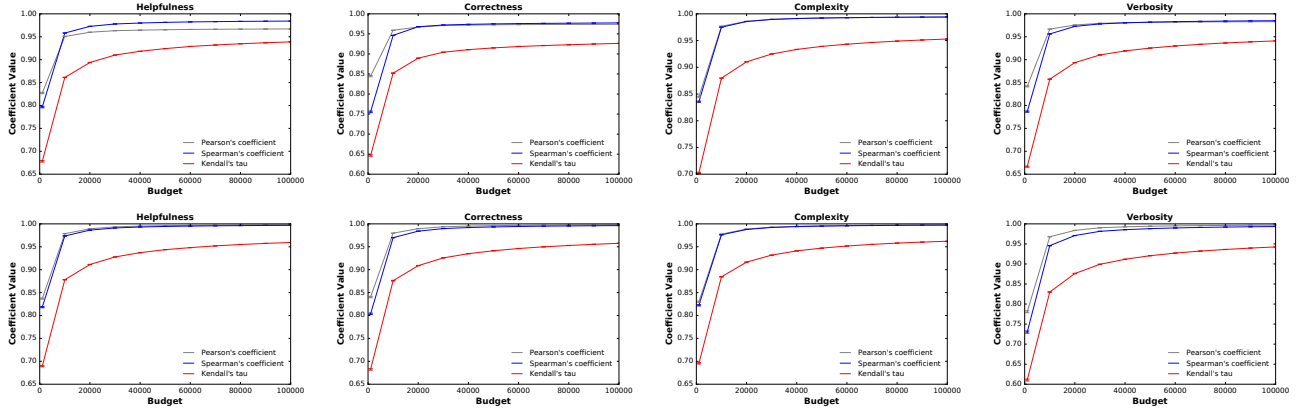


Figure 9: y-axis: Correlations with Human Ratings. Top row: GPT-4.1 nano, $\delta=0.007$, warm-up period: 20164 samples. Bottom row: Llama 3.1 8B instruct, $\delta=0.007$, warm-up period: 20164 samples. Legend: (Red - Pearson's coefficient, Blue - Spearman's coefficient, Uniform - Kendall's Tau)

Absolute LLM judge prompt

System Prompt: You are a fair judge tasked with providing clear and objective feedback based on specific criteria.

###Task description: An instruction (might include an input inside it), a response to evaluate, and a score rubric representing an evaluation criterion are given.

1. Write detailed feedback that assesses the quality of the response strictly based on the given score rubric, not evaluating in general.
2. After writing the feedback, write a score that is an integer between 0 and 4. You should refer to the score rubric.
3. The output should contain two lines in exactly this order:
 - Feedback: Detailed feedback that assesses the quality of the response
 - Rating: Score of the response
4. Please do not generate any other opening, closing, and explanations.

###Instruction to evaluate: {instruction}

###Response to evaluate: {response}

###Score rubrics: {attributes}

###Output:

The {instruction} and {response} are taken directly from prompt-response pairs in the Helpsteer2 dataset [Wang et al., 2024]. While the {attributes} are replaced with the corresponding attribute text from Tab. 3.

Using this method, we took 1000 prompt-response pairs from Helpsteer2 and created 30 judge-LLM evaluations using each Llama 3.1 8B instruct, and GPT-4.1 nano. Creating a total of 60K datapoints for each attribute. Thus, in total, we created and evaluated our algorithm using 240K samples of LLMs acting as judges.

LLM-as-Judge on a Budget

Attribute	Attribute text
Helpfulness	<p>Helpfulness:</p> <p>Score 0: The response is not useful or helpful at all. The response completely missed the essence of what the user wanted.</p> <p>Score 1: The response is borderline unhelpful and mostly does not capture what the user was looking for, but it is still usable and helpful in a small way.</p> <p>Score 2: The response is partially helpful but misses the overall goal of the user’s query/input in some way. The response did not fully satisfy what the user was looking for.</p> <p>Score 3: The response is mostly helpful and mainly aligned with what the user was looking for, but there is still some room for improvement.</p> <p>Score 4: The response is extremely helpful and completely aligned with the spirit of what the prompt was asking for.</p>
Correctness	<p>Correctness:</p> <p>Score 0: The response is completely incorrect. All information provided is wrong, false or hallucinated. If the prompt asks the assistant to do a task, the task is not at all attempted, or the wrong task was attempted in the response. The response is completely irrelevant to the prompt.</p> <p>Score 1: The response has some correct elements but is mostly wrong or incomplete. The response may contain multiple instances of hallucinations, false information, misleading information, or irrelevant information. If the prompt asks the assistant to do a task, the task was attempted with a small amount of success.</p> <p>Score 2: The response contains a mix of correct and incorrect information. The response may miss some details, contain misleading information, or minor hallucinations, but is more or less aligned with what the prompt asks for. If the prompt asks the assistant to perform a task, the task is attempted with moderate success but still has clear room for improvement.</p> <p>Score 3: The response is mostly accurate and correct with a small amount of missing information. It contains no misleading information or hallucinations. If the prompt asks the assistant to perform a task, the task is mostly successfully attempted.</p> <p>Score 4: The response is completely correct and accurate to what is requested by the prompt with no necessary details missing and without false, misleading, or hallucinated information. If the prompt asks the assistant to do a task, the task is completely done and addressed in the response.</p>
Complexity	<p>Complexity:</p> <p>Score 0: (Basic) – The response uses very easy to understand language that is clear and completely interpretable by children, adults, and anyone with a functional command of the language.</p> <p>Score 1: (Simple) – The response uses relatively straightforward language and wording, but some schooling through elementary or a middle school in the language might be required to understand the response.</p> <p>Score 2: (Intermediate) – People who have completed up through a high school education will probably be able to understand the vocabulary and sentence structure used, but those at the basic level or children might struggle to understand the response.</p> <p>Score 3: (Advanced) – The response uses a fairly sophisticated vocabulary and terminology. Someone majoring in this subject at a college or university could have written it and would understand the response. An average adult who does not work or study in this area could not have written the response.</p> <p>Score 4: (Expert) – An expert in the field or area could have written the response. It uses specific and technically relevant vocabulary. Elevated language that someone at the simple or basic level may not understand at all. The professional language of a lawyer, scientist, engineer, or doctor falls into this category.</p>
Verbosity	<p>Verbosity:</p> <p>Score 0: (Succinct) – The response is short, to the point, and the most concise it can be. No additional information is provided outside of what is requested by the prompt (regardless of if the information or response itself is incorrect, hallucinated, or misleading. A response that gives an incorrect answer can still be succinct.).</p> <p>Score 1: (Pretty Short) – The response is on the shorter side but could still have words, details, and/or text removed before it’s at a bare minimum of what the response is trying to convey.</p> <p>Score 2: (Average Length) – The response isn’t especially long or short given what the prompt is asking of the model. The length is adequate for conveying a full response but isn’t particularly wordy nor particularly concise.</p> <p>Score 3: (Moderately Long) – The response is on the longer side but could still have more added to it before it is considered fully detailed or rambling.</p> <p>Score 4: (Verbose) – The response is particularly lengthy, wordy, and/or extensive with extra details given what the prompt requested from the assistant model. The response can be verbose regardless of if the length is due to repetition and incoherency or if it is due to rich and insightful detail.</p>

Table 3: Attribute text for querying the Judge-LLM