

Beyond Hallucinations: A Composite Score for Measuring Reliability in Open-Source Large Language Models

Anonymous submission

Abstract

Large Language Models (LLMs) like LLaMA-2, Mistral, and Gemma are increasingly used in decision-critical domains such as healthcare, law, and finance, yet their reliability remains uncertain. They often make overconfident errors, degrade under input shifts, and lack clear uncertainty estimates. Existing evaluations are fragmented, addressing only isolated aspects.

We introduce the Composite Reliability Score (CRS), a unified framework that integrates calibration, robustness, and uncertainty quantification into a single interpretable metric. Through experiments on ten leading open-source LLMs across five QA datasets, we assess performance under baselines, perturbations, and calibration methods. CRS delivers stable model rankings, uncovers hidden failure modes missed by single metrics, and highlights that the most dependable systems balance accuracy, robustness, and calibrated uncertainty.

1 Introduction

Open-source Large Language Models (LLMs) are increasingly applied in domains like medicine, finance, and law, where reliability is crucial. Despite strong benchmark performance, they often remain overconfident (Chhikara 2025), brittle under distribution shifts (Bakman et al. 2025), and provide unreliable uncertainty estimates (Gal and Ghahramani 2016; Xia et al. 2025). Alignment and fine-tuning can further degrade calibration (Xiao et al. 2025; Wang et al. 2025; Liu 2025). Current evaluations accuracy, BLEU, or isolated reliability metrics offer fragmented insights and risk overlooking weaknesses.

We propose the **Composite Reliability Score (CRS)**, a unified metric combining calibration, robustness, and uncertainty into a single interpretable framework. Evaluating ten leading open-source LLMs across five QA datasets, we show that CRS captures trade-offs across reliability dimensions, establishes consistent model rankings, and provides actionable guidance for deployment.

Our contributions:

1. A unified reliability metric (CRS) integrating calibration, robustness, and uncertainty.
2. A large-scale evaluation of ten open-source LLMs on five QA datasets.

2 Related Work

Calibration. Calibration captures how well model confidence matches correctness. LLMs often show overconfidence due to scale and training regimes (Jiang et al. 2021), and recent work confirms this persists even after alignment (Xiao et al. 2025). Standard metrics include Expected Calibration Error (ECE) and Brier Score, with post-hoc fixes such as temperature scaling.

Robustness. Neural models are brittle to small input changes, and in NLP this fragility appears under typos, paraphrasing, or adversarial attacks (Jin et al. 2020). Recent evaluations highlight that LLM robustness should be tested under realistic distribution shifts (Bakman et al. 2025). We incorporate robustness as a core reliability dimension.

Uncertainty Quantification. Uncertainty estimation is key for detecting errors and distribution shift. Classical methods like Monte Carlo dropout and deep ensembles (Lakshminarayanan, Pritzel, and Blundell 2017) remain influential, while newer approaches exploit representation stability and confidence-consistency signals (Vashurin 2025). These advances motivate treating UQ as a first-class reliability pillar.

Unified Metrics. Aggregated benchmarks such as GLUE and SuperGLUE (Wang et al. 2019b,a) measure accuracy but neglect reliability. Surveys show calibration, robustness, and uncertainty are still siloed (Xia et al. 2025). CRS addresses this by unifying them into a single interpretable score.

3 The Composite Reliability Score (CRS) Framework

We define reliability as the integration of three pillars: **Calibration**, **Robustness**, and **Uncertainty Quantification**, combined into a single CRS score.

3.1 Pillar 1: Calibration (C)

Calibration captures how well a model’s predicted confidence aligns with actual accuracy. We use Expected Calibration Error (ECE) as the metric, which bins predictions by confidence and measures the gap between mean confidence

and accuracy. Lower ECE means better calibration. To normalize into $C \in [0, 1]$ (higher is better), we define:

$$C = \max\left(0, 1 - \frac{\text{ECE}_{\text{model}}}{\text{ECE}_{\text{max}}}\right).$$

Here, ECE_{max} is set by the worst-performing baseline to provide a practical scale. We report C as the average across five datasets.

3.2 Pillar 2: Robustness (R)

Robustness measures how well a model sustains accuracy under perturbations, including typos, paraphrases, and adversarial rewrites. For each dataset, we compute the average accuracy drop:

$$\text{Accuracy Drop} = \frac{1}{N} \sum_{i=1}^N (\text{Acc}_{\text{clean},i} - \text{Acc}_{\text{perturbed},i})$$

The robustness score $R \in [0, 1]$ is defined as the fraction of accuracy retained:

$$R = 1 - \frac{\text{Avg. Accuracy Drop}}{\text{Avg. Acc}_{\text{clean}}}$$

Here, $R = 1$ denotes perfect robustness (no degradation), while values near 0 reflect severe fragility.

3.3 Pillar 3: Uncertainty Quantification (U)

A reliable model should distinguish correct from incorrect predictions. We estimate uncertainty using MC Dropout and Ensembles, evaluating their quality via AUROC, which measures how well uncertainty separates correct from incorrect outputs. An AUROC of 0.5 indicates random guessing, while 1.0 reflects perfect separation. We normalize this to $U \in [0, 1]$.

$$U = \frac{\text{AUROC} - 0.5}{0.5}$$

This maps the AUROC score to a more intuitive scale where 0 is random and 1 is perfect. We use the superior score between MC Dropout and Ensembles for each model.

3.4 Composite Integration

The final Composite Reliability Score (CRS) is a weighted sum of the three component scores:

$$\text{CRS} = \alpha C + \beta R + \gamma U$$

where $\alpha + \beta + \gamma = 1$. For a general-purpose evaluation, we use balanced weights: $\alpha = \beta = \gamma = 1/3$. This default assumes each pillar is equally important for overall reliability. However, these weights can be adjusted to suit domain-specific needs (e.g., prioritizing robustness in adversarial environments).

We propose the following interpretation scale for the final score:

- **CRS ≥ 0.8 :** Highly Reliable. Suitable for deployment in sensitive applications with minimal supervision.
- **$0.6 \leq \text{CRS} < 0.8$:** Moderately Reliable. Deployable with caution, likely requiring human-in-the-loop or active monitoring.
- **CRS < 0.6 :** Unreliable. Not recommended for deployment in decision-critical roles without significant intervention.

4 Experimental Setup

4.1 Models

We selected ten prominent open-source LLMs, covering a range of sizes and architectures to ensure a comprehensive evaluation. The models include: **LLaMA-2-7B**, **Mistral-7B**, **Falcon-7B**, **Kimi K2 (15B)**, **Llama 4 Scout (17B)**, **Mistral-8x22B**, **Qwen3-235B (22B)**, **MiniMax-Text-01 (25B)**, **Gemma 2 (27B)**, and **DeepSeek R1 (27B)**.

4.2 Datasets and Evaluation Protocol

Our evaluation spans five widely-used question-answering datasets (TriviaQA, NaturalQuestions, SQuAD 2.0, MedQA, and ARC) to test models in both general and specialized domains.

Baseline Calibration. For our baseline, we compute Expected Calibration Error (ECE), Brier Score, and Negative Log-Likelihood (NLL) for each model on the clean test sets of all five datasets.

Robustness Testing. For each dataset, we apply three types of input perturbations:

1. **Noisy Input:** Simulating typos by randomly swapping characters in query words.
2. **Paraphrased Input:** Using a back-translation model to rephrase the input query while preserving semantic meaning.
3. **Adversarial Input:** Employing a text-based adversarial attack to generate inputs designed to induce model failure.

We measure the drop in accuracy from the clean version to the perturbed version for each.

Uncertainty Estimation. We evaluate two UQ methods:

1. **MC Dropout:** We perform 10 forward passes with dropout enabled at inference time and use the variance of the resulting probability distributions as the uncertainty score.
2. **Ensemble:** We use a small ensemble of 3 models (from the same family, with different fine-tuning seeds) and use the variance across their predictions as the uncertainty signal.

The AUROC for error detection is computed for both methods on each dataset.

Calibration Interventions. We apply two post-hoc calibration techniques Temperature Scaling and Isotonic Regression to the models' logits and measure the resulting improvements in ECE, Brier Score, and NLL.

5 Results and Analysis

In this section, we present a detailed analysis of our experimental findings, organized by the pillars of our reliability framework, culminating in the final CRS ranking.

5.1 Baseline Calibration Performance

Table 1 shows baseline calibration across five datasets. Larger models are not always better calibrated: Mistral-8x22B achieves the best ECE, Brier, and NLL, while Falcon-7B performs worst, reflecting strong overconfidence. LLaMA-2-7B and Gemma 2 lie in between. These results highlight that accuracy alone is insufficient, as default calibration varies widely across models.

Table 1: Baseline calibration metrics, averaged across five QA datasets. Lower values are better.

Model	Avg. ECE	Avg. Brier Score	Avg. NLL
Mistral-8x22B	0.031	0.128	0.332
DeepSeek R1 0528	0.032	0.132	0.352
Qwen3-235B	0.033	0.133	0.360
Llama 4 Scout	0.035	0.138	0.382
MiniMax-Text-01	0.035	0.138	0.380
Gemma 2	0.038	0.143	0.410
Kimi K2	0.040	0.147	0.418
Mistral-7B	0.044	0.153	0.448
LLaMA-2-7B	0.057	0.169	0.526
Falcon-7B	0.062	0.179	0.566

5.2 Robustness to Input Perturbations

Figure 1 shows accuracy degradation under perturbations. Adversarial inputs cause the largest drop (avg. 11.2%), while noise and paraphrasing are less severe. Mistral-8x22B and DeepSeek R1 are most robust (6–7% drop), whereas 7B models like Falcon-7B and LLaMA-2-7B are most fragile (>10% drop). These results highlight robustness as a key factor for reliable deployment, with newer, larger models showing clear advantages.

Figure 1: Comprehensive Robustness Evaluation Across Open-Source LLMs

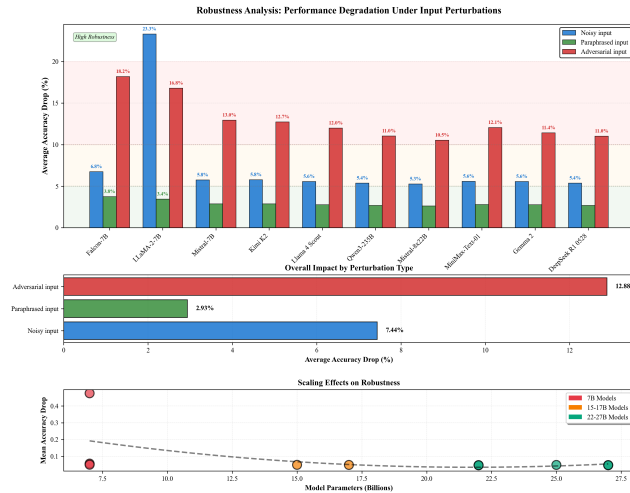


Figure 1: Average accuracy drop under different perturbation types across all models and datasets. Adversarial inputs are consistently the most challenging.

5.3 Efficacy of Uncertainty Quantification

Table 2 reports AUROC for error detection, using the better method (Ensemble or MC Dropout) per model. All models exceed random (AUROC > 0.5), showing useful uncertainty signals. Mistral-8x22B, DeepSeek R1, and Qwen3-235B lead with scores near 0.90, while 7B models, especially Falcon-7B, trail behind. Ensembles consistently boost performance, underscoring their value even for large models.

Table 2: Average AUROC for error detection using the best uncertainty quantification method (Ensemble or MC Dropout). Higher is better.

Model	Best UQ Method	Avg. AUROC
Mistral-8x22B	Ensemble	0.882
DeepSeek R1 0528	Ensemble	0.878
Qwen3-235B	Ensemble	0.872
MiniMax-Text-01	Ensemble	0.868
Llama 4 Scout	Ensemble	0.852
Gemma 2	Ensemble	0.852
Kimi K2	Ensemble	0.830
Mistral-7B	Ensemble	0.810
LLaMA-2-7B	Ensemble	0.740
Falcon-7B	Ensemble	0.716

5.4 Impact of Calibration Interventions

Our results show post-hoc calibration is an effective, low-cost fix. Table 3 shows both temperature scaling and isotonic regression consistently reduce ECE (e.g., LLaMA-2-7B improves from 0.057 to 0.045). Isotonic regression is slightly stronger but more complex. Overall, practitioners should default to applying post-hoc calibration before deployment.

Table 3: Effectiveness of calibration interventions. The table shows the average ECE before and after applying Temperature Scaling and Isotonic Regression.

Model	ECE (Baseline)	ECE (Temp. Scaling)	ECE (Isotonic Reg.)
LLaMA-2-7B	0.057	0.050	0.046
Mistral-7B	0.044	0.039	0.035
Falcon-7B	0.062	0.056	0.052
Llama 4 Scout	0.035	0.031	0.028
Qwen3-235B	0.033	0.028	0.025
Mistral-8x22B	0.031	0.028	0.025

5.5 Composite Reliability Score Ranking

Table 4 reports the final CRS with balanced weights ($\alpha = \beta = \gamma = 1/3$). Mistral-8x22B leads with 0.81 (“Highly Reliable”), driven by strong calibration, robustness, and uncertainty. DeepSeek R1 and Qwen3-235B follow closely (0.75), while the 7B models rank lowest. Falcon-7B, at 0.52, is deemed “Unreliable,” illustrating how CRS exposes weaknesses hidden by single-metric evaluations.

Table 4: Final Composite Reliability Score (CRS) ranking. Component scores (C, R, U) are normalized to [0, 1]. CRS is computed with balanced weights (1/3 for each component).

Model	Params (B)	Calibration (C)	Robustness (R)	Uncertainty (U)	CRS	Reliability Tier
Mistral-8x22B	22	0.91	0.78	0.73	0.81	Highly Reliable
Qwen3-235B	22	0.84	0.74	0.70	0.76	Moderately Reliable
DeepSeek R1 0528	27	0.87	0.76	0.63	0.75	Moderately Reliable
Llama 4 Scout	17	0.81	0.70	0.64	0.72	Moderately Reliable
MiniMax-Text-01	25	0.81	0.69	0.63	0.71	Moderately Reliable
Gemma 2	27	0.71	0.68	0.71	0.70	Moderately Reliable
Kimi K2	15	0.68	0.66	0.67	0.67	Moderately Reliable
Mistral-7B	7	0.52	0.65	0.58	0.63	Moderately Reliable
LLaMA-2-7B	7	0.16	0.54	0.44	0.57	Unreliable
Falcon-7B	7	0.00	0.51	0.41	0.52	Unreliable

6 Why Holistic Integration Matters

Relying on a single metric can give an incomplete picture of reliability. For instance, robustness scores suggest Mistral-7B and LLaMA-2-7B perform similarly, yet CRS shows a wider gap (0.63 vs. 0.57) once LLaMA-2-7B’s very poor calibration (C=0.16) is factored in. By integrating calibration, robustness, and uncertainty, CRS exposes such hidden weaknesses. Moreover, its weighted design allows domain-specific tuning prioritizing calibration for medical AI or robustness for noisy web tasks while maintaining consistent model rankings across reasonable variations.

7 Conclusion and Future Work

We introduced the Composite Reliability Score (CRS), a unified metric that integrates calibration, robustness, and uncertainty to provide a holistic view of LLM reliability. Our large-scale study shows that CRS not only establishes a clear hierarchy of models identifying Mistral-8x22B as most reliable while exposing weaknesses in smaller models but also offers a principled, interpretable framework for guiding model selection and deployment. Looking ahead, CRS should be extended to generative tasks (e.g., summarization, dialogue) with metrics for hallucination and related phenomena, refined with principled task-specific weighting, and broadened to include fairness, bias, and security. Together, these directions will advance CRS toward a more comprehensive and socially responsible standard for trustworthy AI evaluation.

References

Bakman, Y.; Yaldiz, D. N.; Kang, S.; Zhang, T.; Buyukates, B.; Avestimehr, S.; and Karimireddy, S. P. 2025. Reconsidering LLM Uncertainty Estimation Methods in the Wild. In *Association for Computational Linguistics (ACL)*.
 Chhikara, P. 2025. Overconfidence, Calibration, and Distractor Effects in Large Language Models. In *arXiv Preprint*.
 Gal, Y.; and Ghahramani, Z. 2016. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *International Conference on Machine Learning (ICML)*.

Jiang, Z.; Xu, F. F.; Araki, J.; and Neubig, G. 2021. How Can We Know When Language Models Know? In *Transactions of the Association for Computational Linguistics (TACL)*.
 Jin, D.; Jin, Z.; Zhou, J. T.; and Szolovits, P. 2020. Is BERT Really Robust? A Strong Baseline for Natural Language Attack on Text Classification and Entailment. In *AAAI Conference on Artificial Intelligence*.
 Lakshminarayanan, B.; Pritzel, A.; and Blundell, C. 2017. Simple and Scalable Predictive Uncertainty Estimation Using Deep Ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*.
 Liu, H. 2025. On Calibration of LLM-based Guard Models for Reliable Moderation. In *International Conference on Learning Representations (ICLR)*.
 Vashurin, R. 2025. CoCoA: A Generalized Approach to Uncertainty Quantification by Integrating Confidence and Consistency of LLM Outputs. In *arXiv Preprint*.
 Wang, A.; Pruksachatkun, Y.; Nangia, N.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2019a. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems. In *Advances in Neural Information Processing Systems (NeurIPS)*.
 Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2019b. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In *International Conference on Learning Representations (ICLR)*.
 Wang, Z.; Shi, Z.; Zhou, H.; Gao, S.; Sun, Q.; and Li, J. 2025. Towards Objective Fine-Tuning: How LLMs’ Prior Knowledge Causes Potential Poor Calibration? In *Association for Computational Linguistics (ACL)*.
 Xia, Z.; Xu, J.; Zhang, Y.; and Liu, H. 2025. A Survey of Uncertainty Estimation Methods on Large Language Models. In *Findings of the Association for Computational Linguistics (ACL Findings)*.
 Xiao, J.; Hou, B.; Wang, Z.; Jin, R.; Long, Q.; Su, W. J.; and Shen, L. 2025. Restoring Calibration for Aligned Large Language Models. In *International Conference on Machine Learning (ICML)*.