CALIBRATING VERBALIZED CONFIDENCE WITH SELF-GENERATED DISTRACTORS

Anonymous authors
Paper under double-blind review

000

001

002 003 004

010 011

012

013

014

016

017

018

021

023

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Calibrated confidence estimates are necessary for large language model (LLM) outputs to be trusted by human users. While LLMs can express their confidence in human-interpretable ways, verbalized LLM-generated confidence scores have empirically been found to be miscalibrated, reporting high confidence on instances with low accuracy and thereby harming trust and safety. We hypothesize that this overconfidence often stems from a given LLM's heightened suggestibility when faced with claims that it encodes little information about; we empirically validate this hypothesis, finding more suggestibility on lower-accuracy claims. Building on this finding, we introduce Distractor-Normalized Coherence (DINCo), which estimates and accounts for an LLM's suggestibility bias by having the model verbalize its confidence independently across several self-generated distractors (i.e. alternative claims), and normalizes by the total verbalized confidence. To further improve calibration, we leverage generator-validator disagreement, augmenting normalized validator confidence with a consistency-based estimate of generator confidence. Here, we frame the popular approach of self-consistency as leveraging coherence across sampled generations, and normalized verbalized confidence as leveraging coherence across validations on incompatible claims, allowing us to integrate these complementary dimensions of coherence into DINCo. Moreover, our analysis shows that DINCo provides less saturated, and therefore more usable, confidence estimates, and that further sampling alone cannot close the gap between DINCo and baselines, with DINCo at 10 inference calls outperforming self-consistency at 100. We include our code in the supplementary.

1 Introduction

LLMs encode a vast amount of knowledge in their parameters, demonstrating superhuman performance on knowledge-intensive benchmarks (Comanici et al., 2025a; OpenAI, 2023). Users often rely on information obtained from these models to make important decisions, but the information is not always accurate. Thus, we seek to qualify LLM responses with confidence estimates that are *calibrated*, i.e. match the probability of correctness. Users and agentic frameworks often use LLMs in a zero-shot manner without task-specific tuning (Manakul et al., 2023; Geng et al., 2024; Feng et al., 2024; Shorinwa et al., 2025), motivating the development of confidence estimation methods that work in off-the-shelf settings – both gray-box settings with logit access, and black-box settings with only textual input and output.

In these settings, verbalized confidence is a simple and commonly-used approach that prompts the model to report its confidence in an answer (Lin et al., 2022; Xiong et al., 2024; Wei et al., 2024). For brevity, we use *verbalized confidence* as a blanket term for (1) asking the model to decode a numerical confidence like "80%" (Tian et al., 2023b) and (2) asking the model whether an answer is correct and taking P(True) (Kadavath et al., 2022). Verbalized confidence is appealing for several reasons, including resembling the way humans express confidence, making it easy to interpret and integrate into decision-theoretic frameworks (Sun et al., 2025; Steyvers et al., 2025). However, verbalized confidence has several drawbacks. First, it empirically tends to exhibit overconfidence (Tian et al., 2023b; Xiong et al., 2024; Wei et al., 2024; Xu et al., 2025); Fig. 1 (left) shows that verbalized confidence scores generally outstrip average accuracy within a confidence bin.

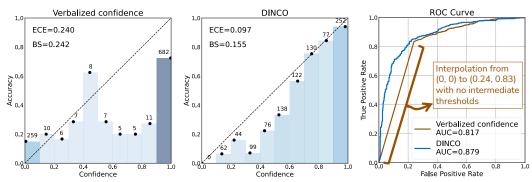


Figure 1: Calibration metrics (Expected Calibration Error \downarrow , Brier score \downarrow , area under the ROC curve \uparrow ; see Appendix B.1) with Qwen3-8B on TriviaQA using P(True) as verbalized confidence. (**Left**) Verbalized confidence is saturated at high confidence. For each bar, we label the number of instances whose confidence falls in the interval and we darken larger bins. (**Center**) DINCO normalizes by the total confidence over candidate answers, relieving saturation and improving calibration. (**Right**) Since verbalized confidence is saturated at high confidence, it is unable to achieve a positive true positive rate (TPR) without incurring a significant false positive rate (FPR) of 0.24. In other words, no rejection threshold can be chosen to reject a high proportion of false claims. Meanwhile, DINCO enjoys better granularity, ranking positives above negatives even among instances with a verbalized confidence of 1.

We highlight a second underexplored factor that makes verbalized confidence suboptimal: **confidence saturation**, wherein the model's reported scores tend to fall into a few bins, making them uninformative. While this might still lead to an acceptable calibration error, it results in "jumpy" curves, as in Fig. 1 (right), where no confidence threshold that accepts at least one claim can avoid accepting a substantial proportion of false claims. To address these shortcomings, we introduce DINCO, which leverages incoherence in verbalized confidence across related claims to detect overconfidence. DINCO is motivated by the intuition that incoherent confidence scores, e.g., a high verbalized confidence in an answer to a question when other distinct answers also have high verbalized confidence, should not be taken at face value. In other words, we should discount high confidence if it does not follow rational coherence norms (Hofweber et al., 2024).

To explain and correct for this kind of incoherence, we first define the notion of suggestibility. Some studies indicate that when LLMs are epistemically uncertain, they tend to rely on their context to resolve the uncertainty (Yadkori et al., 2024; Ahdritz et al., 2024), i.e., the confidence on a claim increases *because* it is in the context. We refer to this phenomenon as *suggestibility* and hypothesize that it contributes to the model's assignment of high confidence to claims it can neither support nor refute. To account for this bias, we propose a method to calibrate verbalized confidence by normalizing by the total verbalized confidence over self-generated distractors (i.e. alternative claims). We generate *minimal pair* distractors using beam search when available, or by directly prompting the model for distractors in the black-box setting. We use an off-the-shelf NLI model to downweight distractors that are similar to other distractors or that do not contradict the main claim.

The approach above for normalizing verbalized confidence with distractors aims to leverage coherence within claim validation, but overlooks another relevant facet of coherence in LLMs. In particular, coherence among sampled generations is correlated with correctness, an observation leveraged by the popular approach of self-consistency (Xiong et al., 2024). Thus, inspired by prior findings on generator-validator disagreement (Li et al., 2023), we integrate these complementary dimensions of coherence into DINCo. Specifically, we use distractor generation and NLI reweighting to estimate and enforce coherence across validations of related claims (e.g. not accepting contradictory claims), while using self-consistency to quantify coherence across sampled generations, upweighting more commonly generated claims.

We test our method on open-source and closed-source models, applied to short-form (TriviaQA and SimpleQA) (Joshi et al., 2017; Wei et al., 2024) and long-form (FActScore; Min et al., 2023) generation domains. DINCo improves ECE over the next best method by an average of 0.099, 0.092, and 0.055, respectively. DINCo effectively extends to long-form biography generation, where it

improves Pearson and Spearman correlation with passage-level FActScore over the best baseline by an average of 0.072 and 0.074, respectively. Further analysis shows that DINCO relieves confidence saturation, and that simply scaling up self-consistency (the strongest baseline overall) does not suffice to match the calibration of DINCO.

2 DISTRACTOR-NORMALIZED COHERENCE (DINCO)

We begin with a motivating hypothesis, supported with preliminary evidence. Then we present the details of our method, illustrated in Fig. 2.

2.1 BACKGROUND AND MOTIVATION

Let $\mathcal C$ be the set of claims with a binary truth value. We denote the truth value of a claim $c \in \mathcal C$ as $v(c) \in \{0,1\}$. A confidence estimation method is a function $f:\mathcal C \to [0,1]$, which is *calibrated* if it correctly predicts the probability of truth. Verbalized confidence is an approach that prompts an LLM to output its confidence $f^{VC}(c)$ in a claim c.¹

For a topic that the model knows little about, it may be willing to adopt the information presented in its context as its prior (Yadkori et al., 2024; Ahdritz et al., 2024), a phenomenon we refer to as *suggestibility*. In other words, the very act of presenting a claim for the model to report its confidence on can bias the reported confidence. For example, if the model does not know who Kang Ji-hwan is, it may assign 60% confidence to both the claim "Kang Ji-hwan was born in 1980." and the claim "Kang Ji-hwan was born in 1990.", even though this seemingly violates the axioms of probability (since the claims are mutually exclusive). Nonetheless, this behavior may not be irrational, as each confidence estimate is conditioned on different information, namely the fact that the respective claim was verbalized in the user prompt; we further discuss the connection between this notion, LLM sycophancy, and human suggestibility in Appendix A.1.

We seek to correct for this bias caused by the model's suggestibility when presented with claims it knows little about. Let f^{VC} be the verbalized confidence, and let f^{lat} be the latent, inaccessible model confidence. We model the bias as a multiplicative scalar $\beta(c)$, which depends on the claim c because the model has varying degrees of uncertainty for different topics: $f^{\text{VC}}(c) = \beta(c) f^{\text{lat}}(c)$. To approximate f^{lat} , we make the assumption that the biases for logically related (e.g. equivalent or contradictory) claims are approximately equal, since they rely on a shared, localized set of knowledge. Let $C \subset \mathcal{C}$ be a set of mutually exclusive and exhaustive claims, e.g. claims for the year a person was born in. Since the claims in C are logically related, we assume that $\beta(c)$ is roughly the same for all $c \in C$ and so there is a scalar $\beta(C)$ with $\beta(C) \approx \beta(c)$ for all $c \in C$. Assuming the latent confidence f^{lat} is probabilistically coherent,

$$1 = \sum_{c \in C} f^{\text{lat}}(c) = \sum_{c \in C} \frac{f^{\text{VC}}(c)}{\beta(c)} \approx \sum_{c \in C} \frac{f^{\text{VC}}(c)}{\beta(C)}.$$
 (1)

Thus, we can approximate $\beta(C)$ and then f^{lat} :

$$\beta(C) \approx \sum_{c \in C} f^{\text{VC}}(c), \qquad f^{\text{NVC}}(c) = \frac{f^{\text{VC}}(c)}{\beta(C)} \approx f^{\text{lat}}(c)$$
 (2)

In practice, we set $\beta(C) \leftarrow \max(1, \beta(C))$ to account for the case where C fails to contain a true claim, or more precisely, a claim that the model believes to be true.

2.2 PRELIMINARY STUDY

To empirically support our motivation above, we plot the distributions of the total confidence ($\beta(C)$ in Eq. 2) over correctly and incorrectly answered questions. In this preliminary study, we take correctness as a proxy for epistemic certainty, although in reality LLMs may still be uncertain about questions they answered correctly. In other words, we treat correct and incorrect instances as epistemically certain and uncertain instances, respectively. Thus, our hypothesis predicts that on incorrect instances, the model would be *more suggestible* and assign high confidence to more answers,

¹We talk about what the model knows, believes, has confidence in, etc. as short-hand notation for the latent ability to produce language similar to that which a human would to demonstrate such knowledge, etc. (Piantadosi & Hill, 2022; West et al., 2023; Hofweber et al., 2024).

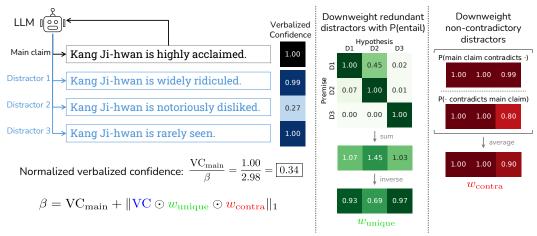


Figure 2: Normalizing verbalized confidence with DINCo. (**Left**) The LLM generates a claim along with several distractors and reports its confidences on them independently. To calibrate the main claim's confidence, we divide it by β , the sum over each distractor's confidence, weighted by uniqueness (**center**) and counterfactuality (**right**). Details in Section 2.3.

leading to *higher* total confidences. On the other hand, if the model is calibrated (in particular, on uncertain instances it exhibits epistemic humility, i.e. recognizing its lack of knowledge), then the total confidence would tend to be around 1 and 0 for correct and incorrect instances, respectively, so the total confidence would be *lower* on incorrect instances. Even if the model is confident in its incorrect answers (e.g. due to misconceptions), we would expect *similar* behavior between correct and incorrect instances.

Experimental Setup. We use TriviaQA (Joshi et al., 2017), a dataset evaluating real-world knowledge with short-form answers. We sample 1000 questions from the validation split of the rc.nocontext subset. We use Qwen3-8B (Yang et al.,)² and take verbalized confidence using P(True) (Kadavath et al., 2022), computed as P(Yes)/(P(Yes) + P(No)) when asking the model whether a given answer is correct (see Appendix B.2 for prompts). In Section 2.3, we specify how we generate dis-

Results. In Fig. 3, we first observe that many incorrect answers are assigned a confidence near 0, suggesting epistemic humility. Even so, the incorrect distribution is heavy-tailed, resulting in a higher mean and median than the correct distribution. These results are consistent with our hypothesis that LLMs are more prone to accepting claims that they are epistemically uncertain about.

tractors and account for answer redundancy.

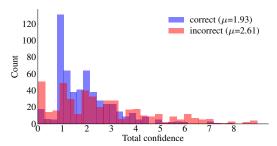


Figure 3: Distributions of total confidence over correct and incorrect answers.

2.3 METHOD

Our preliminary result (Fig. 3) showed that LLMs can produce incoherent probability judgments (i.e. the total confidence $\beta(C)$ exceeds 1), especially when epistemically uncertain, suggesting a need to estimate and correct for this bias. As proposed in Section 2.1, we normalize verbalized confidence by the total confidence $\beta(C)$ over a distractor set C. We now describe in detail how we generate these distractors, and how we account for redundancy among distractors, a process illustrated in Fig. 2. Finally, we discuss how the phenomenon of generator-validator disagreement motivates incorporating self-consistency into DINCO.

²Throughout, we use the instruction-tuned versions of Qwen3 models.

Distractor generation. The distractor set should contain enough plausible distractors to avoid underestimating the normalization factor $(\beta(C))$ in Eq. 2), while remaining small enough to be tractably computed. Thus, we frame the problem of choosing the optimal distractor set containing an original claim c_0 as maximizing the total acceptance probability $\sum_{c \in C} f^{\text{VC}}(c)$ subject to a size constraint $|C| \leq K$; shortly we address relaxing the requirement of mutual exclusivity. Unfortunately, the validation probability $f^{\text{VC}}(c)$ can only be elicited on a per-claim basis, leading to intractable sample complexity. Motivated by the intuition that an LLM tends to generate claims that it would find plausible during validation, as a proxy for the set of claims with high verbalized confidence, we use the set of claims with high generation probability (under an appropriate prompt, explained below).

To encourage mutual exclusivity among claims, we set up the claims to be minimal pairs (Fig. 2 left). For a given short-form question, we can simply sample many answers, but independent sampling is inefficient as it overrepresents high probability generations (Gekhman et al., 2025), limiting the number of unique distractors. Instead, we use beam search (Sutskever et al., 2014) when available to efficiently identify unique sequences that have high probability mass coverage. For API-access models, we use top token probabilities if available to implement a pseudo-beam search (see Appendix C.1 for details). Otherwise in the black-box setting, we directly prompt the model to generate a list of candidate answers (see Appendix B.2 for prompts). For long-form QA, we follow Min et al. (2023) in decomposing a long generation into claims. We separately prompt the model to generate one distractor for a given claim, and again use beam search to create multiple distractors.

Addressing Claim Redundancy. Although the heuristic of generating minimal pairs encourages mutually exclusive claims, we have little guarantee of this mutual exclusivity (1) within the distractor set C and (2) between distractors and the original claim. Assuming such mutual exclusivity when there are actually redundant claims can lead to overcounting in the normalization factor $\beta(C)$. Thus, we use an NLI model to quantify entailment and contradiction relationships between claims. We address (1) and (2) with w_{unique} and w_{contra} , respectively (Fig. 2):

address (1) and (2) with
$$w_{\text{unique}}$$
 and w_{contra} , respectively (Fig. 2):
$$w_{\text{unique}}(c) = \frac{1}{\sum_{c' \in C} P(\text{entail} \mid c', c)}, \quad w_{\text{contra}}(c) = \frac{P(\text{contra} \mid c_0, c) + P(\text{contra} \mid c, c_0)}{2} \quad (3)$$

Intuitively, $w_{\rm unique}$ downweights a claim if it is entailed by other claims, and $w_{\rm contra}$ downweights a claim if it is not contradictory with the main claim. We now normalize the verbalized confidence of the original claim as

$$f^{\text{NVC}}(c_0) = \frac{f^{\text{VC}}(c_0)}{\beta(C)}, \quad \beta(C) = \max\left(1, f^{\text{VC}}(c_0) + \sum_{c \in C} f^{\text{VC}}(c) \cdot w_{\text{unique}}(c) \cdot w_{\text{contra}}(c)\right), \tag{4}$$

generalizing the mutually exclusive case in Eq. 2. The maximization with 1 allows for defaulting back to the vanilla verbalized confidence in the case where C fails to contain claims that the model considers plausible.

Combining Coherence within Generation and Validation. Our approach so far (summarized in Fig. 2) for normalizing verbalized confidence across distractors focuses on coherence within claim *validation*. Previous studies disagree on the question of whether models are better at the discriminative or generative counterparts of a given task (West et al., 2023; Gekhman et al., 2025), but generally agree on the presence of generator-validator inconsistency (Li et al., 2023), wherein a model may produce inconsistent results between the generation and validation stages. Indeed, we find that in the preliminary study setting in Section 2.2, the answer with the highest generation probability and the answer with the highest validation probability (over 10 answers obtained from beam search) agree⁴ on only 592 out of 1000 questions.

We integrate these complementary aspects of generation and validation into DINCo. In particular, we draw on a distributional view of the generator-validator gap (Rodriguez et al., 2025), in which the generation probability distribution over candidate answers (which we approximate with self-consistency sampling, f^{SC}) is distinct from the validation probability distribution over

³Access to an NLI model poses only a minimal departure from the zero-resource setting, since NLI is a generic task for which there are off-the-shelf models. NLI is a subset of the tasks that LLMs are capable of, so the usage of a separate NLI model is motivated merely by efficiency (Kuhn et al., 2023; Lin et al., 2024).

⁴We consider c and c' equivalent answers to question q if $\frac{1}{2}P(\text{entail} \mid c, c'; q) + \frac{1}{2}P(\text{entail} \mid c', c; q) > 0.9$.

them (which we approximate with normalized verbalized confidence, $f^{\rm NVC}$). DINCO thus incorporates confidence (i.e. probability mass) in both the generator and validator distributions: $f^{\rm DiNCo}(c) = \frac{1}{2} f^{\rm SC}(c) + \frac{1}{2} f^{\rm NVC}(c)$. We leave a full description of the self-consistency component to Appendix C.2.

3 EXPERIMENTS

3.1 EXPERIMENTAL SETUP

Short-form Datasets. Short-form QA serves as a testbed for evaluating factuality as well as calibration because of its tractable evaluation and adjustable difficulty. The task is relevant in practice because it assesses models' ability to respond to information-seeking users. TriviaQA contains trivia questions requiring diverse world knowledge (Joshi et al., 2017). SimpleQA similarly contains short, fact-seeking questions, curated with the criterion of challenging frontier models (Wei et al., 2024). We sample 1000 questions from each dataset.⁵ We use LLM-as-a-judge to evaluate binary correctness, following best practices for robust evaluation (Wei et al., 2024); in Appendix B.3 we confirm high human agreement.

Long-form Datasets. While short-form settings are appealing for their simple evaluation, many real-world tasks require longer generations, for which calibrated confidence estimation remains critical. The long-form setting comes with the evaluation challenge that responses generally contain both correct and incorrect parts, complicating the assignment of a single correctness score. In our experiments, we evaluate long-form calibration on biography generation using FActScore (Min et al., 2023). FActScore decomposes a generated biography into atomic claims and verifies each claim against Wikipedia (see Appendix B.4 for an example), thus enabling us to evaluate calibration at the claim level. We use the labeled subset containing 183 entities from Min et al. (2023).

Models. Since TriviaQA and SimpleQA are adequately challenging for smaller and larger models, respectively, we focus their evaluation accordingly. TriviaQA is largely solved by frontier models (Wei et al., 2024), and SimpleQA is too difficult for smaller models.⁶ On TriviaQA, we use popular open-source models: Qwen3-8B and Qwen3-1.7B (Yang et al., 2025), Llama-3.2-3B-Instruct (Grattafiori et al., 2024), and Gemma-3-4B-IT (Team et al., 2025). On SimpleQA, we use popular frontier models: GPT-4.1 (2025-04-14; OpenAI, 2025) and Gemini-2.5-Flash (Comanici et al., 2025b). For SimpleQA evaluated on frontier models, we also consider the black-box setting where no logit access is assumed; here, rather than using pseudo-beam search to generate distractors, we prompt the model directly to generate diverse distractors. Moreover, we replace P(True) with verbalized numerical confidence to forgo logit access. For the long-form task of biography generation, we limit our evaluation to Qwen3-8B and Gemma-3-4B-IT, due to the cost of FActScore evaluation with GPT-4.1. We use the NLI model DeBERTa-v3-base-mnli-fever-anli (He et al., 2021) for our methods and self-consistency.

Evaluation Metrics. We evaluate Expected Calibration Error (**ECE** \downarrow ; Naeini et al., 2015) with 10 bins, Brier score (**BS** \downarrow ; Brier, 1950), and area under the ROC curve (**AUC** \uparrow ; Hanley & McNeil, 1982). See Appendix B.1 for descriptions. For biography generation, we also evaluate Pearson and Spearman correlation between average claim-level confidence and passage-level FActScore (without length penalty), where the latter measures the proportion of claims that are correct.

Baselines. Following past work in zero-resource calibration, we compare against training-free methods that produce probabilities without post-hoc calibration (Tian et al., 2023b; Xiong et al., 2024; Steyvers et al., 2025). We provide prompts for our methods and baselines in Appendix B.2. Verbalized confidence (\mathbf{VC} ; P(True) from Kadavath et al., 2022) asks the model whether its answer is correct and computes $P(\mathrm{Yes})/(P(\mathrm{Yes})+P(\mathrm{No}))$. It is straightforward to replace P(True) with verbalized numerical confidence; in Appendix D.1 we show that the latter similarly benefits from our methods, showing robustness to the format of verbalized confidence. Top-K prompting (K-VC; Verb. 1S from Tian et al., 2023a) prompts the model to provide its top K guesses along with verbalized numerical confidences. For our black-box setting, we use the candidate answers generated

⁵We use the validation split of the rc.nocontext subset for TriviaQA. On SimpleQA, Gemini-2.5-Flash produced a refusal error with no output on 86 questions, so we exclude them from all experiments.

⁶https://www.kaggle.com/benchmarks/openai/simpleqa

Table 1: TriviaQA results. We evaluate Expected Calibration Error (ECE), Brier score (BS), and area under the ROC curve (AUC). In each column, we bold the best result and underline results not significantly worse under a paired test ($\alpha = 0.05$; see Appendix B.6 for tests).

	Ç	wen3-8	BB	Q	Qwen3-1.7B			Llama-3.2-3B-Instruct			Gemma-3-4B-IT		
Method	ECE ↓	BS ↓	AUC ↑	ECE ↓	BS ↓	AUC ↑	ECE ↓	BS ↓	AUC ↑	ECE ↓	BS ↓	AUC ↑	
VC (Kadavath et al., 2022)	0.240	0.242	0.817	0.387	0.383	0.720	0.189	0.208	0.826	0.300	0.299	0.702	
K-VC (Tian et al., 2023a)	0.341	0.348	0.604	0.538	0.524	0.596	0.146	0.228	0.678	0.254	0.262	0.786	
MSP (Fadeeva et al., 2023)	0.149	0.203	0.819	0.104	0.186	0.774	0.243	0.253	0.764	0.252	0.268	0.790	
SC-VC (Xiong et al., 2024)	0.299	0.325	0.704	0.451	0.474	0.559	0.122	0.211	0.761	0.362	0.378	0.653	
SC (Xiong et al., 2024)	0.236	0.244	0.785	0.233	0.229	0.780	0.065	0.177	0.808	0.303	0.304	0.713	
NVC	0.171	0.190	0.853	0.084	0.164	0.806	0.168	0.192	0.845	0.218	0.236	0.791	
DINCO	0.097	0.155	0.879	0.177	0.179	0.835	0.044	0.148	0.864	0.121	0.191	0.817	

using the K-VC prompt as distractors, but we discard the verbalized confidences contained in the same generation, and instead separately collect verbalized numerical confidence on each distractor independently. Maximum sequence probability (MSP; Fadeeva et al., 2023) is the LLM's probability of generating its answer. Self-consistency (SC; Xiong et al., 2024) samples several answers (we use temperature 1) and computes the proportion that match the main answer. Following Kuhn et al. (2023), we use an NLI model to determine semantic equivalence. For biography generation, we modify self-consistency to sample several biographies and measure entailment of each claim (Zhang et al., 2024). SC-VC Xiong et al. (2024) is SC weighted by verbalized confidence (we use P(True) following Taubenfeld et al., 2025).

Inference Budget. As our methods and several baseline methods (K-VC, SC, SC-VC) operate on a variable inference-time budget, we control this budget at K=10. For DINCO, we use 5 samples for self-consistency and 5 distractors for normalized verbalized confidence.

3.2 RESULTS

Short-form QA. Tables 1 and 2 report results for TriviaQA and SimpleQA. On TriviaQA, DINCO outperforms the best baseline, SC, by an average ECE of 0.099. On SimpleQA, DINCo outperforms the best baseline, MSP, by an average ECE of 0.092. While the better baseline varies by dataset, DINCo consistently achieves the strongest calibration. MSP is sometimes a competitive baseline (e.g. AUC 0.800 surpasses DINCO at 0.786 on SimpleQA with GPT-4.1), but this often does not hold across metrics (e.g. MSP has an ECE of 0.263 in the same setting, heavily underperforming DINCo at 0.089) or across settings (e.g. MSP underperforms DINCo on AUC in TriviaQA by an average of 0.062). Moreover, we note that the effectiveness of MSP relies on answers having a canonical form, preventing its generalization to long-form settings (Farquhar et al., 2024). We highlight that NVC outperforms SC (e.g. by an ECE of 0.049 and 0.060 on TriviaQA and SimpleQA, respectively) despite only leveraging coherence in validation and not in generation (Section 2.3). Nonetheless, DINCo is more consistently calibrated than NVC (in particular on AUC, e.g. 0.786 DINCO vs. 0.729 NVC with GPT-4.1 on SimpleQA), empirically supporting our motivation in Section 2.3 for integrating coherence in generation (SC) and validation (NVC) into DINCo. In the black-box setting on SimpleQA, DINCo continues to do well (e.g. outperforming the baselines on ECE), but it tends to fall behind DINCo with logit access, underscoring the benefit of leveraging token probabilities for calibration.

Long-form QA. Table 3 reports results on FActScore. While VC is extremely miscalibrated (e.g. ECE of 0.433 with Qwen3-8B), DINCo is able to leverage incoherence across related claims to normalize verbalized confidence and achieve strong calibration. Whether SC or NVC performs better varies by the model Qwen3-8B or Gemma-3-4B-IT, but DINCo continues to outperform SC (0.076 vs. 0.162 ECE with Qwen3-8B, and 0.172 vs. 0.197 ECE with Gemma-3-4B-IT). Furthermore, DINCo is the method most strongly correlated with passage-level FActScore (e.g. improving Pearson and Spearman correlation over SC by an average of 0.072 and 0.074, respectively), demonstrating that the effectiveness of DINCo extends to the long-form setting. Taken together with the short-form results in Tables 1 and 2, these results indicate that DINCo is applicable to open- and

Table 2: SimpleQA results. The black-box variants of our methods assume no logit access. Metrics and text styling follow Table 1.

		GPT-4.1		Gemini-2.5-Flash			
Method	ECE ↓	BS ↓	AUC↑	ECE ↓	BS ↓	AUC ↑	
VC (Kadavath et al., 2022)	0.547	0.549	0.644	0.409	0.393	0.617	
K-VC (Tian et al., 2023a)	0.338	0.337	0.632	0.535	0.511	0.566	
MSP (Fadeeva et al., 2023)	0.263	0.255	0.800	0.098	0.177	0.773	
SC-VC (Xiong et al., 2024)	0.223	0.252	0.761	0.186	0.221	0.755	
SC (Xiong et al., 2024)	0.220	0.252	0.750	0.170	0.212	0.748	
NVC _{black-box}	0.213	0.270	0.607	0.208	0.262	0.595	
DINCO _{black-box}	0.161	0.251	0.605	0.079	0.199	0.697	
NVC	0.164	0.222	0.729	0.105	0.199	0.662	
DINCO	0.089	0.183	0.786	0.088	0.174	0.762	

Table 3: FActScore results. In addition to the claim-level metrics, we report Pearson (r) and Spearman (ρ) correlation with passage-level FActScore. Text styling follows Table 1, and we bold the best r and ρ .

	Qwen3-8B					Gemma-3-4B-IT					
Method	ECE ↓	BS ↓	AUC ↑	$r \uparrow$	ρ \uparrow	ECE ↓	BS ↓	AUC ↑	$r \uparrow$	ρ \uparrow	
VC (Kadavath et al., 2022)	0.433	0.431	0.625	0.073	0.122	0.527	0.527	0.683	-0.081	-0.129	
SC (Zhang et al., 2024)	0.162	0.226	0.771	0.468	0.494	0.197	0.233	0.787	0.629	0.607	
NVC	0.191	0.263	0.681	0.444	0.443	0.123	0.230	0.726	0.695	0.704	
DiNCo	0.076	0.202	0.767	0.518	0.538	0.172	0.210	0.793	0.724	0.712	

closed-source models, and crucially can transfer seamlessly between short-form QA and long-form generation settings.

DISCUSSION AND ANALYSIS

We conduct further analysis using Qwen3-8B on TriviaQA with P(True), which appeared in Table 1 in the main experiments.

Scaling Self-Consistency. While our main experiments in Section 3 were controlled for the inference budget of each method, here we examine whether simply scaling the inference budget of self-consistency allows it to recover the calibration of DINCo. Fig. 4 shows that self-consistency alone is unable to reach the performance of DINCo (using 5 distractors and 5 self-consistency samples as in Section 3) even when scaling up to 100 samples, demonstrating that the effectiveness of DINCo comes from leveraging coherence in both generation and validation, which is not matched by scaling the generation axis alone.

Quantifying Saturation. A core motivation for our method is the notion (shown in Fig. 1) that verbalized confidence exhibits saturation at high confidence. To better quantify this notion, we

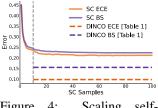


Figure 4: Scaling selfconsistency does not close the gap with DINCO.

(higher Δ means lower sat- weighting. uration). DINCO alleviates saturation.

Method	Δ_0	$\Delta_{0.001}$
VC	0.670	0.605
SC	0.734	0.734
SC@100	0.832	0.832
DINCO	0.998	0.984

Table 4: Saturation analysis Table 5: Ablation of NLI-based

Method	ECE ↓	BS ↓	AUC ↑
NVC	0.171	0.190	0.853
w/o NLI	0.358	0.335	0.778
DiNCo	0.097	0.155	0.879
w/o NLI	0.130	0.185	0.810

introduce a metric Δ_ϵ that measures the absence of saturation. We define Δ_ϵ as the proportion of pairs of distinct instances that have a confidence difference exceeding ϵ . For example, if all confidence scores are the same, then $\Delta_0=0$, and if all confidence scores are distinct, then $\Delta_0=1$. We consider $\epsilon\in\{0,0.001\}$. Table 4 shows that DINCo leads to substantially higher rates of distinct confidence, indicating lower saturation. In particular, self-consistency scaled to 100 samples (as above) continues to be more saturated than DINCo. While the absence of saturation alone means little without calibration (as evaluated in Section 3), this analysis helps explain DINCo's calibration improvement, as hinted at by Fig. 1. Moreover, we argue that lower saturation leads to more usable confidence estimates: a saturated distribution is inherently less controllable, with large jumps in error between thresholds.

Ablating NLI. Our main experiments in Tables 1 to 3 showed comparisons with SC and NVC, which ablate NVC and SC, respectively, from DINCO. Here, to understand how necessary access to an NLI model is, we ablate the NLI component, which is used to downweight distractors that overlap with the main claim or other distractors (Section 2.3, Fig. 2). In Table 5, we see that performance substantially decreases without NLI-based weighting, emphasizing the utility of an off-the-shelf NLI model to account for claim redundancy.

5 RELATED WORK

Considering Multiple Answers for Verbalized Confidence. The approach most related to DINCo is to have the model consider several candidates within a single prompt and assign confidences to them (Tian et al., 2023b; Kadavath et al., 2022; Wang et al., 2024). Similarly, Chhikara (2025) samples answers and presents them in the context as distractors. A subtle but crucial distinction between these methods and DINCo is that if we present all the candidates together, we become unable to gauge the probabilistic coherence of the confidence estimates, i.e. whether they form a valid probability distribution, since LLMs can satisfy probabilistic coherence via simple arithmetic. In Section 2.2, we show that the probabilistic coherence of confidence estimates is correlated with answer correctness. Instruction-tuned LLMs have a tendency to assert confidence even when it is undue (OpenAI, 2023; Leng et al., 2025; Sun et al., 2025; Xu et al., 2025), leading joint prompting to suffer similar issues of overconfidence as vanilla prompting. With independent prompting, we can expose and account for inconsistencies in self-declared knowledge. In Section 3, we empirically verify that our method leads to better calibration than joint prompting for verbalized confidence.

Confidence Estimation in Long-form Generation. While historically confidence estimation has mostly been applied to classification (Houlsby et al., 2011), multiple-choice (Jiang et al., 2021), and short-form QA (Xiong et al., 2024), with the advent of LLMs it has increasingly been considered for long-form generation. Since token-level uncertainty is often ill-suited for representing claim-level uncertainty, the primary approaches have been self-consistency and verbalized confidence (Manakul et al., 2023; Zhang et al., 2024; 2025). We propose a method to normalize verbalized confidence, enabling DINCo to combine the complementary confidence signals in these two prior approaches.

Reconciling Inconsistent LLM Probability Judgments. In Appendix A.2, we discuss a related line of work demonstrating the benefits of reconciling inconsistent LLM probability judgments. In this work, we propose a zero-resource confidence estimator that normalizes verbalized confidence over self-generated distractors, motivated by the suggestibility of LLMs in unfamiliar topics.

6 CONCLUSION

We present DINCo, which estimates LLM confidence by leveraging coherence in generation as well as validation (through verbalized confidence). Verbalized confidence tends to be saturated at high confidence. We show evidence that this behavior is correlated with *suggestibility*, where the LLM is more likely to accept claims that it knows less about. Motivated by this finding, DINCo has the LLM verbalize its confidence independently on several self-generated distractors to estimate and correct for the bias caused by suggestibility. DINCo outperforms existing methods in the zero-resource setting on short-form QA (TriviaQA, SimpleQA) and long-form QA (FActScore) and mitigates saturation.

ETHICS STATEMENT

Our work aims to make AI safer by addressing calibration, an important component of safety. We do not foresee any additional ethical implications beyond standard ethical and safety considerations that apply to AI research generally.

REPRODUCIBILITY STATEMENT

To further reproducibility, we have included our code in the supplementary. We have also provided descriptions of our method and our prompts in Section 2 and Appendix B.

REFERENCES

- Gustaf Ahdritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L. Edelman. Distinguishing the knowable from the unknowable with language models, 2024. URL https://arxiv.org/abs/2402.03563.
- Glenn W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78:1–3, 1950. URL https://api.semanticscholar.org/CorpusID:122906757.
- Maggie Bruck and Stephen Ceci. The suggestibility of children's memory. *Annual review of psychology*, 50:419–39, 02 1999. doi: 10.1146/annurev.psych.50.1.419.
- Prateek Chhikara. Mind the confidence gap: Overconfidence, calibration, and distractor effects in large language models, 2025. URL https://arxiv.org/abs/2502.11028.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, 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, 2025a. URL https://arxiv.org/abs/2507.06261.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, 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, 2025b. URL https://arxiv.org/abs/2507.06261.
- Elizabeth R. DeLong, David M. DeLong, and Daniel L. Clarke-Pearson. Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach. *Biometrics*, 44(3):837–845, 1988. ISSN 0006341X, 15410420. URL http://www.jstor.org/stable/2531595.
- Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, Timothy Baldwin, and Artem Shelmanov. Lm-polygraph: Uncertainty estimation for language models, 2023. URL https://arxiv.org/abs/2311.07383.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630:625 630, 2024. URL https://api.semanticscholar.org/CorpusID:270615909.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration, 2024. URL https://arxiv.org/abs/2402.00367.
- Yu Feng, Ben Zhou, Weidong Lin, and Dan Roth. Bird: A trustworthy bayesian inference framework for large language models, 2025. URL https://arxiv.org/abs/2404.12494.
- Zorik Gekhman, Eyal Ben David, Hadas Orgad, Eran Ofek, Yonatan Belinkov, Idan Szpektor, Jonathan Herzig, and Roi Reichart. Inside-out: Hidden factual knowledge in llms, 2025. URL https://arxiv.org/abs/2503.15299.

- Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koeppl, Preslav Nakov, and Iryna Gurevych. A survey of confidence estimation and calibration in large language models, 2024. URL https://arxiv.org/abs/2311.08298.
- 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, 2024. URL https://arxiv.org/abs/2407.21783.
- James A. Hanley and Barbara J. McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143:29–36, 05 1982. doi: 10.1148/radiology.143.1. 7063747.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention, 2021. URL https://arxiv.org/abs/2006.03654.
- Thomas Hofweber, Peter Hase, Elias Stengel-Eskin, and Mohit Bansal. Are language models rational? the case of coherence norms and belief revision. *arXiv preprint arXiv:2406.03442*, 2024.
- Bairu Hou, Yang Zhang, Jacob Andreas, and Shiyu Chang. A probabilistic framework for llm hallucination detection via belief tree propagation, 2025. URL https://arxiv.org/abs/2406.06950.
- Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning, 2011. URL https://arxiv.org/abs/1112.5745.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How can we know when language models know? on the calibration of language models for question answering, 2021. URL https://arxiv.org/abs/2012.00955.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017. URL https://arxiv.org/abs/1705.03551.
- Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, Ronan Le Bras, and Yejin Choi. Maieutic prompting: Logically consistent reasoning with recursive explanations, 2022. URL https://arxiv.org/abs/2205.11822.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know, 2022. URL https://arxiv.org/abs/2207.05221.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation, 2023. URL https://arxiv.org/abs/2302.09664.
- Philippe Laban, Lidiya Murakhovs'ka, Caiming Xiong, and Chien-Sheng Wu. Are you sure? challenging llms leads to performance drops in the flipflop experiment, 2024. URL https://arxiv.org/abs/2311.08596.
- Jixuan Leng, Chengsong Huang, Banghua Zhu, and Jiaxin Huang. Taming overconfidence in llms: Reward calibration in rlhf, 2025. URL https://arxiv.org/abs/2410.09724.
- Xiang Lisa Li, Vaishnavi Shrivastava, Siyan Li, Tatsunori Hashimoto, and Percy Liang. Benchmarking and improving generator-validator consistency of language models, 2023. URL https://arxiv.org/abs/2310.01846.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words, 2022. URL https://arxiv.org/abs/2205.14334.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models, 2024. URL https://arxiv.org/abs/2305.19187.

- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models, 2023. URL https://arxiv.org/abs/2303.08896.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation, 2023. URL https://arxiv.org/abs/2305.14251.
- Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. *Proceedings of the ... AAAI Conference on Artificial Intelligence*. AAAI Conference on Artificial Intelligence, 2015:2901–2907, 2015. URL https://api.semanticscholar.org/CorpusID:6292807.
- Aliakbar Nafar, Kristen Brent Venable, Zijun Cui, and Parisa Kordjamshidi. Extracting probabilistic knowledge from large language models for bayesian network parameterization, 2025. URL https://arxiv.org/abs/2505.15918.
- OpenAI. Gpt-4 technical report, 2023. URL https://arxiv.org/abs/2303.08774.
- OpenAI. Introducing gpt-4.1 in the api. https://openai.com/index/gpt-4-1/, April 2025. Accessed: 2025-09-23.
- Steven T. Piantadosi and Felix Hill. Meaning without reference in large language models, 2022. URL https://arxiv.org/abs/2208.02957.
- Juan Diego Rodriguez, Wenxuan Ding, Katrin Erk, and Greg Durrett. Rankalign: A ranking view of the generator-validator gap in large language models, 2025. URL https://arxiv.org/abs/2504.11381.
- Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z. Ren, and Anirudha Majumdar. A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions, 2025. URL https://arxiv.org/abs/2412.05563.
- Mark Steyvers, Heliodoro Tejeda, Aakriti Kumar, Catarina Belem, Sheer Karny, Xinyue Hu, Lukas W. Mayer, and Padhraic Smyth. What large language models know and what people think they know. *Nature Machine Intelligence*, 7(2):221–231, January 2025. ISSN 2522-5839. doi: 10.1038/s42256-024-00976-7. URL http://dx.doi.org/10.1038/s42256-024-00976-7.
- Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. Large language models are overconfident and amplify human bias, 2025. URL https://arxiv.org/abs/2505.02151.
- Xu Sun and Weichao Xu. Fast implementation of delong's algorithm for comparing the areas under correlated receiver operating characteristic curves. *IEEE Signal Processing Letters*, 21(11):1389–1393, 2014. doi: 10.1109/LSP.2014.2337313.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks, 2014. URL https://arxiv.org/abs/1409.3215.
- Amir Taubenfeld, Tom Sheffer, Eran Ofek, Amir Feder, Ariel Goldstein, Zorik Gekhman, and Gal Yona. Confidence improves self-consistency in llms. In *Findings of the Association for Computational Linguistics:* ACL 2025, pp. 20090–20111. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.findings-acl.1030. URL http://dx.doi.org/10.18653/v1/2025.findings-acl.1030.
- 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, 2025. URL https://arxiv.org/abs/2503.19786.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea Finn. Fine-tuning language models for factuality, 2023a. URL https://arxiv.org/abs/2311.08401.

- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback, 2023b. URL https://arxiv.org/abs/2305.14975.
 - D.N. Walton. *Arguments from Ignorance*. Arguments from Ignorance. Pennsylvania State University Press, 1996. ISBN 9780271014746. URL https://books.google.com/books?id=peTWAAAAMAAJ.
 - Cheng Wang, Gyuri Szarvas, Georges Balazs, Pavel Danchenko, and Patrick Ernst. Calibrating verbalized probabilities for large language models, 2024. URL https://arxiv.org/abs/2410.06707.
 - Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023. URL https://arxiv.org/abs/2203.11171.
 - Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models, 2024. URL https://arxiv.org/abs/2411.04368.
 - Peter West, Ximing Lu, Nouha Dziri, Faeze Brahman, Linjie Li, Jena D. Hwang, Liwei Jiang, Jillian Fisher, Abhilasha Ravichander, Khyathi Chandu, Benjamin Newman, Pang Wei Koh, Allyson Ettinger, and Yejin Choi. The generative ai paradox: "what it can create, it may not understand", 2023. URL https://arxiv.org/abs/2311.00059.
 - Shepard Xia, Brian Lu, and Jason Eisner. Let's think var-by-var: Large language models enable ad hoc probabilistic reasoning, 2024. URL https://arxiv.org/abs/2412.02081.
 - Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can Ilms express their uncertainty? an empirical evaluation of confidence elicitation in Ilms, 2024. URL https://arxiv.org/abs/2306.13063.
 - Chenjun Xu, Bingbing Wen, Bin Han, Robert Wolfe, Lucy Lu Wang, and Bill Howe. Do language models mirror human confidence? exploring psychological insights to address overconfidence in llms, 2025. URL https://arxiv.org/abs/2506.00582.
 - Yasin Abbasi Yadkori, Ilja Kuzborskij, András György, and Csaba Szepesvári. To believe or not to believe your llm, 2024. URL https://arxiv.org/abs/2406.02543.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
 - Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. Luq: Long-text uncertainty quantification for llms, 2024. URL https://arxiv.org/abs/2403.20279.
 - Caiqi Zhang, Xiaochen Zhu, Chengzu Li, Nigel Collier, and Andreas Vlachos. Reinforcement learning for better verbalized confidence in long-form generation, 2025. URL https://arxiv.org/abs/2505.23912.
 - Jian-Qiao Zhu and Thomas L. Griffiths. Incoherent probability judgments in large language models, 2025. URL https://arxiv.org/abs/2401.16646.
 - Jian-Qiao Zhu, Adam Sanborn, and Nick Chater. Bayesian inference causes incoherence in human probability judgments, 11 2018.

A RELATED WORK

A.1 RELATED BEHAVIOR IN LLMS AND HUMANS

Humans are also known to be susceptible to suggestibility. They can alter their memories to match the suggestions of other people, especially at a young age (Bruck & Ceci, 1999). Sycophancy is a similar behavior observed in LLMs, where an epistemically vacuous prompt such as "Are you sure?" often leads the model to change its answer (Laban et al., 2024). The user expressing doubt suggests to the model that its answer may be incorrect, since the user has some assumed level of credibility and would be unlikely to ask again if they agreed. Since instruction-tuned models aim to adhere to user preferences, it is plausible that they would employ an argument from ignorance (Walton, 1996) to accept a user claim that they cannot refute.

Zhu & Griffiths (2025) provide evidence of probabilistic incoherence in LLMs and attribute this finding to the prior from the Bayesian Sampler model, which has been used to explain incoherence in human probability judgments (Zhu et al., 2018). In particular, if the same prior is used for every probability judgment, the sum of the probability judgments for mutually exclusive events can exceed 1 (Zhu & Griffiths, 2025). This failure to satisfy the axioms of probability is consistent with our empirical evidence in Section 2.2.

A.2 RECONCILING INCONSISTENT LLM PROBABILITY JUDGMENTS.

Prior works have demonstrated the benefits of reconciling inconsistent LLM probability judgments instead of taking them at face value. Jung et al. (2022) improve factuality by selecting claims to which the LLM assigns coherent truth values upon negation. Hou et al. (2025) use belief tree propagation with logically related claims to detect hallucinations. Nafar et al. (2025) find that independent prompting followed by normalization outperforms joint prompting for Bayesian network parameter estimation. Feng et al. (2025); Xia et al. (2024) optimize a probability distribution to approximately satisfy LLM-generated probability constraints. In this work, we propose a zero-resource confidence estimator that normalizes verbalized confidence over self-generated distractors, motivated by the suggestibility of LLMs in unfamiliar topics.

B EXPERIMENTAL SETUP

B.1 EVALUATION METRICS

We adopt the notation from Section 2.1. For a claim c, the truth value is $v(c) \in \{0,1\}$ and the assigned confidence is $f(c) \in [0,1]$.

Expected Calibration Error (ECE; Naeini et al., 2015). The confidence space [0,1] is partitioned into K intervals of equal length. Out of the N claims in the dataset, let B_k be the list of claims assigned a confidence in the interval $I_k = \left(\frac{k-1}{K}, \frac{k}{K}\right]$ (with I_1 including 0). We compute bin-level truthfulness and confidence as

bin-level truthfulness and confidence as
$$\bar{v}_k = \frac{1}{|B_k|} \sum_{c \in B_k} v(c), \qquad \qquad \bar{f}_k = \frac{1}{|B_k|} \sum_{c \in B_k} f(c). \tag{5}$$

ECE is the bin size-weighted average of the absolute differences:

$$ECE = \sum_{k=1}^{K} \frac{|B_k|}{N} |\bar{v}_k - \bar{f}_k|$$
 (6)

Brier score (BS; Brier, 1950). For a dataset of claims c_1, \ldots, c_N , the Brier score is the mean squared error between the truth values and the confidence estimates:

$$BS = \frac{1}{N} \sum_{i=1}^{N} (v(c_i) - f(c_i))^2$$
(7)

Area under ROC curve (AUC; Hanley & McNeil, 1982) The ROC curve (depicted in Fig. 1 right) captures the tradeoffs between true and false positive rate (TPR, FPR) that we can obtain

with selective prediction, i.e. setting a confidence threshold above which to accept claims. We take correct and incorrect claims to be labeled positive and negative, respectively. The TPR is the proportion of positive instances that are accepted, and the FPR is the proportion of negative instances that are accepted. We want the TPR to be high but the FPR to be low. By setting a lower confidence threshold, TPR will be higher, but FPR may also be higher. By setting a higher confidence threshold, FPR will be lower, but TPR may also be lower. As not every TPR (or FPR) in the interval [0,1] may be achievable, the ROC fills in the gaps between achievable (TPR, FPR) tradeoffs with linear interpolations, as seen with verbalized confidence in Fig. 1. As a summary statistic for the selective predictive power that a confidence estimator grants us, we compute the area under the ROC curve (AUC). For example, if all positive instances are assigned a higher confidence than all negative instances, the AUC is 1. Meanwhile, if confidences are sampled independently at random from the same distribution, the expected AUC is 0.5.

The AUC can also be characterized as the probability that a random positive instance is assigned higher confidence than a random negative instance, with ties randomly broken. Denoting C_+ and C_- as the list of correct and incorrect claims in the dataset, respectively,

$$AUC = \frac{1}{|C_{+}||C_{-}|} \sum_{c_{+} \in C_{+}} \sum_{c_{-} \in C_{-}} \frac{\mathbf{1}\{f(c_{+}) \ge f(c_{-})\} + \mathbf{1}\{f(c_{+}) > f(c_{-})\}}{2}.$$
 (8)

B.2 PROMPTS

B.2.1 SHORT-FORM QA

Prompt to generate main answer. Also used for beam search in DINCo and sampling for self-consistency.

Here are 2 sets of example prompt and answer.

Example Prompt: Which American-born Sinclair won the Nobel Prize for Literature in

→ 1930?

Example Answer: Sinclair Lewis

Example Prompt: Where in England was Dame Judi Dench born?

Example Answer: York

Now, here is a new prompt to answer. Answer with a concise phrase, as in the examples.

Prompt: {question}

Answer:

P(True)

Below is a question and a candidate answer. Your task is to determine whether the answer is correct or not. Only output "Yes" (correct) or "No" (incorrect).

Question: {question}

Candidate answer: {candidate_answer}

Verbalized numerical confidence

Below is a question and a candidate answer. State your confidence that the candidate answer is correct. Only output an integer followed by "%".

Question: {question}

Candidate answer: {candidate_answer}

K-VC (Verb. 1S from Tian et al., 2023b). We use K = 10.

Provide your {K} best guesses and the probability that each is correct (0.0 to 1.0) for the following question. Give ONLY the guesses and probabilities, no other words or

→ explanation. For example:

G1: <first most likely guess, as short as possible; not a complete sentence, just the guess!>

P1: <the probability between 0.0 and 1.0 that G1 is correct, without any extra commentary

→ whatsoever; just the probability!>

 $G\{K\}: <\{K\}$ th most likely guess, as short as possible; not a complete sentence, just the \hookrightarrow guess!>

 $P\{\texttt{K}\}: < \text{the probability between } 0.0 \text{ and } 1.0 \text{ that } G\{\texttt{K}\} \text{ is correct, without any extra} \\ \hookrightarrow \text{ commentary whatsoever; just the probability!} >$

The question is: {question}

Follow-up after main answer for SC-VC

Is your answer correct? Only output "Yes" or "No".

B.2.2 BIOGRAPHY GENERATION

Prompt to generate main biography. Also used to sample biographies for self-consistency.

Write me a paragraph biography on {entity}.

Prompt to generate one distractor. We use beam search to extract multiple.

You will be given a fact about a person. Assuming the fact is accurate, your task is to

- → generate a plausible but inaccurate statement of a similar nature. The distractor
- ⇒ statement should form a minimal pair with the original statement, i.e. the distractor
- → should be as similar to the original as possible while ensuring that the distractor is not
- → factual. The distractor should be crafted so that someone with only superficial
- → knowledge about the topic is likely to be fooled.

Let's see some examples before the real task.

Topic: Barack Obama

 Fact: Barack Obama was born in Hawaii. Distractor: Barack Obama was born in Kenya.

Topic: Wright brothers

Fact: Wright airplanes were involved in fatal crashes. Distractor: Wright airplanes were praised for their safety.

Topic: John Clempert

Fact: John Clempert was inspired by Houdini when developing acts.

Distractor: John Clempert was inspired by Penn and Teller when developing acts.

Now for the real task. Only output a distractor as in the examples.

Topic: {entity}
Fact: {claim}
Distractor:

P(True)

Your task is to determine whether the following claim related to {entity} is correct.

→ Only output "Yes" (correct) or "No" (incorrect).

Claim: {claim}

Yes or No:

Verbalized numerical confidence

The claim below was found in a passage about {entity}. State your confidence that the claim is correct. Only output an integer followed by "%".

Claim: {claim}

Prompt for the LLM to determine whether a sampled biography entails a claim for self-consistency

You will be given a passage and a claim. Your task is to determine whether the passage supports, refutes, or does not mention the claim. Output only "Support", "Refute", or

→ "No Mention".

918

919

920 921

922

923

924

925 926

927

928

929

930

931

932

933 934

935

936

937 938

939

940

941

942

943

944

945

946

947

948

949 950

951

952

953

954

955

956

957

958 959

960 961

962

963

Let's see some examples before the real task.

Passage: Barack Obama was the 44th President of the United States, serving from 2009 to

- → 2017. Born on August 4, 1961, in Honolulu, Hawaii, he was the first African
- → American to hold the office. Before his presidency, Obama served as a state senator in
- → Illinois and later as the 47th Governor of Illinois. A former constitutional law
- → professor, he was known for his eloquence, bipartisan approach, and focus on issues
- ⇒ such as healthcare reform, climate change, and foreign policy. His presidency was
- → marked by significant legislative achievements, including the Affordable Care Act,
- → and a commitment to diplomacy and international cooperation. After leaving office, he
- → authored memoirs and remained active in public life, advocating for social justice and

Claim: Barack Obama was born in Hawaii.

Relationship: Support

Passage: Tiger Woods is one of the most iconic and accomplished golfers in history, known

- → for his extraordinary talent, dominance on the course, and global influence on the
- → sport. Born on December 30, 1975, in Cypress, Florida, Woods rose to fame in the
- → mid-1990s and quickly became a household name, winning his first major
- → 15 major titles, the most in PGA Tour history, and has consistently ranked among the
- → world's top golfers for over two decades. His aggressive playing style, precision, and
- → mental toughness set him apart, making him a symbol of excellence in golf. Despite
- → personal challenges and setbacks, Woods has remained a dominant force in the sport,
- → inspiring millions of fans around the world.

Claim: Tiger Woods won a major championship at 19 years old.

Relationship: Refute

Passage: Albert Einstein was a theoretical physicist renowned for developing the theory of

- → relativity, which revolutionized the understanding of space, time, and gravity. Born in
- → 1879 in Ulm, Germany, he later moved to Switzerland and eventually to the United
- → States. Einstein's work, including the famous equation E=mc², laid the foundation for
- → modern physics and contributed to the development of nuclear energy. Despite his
- ⇒ scientific achievements, he was also a passionate advocate for peace, civil rights, and
- education. His legacy endures as one of the most influential scientists in history.

Claim: Albert Einstein became a US citizen.

Relationship: No Mention

Now for the real task.

Passage: {sampled biography}

Claim: {claim} Relationship:

B.2.3 PSEUDO-BEAM SEARCH

Prompt for the LLM to complete the prefix of an answer. Used for pseudo-beam search (Appendix C.1).

You will be given a prompt along with a prefix to begin your answer with. Your answer

→ should start with the given prefix. If the prefix itself is your final answer, you can

→ simply output just the prefix.

Let's look at 2 examples before the real task.

Example Prompt: Which American-born Sinclair won the Nobel Prize for Literature in

 \hookrightarrow 1930?

Example Answer Prefix: Sin Example Answer: Sinclair Lewis

Example Prompt: Where in England was Dame Judi Dench born?

Example Answer Prefix: York Example Answer: York

972

973 974

975

976

977

978

979 980

981 982

983

984

985

986

987

988

990 991 992

993

994

995

996

997 998 999

1000 1001

1002

1003

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017 1018

1019 1020

1021

1023 1024

1025

Now, here is a new prompt to answer. Answer with a concise phrase starting with the given \hookrightarrow prefix, as in the examples.

Prompt: {question}
Prefix: {prefix}
Answer:

B.3 LLM-AS-A-JUDGE

We score model responses for short-form QA (TriviaQA and SimpleQA) using an LLM judge rather than lexical matching for robust evaluation (Wei et al., 2024). Our biography generation task is evaluated with FActScore (Min et al., 2023), which uses a strong LLM for atomic claim decomposition and verification; we use GPT-4.1. We use Llama-3.1-8B-Instruct (Grattafiori et al., 2024) as the judge on TriviaQA, and GPT-4.1 as the judge on SimpleQA. On a sample of 100 questions from each dataset, we compared LLM judgments with human judgments. We used responses from Qwen3-8B on TriviaQA, and responses from GPT-4.1 on SimpleQA. The rate of agreement was 96/100 and 99/100, respectively. Upon reviewing disagreements, we found that in most cases there was genuine room for interpretation. For example, for the question "In The Living Daylights what did Carla keep in her cello case?", the target answer is "A machine gun" while the model answer was "A gun". It is unclear whether the model answer has the desired specificity. As another example, for the question "In Charles Dickens' "Great Expectations", who or what was Abel Magwitch?", the target answer is "Convict" while the model answer was "A convict and the main benefactor of Pip". The model answer contains the answer but contains more information, and without external information, it is impossible to determine whether the model's answer is correct. Overall, given the high agreement on unambiguously gradable questions, we deem it safe to adopt LLM-as-a-judge as a reliable evaluator in our experiments.

B.4 FACTSCORE EXAMPLE

Table 6 presents an example of atomic claim decomposition and verification with FActScore.

B.5 NLI MODEL

We use the NLI model DeBERTa-v3-base-mnli-fever-anli (He et al., 2021).⁷

⁷https://huggingface.co/MoritzLaurer/DeBERTa-v3-base-mnli-fever-anli

Table 6: Example of FActScore atomic claim decomposition and verification.

Generation

Kang Ji-hwan is a renowned South Korean actor and singer, best known for his role as the lead vocalist of the popular K-pop group BE:FIRST. Born on April 15, 2001, in Seoul, South Korea, Kang began his career in the entertainment industry at a young age, show-casing his talent through various music projects and performances. His distinctive voice and charismatic stage presence quickly earned him a loyal fanbase. Beyond his work in music, Kang has also ventured into acting, appearing in television dramas and variety shows, further solidifying his status as a multifaceted entertainer. With his dedication and natural talent, Kang Ji-hwan continues to make a significant impact in the K-pop and entertainment world.

Extracted claim	Correct?
Kang Ji-hwan is a South Korean actor.	Yes
Kang Ji-hwan is a South Korean singer.	No
Kang Ji-hwan is renowned.	Yes
Kang Ji-hwan is best known for his role as the lead vocalist of BE:FIRST.	No
BE:FIRST is a K-pop group.	Yes
BE:FIRST is a popular group.	Yes
Kang was born on April 15, 2001.	No
Kang was born in Seoul, South Korea.	Yes
Kang began his career in the entertainment industry at a young age.	No
Kang has showcased his talent through various music projects.	No
Kang has showcased his talent through various performances.	Yes
He has a distinctive voice.	No
He has a charismatic stage presence.	No
His distinctive voice quickly earned him a loyal fanbase.	No
His charismatic stage presence quickly earned him a loyal fanbase.	No
He quickly earned a loyal fanbase.	Yes
Kang has worked in music.	No
Kang has ventured into acting.	Yes
Kang has appeared in television dramas.	Yes
Kang has appeared in variety shows.	Yes
Kang is a multifaceted entertainer.	Yes
Kang's status as a multifaceted entertainer has been further solidified.	Yes
Kang Ji-hwan is dedicated.	Yes
Kang Ji-hwan has natural talent.	No
Kang Ji-hwan continues to make a significant impact in the K-pop world.	No
Kang Ji-hwan continues to make a significant impact in the entertainment world.	No

B.6 SIGNIFICANCE TESTING

ECE. We subsample 10k subsets of size 0.9 times the original dataset, where sampling is done without replacement. We construct an upper one-sided confidence interval with confidence level 0.95 for the tested method's ECE minus the best method's ECE and check whether the interval contains 0.

BS. As the Brier score is simply the mean squared error between confidences and truth values, it is well behaved and amenable to bootstrapping. We sample 10k subsets with the same size as the original dataset, where sampling is done with replacement. We construct an upper one-sided confidence interval with confidence level 0.95 for the tested method's BS minus the best method's BS and check whether the interval contains 0.

AUC. As AUC is a U-statistic, we use a one-sided DeLong test (DeLong et al., 1988; Sun & Xu, 2014) with confidence level 0.95.

Table 7: TriviaQA results. Building on Table 1, we report results when replacing P(True) with verbalized numerical confidence. Text styling follows Table 1.

	Qwen3-8B			Qwen3-1.7B			Llama-3.2-3B-Instruct			Gemma-3-4B-IT		
Method	ECE ↓	$BS\downarrow$	AUC ↑	ECE ↓	$BS\downarrow$	AUC ↑	ECE ↓	$BS\downarrow$	AUC ↑	ECE ↓	$BS\downarrow$	AUC ↑
VC _{numerical} (Tian et al., 2023b)	0.310	0.304	0.791	0.364	0.342	0.697	0.226	0.226	0.773	0.311	0.321	0.731
VC _{P(True)} (Kadavath et al., 2022)	0.240	0.242	0.817	0.387	0.383	0.720	0.189	0.208	0.826	0.300	0.299	0.702
K-VC (Tian et al., 2023a)	0.341	0.348	0.604	0.538	0.524	0.596	0.146	0.228	0.678	0.254	0.262	0.786
MSP (Fadeeva et al., 2023)	0.149	0.203	0.819	0.104	0.186	0.774	0.243	0.253	0.764	0.252	0.268	0.790
SC-VC (Xiong et al., 2024)	0.299	0.325	0.704	0.451	0.474	0.559	0.122	0.211	0.761	0.362	0.378	0.653
SC (Xiong et al., 2024)	0.236	0.244	0.785	0.233	0.229	0.780	0.065	0.177	0.808	0.303	0.304	0.713
NVC _{numerical}	0.232	0.225	0.870	0.069	0.162	0.812	0.183	0.202	0.825	0.267	0.272	0.787
DINCO _{numerical}	0.117	0.171	0.857	0.170	0.183	0.825	0.037	<u>0.150</u>	0.857	0.101	0.202	0.791
$NVC_{P(True)}$	0.171	0.190	0.853	0.084	0.164	0.806	0.168	0.192	0.845	0.218	0.236	0.791
$DINCO_{P(True)}$	0.097	0.155	0.879	0.177	0.179	0.835	0.044	0.148	0.864	0.121	0.191	0.817

C METHODS

C.1 PSEUDO-BEAM SEARCH

DINCO uses beam search to generate distinct distractors. However, beam search may not be available for API-access models. Here, we describe an approximation of beam search with similar inference cost that can be implemented if top token probabilities are provided, as they are in many API-access models such as GPT and Gemini.

We first run one inference pass to generate the main answer, which also gives us the highest probability tokens (and their probabilities) at each token position in the main answer. We consider the set of sequences that are a prefix of the main answer followed by a top token that is not the token in the main answer. We sort the sequences by their probabilities (computed with the chain rule). For each of the highest probability sequences, we present it to the LLM as a prefix to complete as an answer to the question (See Appendix B.2.3 for the prompt). This procedure chooses distinct prefixes with relatively high probability, shaping their completions to be distinct as well as plausible answers.

C.2 Self-consistency

Self-consistency estimates confidence in a main answer by sampling many answers and counting the proportion of answers that match the main answer (Xiong et al., 2024; Wang et al., 2023). Given a short-form question, a main answer c_0 , and sampled answers c_1, \ldots, c_K , we compute the confidence of c_0 as

$$f^{SC}(c_0) = \frac{1}{K+1} \sum_{k=0}^{K} \mathbf{1} \{ c_0 = c_k \}.$$
 (9)

Following Kuhn et al. (2023), we use an NLI model to determine semantic equivalence.⁴ For the long-form setting, we follow Zhang et al. (2024) and decompose a long-form response r_0 into claims c_1, \ldots, c_n and measure entailment from other sampled responses r_1, \ldots, r_K . We compute the confidence $f^{SC}(c_i)$ of a claim c_i by Eq. 9 with the summand replaced with $P(\text{entail} \mid r_k, c_i)$. Since the NLI task here involves a long-form text, we use the original LLM instead of an NLI model (see Appendix B.2 for prompts).

D RESULTS

D.1 DINCO WITH VERBALIZED NUMERICAL CONFIDENCE

In Tables 7 and 8, we expand on Tables 1 and 3 and report results on TriviaQA and biography generation when replacing P(True) with verbalized numerical confidence (see Appendix B.2 for prompts). For SimpleQA, the black-box variants of our methods in Table 2 use verbalized numerical confidence rather than P(True).

Table 8: FActScore results. Building on Table 3, we report results when replacing P(True) with verbalized numerical confidence. Text styling follows Table 3.

	Qwen3-8B					Gemma-3-4B-IT					
Method	ECE ↓	BS ↓	AUC ↑	$r \uparrow$	ρ \uparrow	ECE ↓	BS ↓	AUC ↑	$r \uparrow$	ρ \uparrow	
VC _{numerical} (Tian et al., 2023b)	0.327	0.328	0.749	0.508	0.565	0.482	0.465	0.648	0.186	0.169	
VC _{P(True)} (Kadavath et al., 2022)	0.433	0.431	0.625	0.073	0.122	0.527	0.527	0.683	-0.081	-0.129	
SC (Zhang et al., 2024)	0.162	0.226	0.771	0.468	0.494	0.197	0.233	0.787	0.629	0.607	
NVC _{numerical}	0.181	0.259	0.693	0.555	0.566	0.090	0.208	0.744	0.683	0.698	
DINCO _{numerical}	0.045	0.193	0.779	0.559	0.574	0.120	0.188	0.808	0.713	0.696	
$NVC_{P(True)}$	0.191	0.263	0.681	0.444	0.443	0.123	0.230	0.726	0.695	0.704	
$DINCO_{P(True)}$	0.076	0.202	0.767	0.518	0.538	0.172	0.210	0.793	0.724	0.712	

We find that the calibration benefits of DINCo generalize to the format of verbalized numerical confidence. For example, DINCo_{numerical} compared to $VC_{numerical}$ lowers ECE by 0.196 on average on TriviaQA. DINCo_{numerical} outperforms the stronger baselines as well, e.g. a lower ECE than SC by 0.103 on average on TriviaQA.

In the long-form setting of biography generation (Table 8), DINCO_{numerical} continues to achieve strong calibration, even outperforming DINCO $_{P(\mathrm{True})}$ in many cases, e.g. 0.045 vs. 0.076 ECE with Qwen3-8B.