

# Confidence Estimation for LLMs in Multi-turn Interactions

Caiqi Zhang<sup>1\*</sup>, Ruihan Yang<sup>2\*</sup>, Xiaochen Zhu<sup>1</sup>, Chengzu Li<sup>1</sup>, Tiancheng Hu<sup>1</sup>,  
Yijiang Dong<sup>1</sup>, Deqing Yang<sup>2</sup>, Nigel Collier<sup>1</sup>

<sup>1</sup>University of Cambridge <sup>2</sup>Fudan University

{cz391, nhc30}@cam.ac.uk<sup>1</sup>, {rhyang17, yangdeqing}@fudan.edu.cn<sup>2</sup>,

## Abstract

While confidence estimation is a promising direction for mitigating hallucinations in Large Language Models (LLMs), current research dominantly focuses on single-turn settings. The dynamics of model confidence in multi-turn conversations, where context accumulates and ambiguity is progressively resolved, remain largely unexplored. Reliable confidence estimation in multi-turn settings is critical for many downstream applications, such as autonomous agents and human-in-the-loop systems. This work presents the first systematic study of confidence estimation in multi-turn interactions, establishing a formal evaluation framework grounded in two key desiderata: per-turn calibration and monotonicity of confidence as more information becomes available. To facilitate this, we introduce novel metrics, including a length-normalized Expected Calibration Error (*InfoECE*), and a new "Hinter-Guesser" paradigm for generating controlled evaluation datasets. Our experiments reveal that widely-used confidence techniques struggle with calibration and monotonicity in multi-turn dialogues. We propose P(SUFFICIENT), a logit-based probe that achieves comparatively better performance, although the task remains far from solved. Our work provides a foundational methodology for developing more reliable and trustworthy conversational agents.

## 1 Introduction

Large Language Models (LLMs) have shown remarkable capabilities in multi-turn dialogue, collaborating with users on complex tasks (Wang et al., 2023; Yi et al., 2024; Laban et al., 2025). Yet their tendency to "hallucinate" (*i.e.*, producing incorrect statements with high apparent certainty) remains a major obstacle for high-stakes use (Manakul et al., 2023; Zhang et al., 2024a; Shelmanov et al., 2025; Hu et al., 2025). Confidence estimation, which

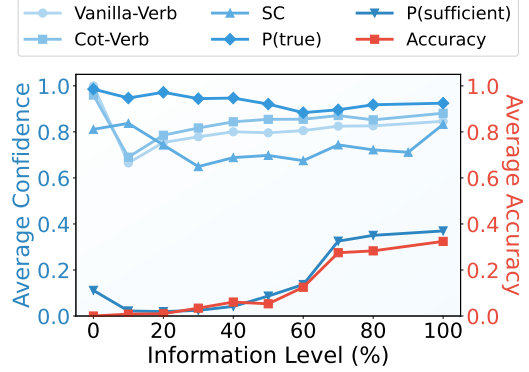


Figure 1: InfoECE on GUESS (Llama3.1-70B). Ideally, confidence (blue curves) increases as more information is provided. Calibration improves when the confidence curves (blue) are closer to the accuracy curve (red). In this setting, P(SUFFICIENT) best satisfies both monotonicity and calibration.

aims to predict the likelihood that a model’s answer is correct, has accordingly become a promising direction for identifying and mitigating such failures (Zhang et al., 2025b; Yang et al., 2025a,b).

Despite recent progress, most prior work studies confidence in single-turn question answering (Tian et al., 2023; Xiong et al., 2024), a static setup that overlooks the inherently **dynamic nature** of real human–AI interaction. In multi-turn conversations, information arrives incrementally: users refine their queries, models ask clarifying questions, and the hypothesis space narrows turn by turn. In such settings, confidence should not be a fixed attribute of a solitary response but a signal that evolves with the dialogue—ideally increasing as ambiguity is resolved and evidence accumulates. Reliable confidence estimation in this progression is therefore critical, as it serves as a decision-making heuristic for when to ask clarifying questions, invoke tools, or commit to actions in agentic workflows and human–AI collaboration. However, *how well current methods track this progression is largely unknown.*

To address this gap, we present the first system-

\*Equal contribution. Codes and data can be found in [GitHub](#).

atic study of confidence estimation in multi-turn conversations. We introduce a novel evaluation framework in which the model receives progressively more task-relevant information. We argue that in this controlled setting, a reliable confidence signal should satisfy two desiderata: (1) **Calibration**, where the confidence accurately reflects empirical correctness at any given turn, and (2) **Monotonicity**, where confidence consistently increases as more useful information becomes available.

Guided by these desiderata, we develop new metrics and datasets tailored to confidence estimation in multi-turn settings. To measure calibration across dialogues of varying lengths, we introduce a length-normalized Expected Calibration Error at information level (*InfoECE*). To quantify monotonicity, we employ Kendall’s  $\tau$ , a non-parametric rank correlation coefficient. We establish evaluation testbeds for two distinct regimes: (1) *under-specified initial queries*, for which we introduce a novel “Hinter–Guesser” paradigm to generate dialogues with progressively revealed clues; and (2) *difficult but fully-specified queries*, for which we adapt existing incremental QA benchmarks (Wallace et al., 2019; Sung et al., 2025) that provide sequential hints toward the correct answer.

Our experiments evaluate a suite of confidence estimation methods across four open-source models (§ 5.1), revealing key insights into their multi-turn behavior. **First** (§ 5.2), we find that widely used techniques struggle to maintain calibration or exhibit consistent monotonicity as conversations progress, as illustrated in Figure 1. Our proposed method, P(SUFFICIENT), proves comparatively more effective; however, substantial room for improvement remains. Meanwhile, we find that models exhibit stronger monotonicity when confidence is evaluated against the ground-truth answer rather than the model’s provisional answer at each turn. **Second** (§ 5.3), we examine whether confidence increases are driven by added information or merely by turn count. P(SUFFICIENT) more effectively distinguishes meaningful information gains from conversational filler. **Finally** (§ 5.4), our analysis reveals that while model accuracy is comparable between multi-turn dialogues and single-turn summaries, the confidence signals behave very differently, underscoring that the interactive structure of the dialogue is crucial to models’ confidence estimation. Overall, our findings highlight multi-turn confidence as a distinct and necessary target for reliable, decision-oriented LLM behavior.

## 2 Related Work

**Confidence Estimation in LLMs.** Confidence and uncertainty estimation has been extensively studied in LLMs (Geng et al., 2024). More specifically, uncertainty reflects the variability in the model’s predictions given *only the input query*, while confidence is defined with respect to *both the input and the specific generated output*, indicating how certain the model is about that particular response (Lin et al., 2023; Zhang et al., 2024b, 2025a). Mainstream confidence estimation approaches include prompting-based (verbalized) methods (Tian et al., 2023; Dong et al., 2024), consistency-based methods (Manakul et al., 2023; Zhang et al., 2024b), and logit-based methods (Kadavath et al., 2022). These methods have been applied to various tasks, such as short-form factual QA (Tian et al., 2023; Lin et al., 2023), long-form factual QA (Zhang et al., 2024b,c, 2025b), and reasoning tasks (Zhang et al., 2025a; Zhang and Zhang, 2025). However, a major limitation of existing works is their focus on single-turn settings. The effectiveness of confidence estimation in multi-turn conversations remains underexplored (Kirchhof et al., 2025), where model confidence may **evolve dynamically** throughout the interaction. Our work aims to fill this gap by systematically evaluating existing confidence estimation methods and proposing novel approaches in multi-turn contexts.

**LLMs in Multi-turn Interactions.** There has been growing interest in studying LLMs in multi-turn scenarios (Laban et al., 2025; Zhu et al., 2025). Modern LLMs support interactive dialogue, enabling users to collaborate with the model across multiple turns to accomplish complex tasks. However, recent studies show that LLMs often perform significantly worse on the same tasks when framed in a multi-turn context compared to a single-turn setting (Laban et al., 2025). Laban et al. (2025) also point out that many prior works (Bai et al., 2024; Kwan et al., 2024; Duan et al., 2024) simulate *episodic conversations*, where each turn introduces a subtask related to previous turns but can be evaluated in isolation. Under this framing, multi-turn tasks differ structurally from single-turn tasks and are not evaluated on the same set of questions. Laban et al. (2025) argue that episodic tasks tend to overestimate LLM performance in multi-turn settings and construct a *sharded* data construction method. In our work, we follow a similar sharded question construction strategy. For each question,

we create multiple variants with increasing levels of contextual information provided across turns. This allows us to directly compare confidence estimation methods under varying levels of complexity.

### 3 Methodology

#### 3.1 Notations

We study confidence estimation in *multi-turn* dialogue between a *user* and an *LLM*. Dialogs are indexed by  $d \in \{1, \dots, N\}$  and have  $L_d$  turns, where  $L_d$  may differ across dialogs. At turn  $i \in \{1, \dots, L_d\}$ , let the dialogue history be

$$h_{d,i} = \{q_{d,1}, a_{d,1}, \dots, q_{d,i-1}, a_{d,i-1}\}.$$

The turn- $i$  prompt consists of the task description  $T$  and the history  $h_{d,i}$ . Model  $M$  returns an answer  $\hat{y}_{d,i}$  and a confidence  $c_{d,i} \in [0, 1]$ . Each dialog has one gold label  $y_d$ , and we record correctness

$$z_{d,i} = \mathbb{I}[\hat{y}_{d,i} = y_d].$$

**Task characteristics.** We design a controlled task that should exhibit three key properties: *C1: Progressive Information Acquisition*. Each turn reveals additional task-relevant information that narrows the hypothesis space or supports step-by-step reasoning. *C2: Step-wise Answerability and Evaluation*. At every turn the model outputs an answer and a confidence, enabling per-turn accuracy and calibration assessment. *C3: Monotonic Confidence Progression*. Confidence should increase with the turn index and align more closely with true accuracy, providing a usable reliability signal. These properties yield a controllable testbed in which information strictly accumulates across turns, avoiding the limitations of episodic, non-progressive interactions.

#### 3.2 Two Initial-Question Regimes

Based on the task characteristics, we consider two regimes defined by the completeness of the initial question. **(1) Under-specified:** The initial question  $q_{d,1}$  admits many plausible answers. Hints progressively *prune* the candidate set. **(2) Fully-specified but difficult:** The initial question  $q_{d,1}$  pinpoints a unique answer in principle, but is hard to answer due to the models' knowledge limitation or reasoning ability. Hints make solving *easier*.

#### 3.3 Evaluation Protocol and Metrics

As a dialogue progresses, later turns include at least as much information as earlier ones. Therefore, a useful confidence signal should have at least:

a) *Per-level calibration:* within the same (normalized) information level, average confidence matches empirical accuracy.

b) *Monotonicity:* typically  $c_{d,i+1} \geq c_{d,i}$ ;

Because dialogue lengths vary, we first normalize turn  $i$  of dialogue  $d$  to a fractional information level

$$s_{d,i} = \frac{i}{L_d} \in (0, 1].$$

where  $L_d$  is the number of turns in dialogue  $d$ . We then partition  $[0, 1]$  into  $B$  bins  $\{S_b\}_{b=1}^B$  (either equal-width or equal-mass) and index all turn positions by  $\mathcal{I} = \{(d, i) : 1 \leq i \leq L_d\}$ . The subset of indices whose normalized level falls into bin  $b$  is  $\mathcal{I}_b = \{(d, i) \in \mathcal{I} : s_{d,i} \in S_b\}$ . For each information level  $b$ , the average confidence and accuracy are

$$\text{conf}_b = \frac{1}{|\mathcal{I}_b|} \sum_{(d,i) \in \mathcal{I}_b} c_{d,i}, \quad \text{acc}_b = \frac{1}{|\mathcal{I}_b|} \sum_{(d,i) \in \mathcal{I}_b} z_{d,i},$$

where  $c_{d,i} \in [0, 1]$  is the model's per-turn confidence and  $z_{d,i} \in \{0, 1\}$  indicates correctness of the turn- $i$  answer in dialogue  $d$ .

**Information-level ECE (InfoECE).** We compute an information-level ECE that groups predictions by normalized information exposure, enabling fair calibration comparisons across dialogues of different lengths.

$$\text{InfoECE} = \frac{1}{B} \sum_{b=1}^B |\text{acc}_b - \text{conf}_b|.$$

**Kendall's  $\tau$ .** Kendall's  $\tau$  measures the pairwise monotonic trend of confidence over turns. For dialog  $d$ , consider all  $\binom{L_d}{2}$  pairs  $(i < j)$ : a pair is *concordant* if  $c_{d,j} > c_{d,i}$  and *discordant* if  $c_{d,j} < c_{d,i}$  (ties ignored).

$$\tau^{(d)} = \frac{N_{\text{con}}^{(d)} - N_{\text{dis}}^{(d)}}{\binom{L_d}{2}}, \quad \bar{\tau} = \frac{1}{N} \sum_{d=1}^N \tau^{(d)}.$$

Values lie in  $[-1, 1]$ : 1 means strictly increasing confidences, and 0 no overall trend.

#### 3.4 Confidence Estimation Methods

Given the diversity of confidence estimation methods, we focus on the following three representative categories. We explicitly exclude post-hoc calibration techniques in this study (e.g., Platt scaling and temperature scaling), as they represent orthogonal research directions, which aim to rescale the confidence scores using statistical procedures (Zhou et al., 2025; Zhang et al., 2024c).

**Verbalized confidence.** We apply two verbalized prompting strategies (Tian et al., 2023) to elicit confidence scores directly from the model (see prompts in Appendix A):

- a) VANILLA-VERB: Given a candidate answer, the model is required to self-report confidence in  $[0, 100]$ ; rescale to  $[0, 1]$  for  $c_{d,i}$ .
- b) COT-VERB: Different from VANILLA-VERB, the model is now required to think step by step before given the self-reported confidence in  $[0, 100]$ ; rescale to  $[0, 1]$ .

**Self-consistency (SC).** Given the question, we independently sample  $m$  (e.g., 20) answers  $a_{d,i}^{(1)}, \dots, a_{d,i}^{(m)}$ . For any answer  $a$ , we define the confidence as the fraction of samples that match  $a$ :

$$c_{d,i} = \frac{1}{m} \sum_{j=1}^m \mathbb{I}[a_{d,i}^{(j)} = a].$$

**Logit-based probes.** We leverage internal model signals to estimate prediction confidence (see prompts in Appendix A). (1) P(TRUE) (Kadavath et al., 2022): At step  $i$ , given the prompt  $p_i$ , we first elicit the answer  $a_i$  from model  $M$ . We force a binary choice **A. True** vs. **B. False** with output constrained to a single uppercase letter. The confidence score is the model’s softmax probability assigned to label **A**.

Unlike P(TRUE), which asks if the answer is correct, we propose a new method that probes the confidence by asking model if the current information is **sufficient** (P(SUFFICIENT)) to entail that answer  $a$  is the only correct answer. P(SUFFICIENT) works particular well in our under-specified settings, where the set of plausible answers shrinks with each turn. This method allows the model to express low confidence even if its current guess happens to be correct, as long as other candidates have not yet been ruled out by the provided hints. This aligns the confidence score more closely with the true identifiability of the answer from the accumulated evidence, rather than mere incidental correctness.

Set  $c_{d,i} = P_T(d, i)$  or  $c_{d,i} = P_S(d, i)$ , we ask the model two binary probes about  $\hat{y}_{d,i}$  under  $p_{d,i}$ :

$$\begin{aligned} P_T(d, i) &= \Pr[A \mid p_{d,i}, \hat{y}_{d,i}; \text{P(TRUE)}], \\ P_S(d, i) &= \Pr[A \mid p_{d,i}, \hat{y}_{d,i}; \text{P(SUFFICIENT)}]. \end{aligned}$$

## 4 Dataset Construction

For *under-specified* (initially many plausible answers) regime, we construct our own datasets with 20Q and GUESS (as shown in Table 1). For *fully-specified* (a unique answer exists from the start but is hard to infer until sufficient evidence accumulates) regime, we directly apply existing datasets: GRACE (Sung et al., 2025) and TRICKME (Wallace et al., 2019) (examples in Table 7 in Appendix).

### 4.1 Under-specified Datasets

We primarily leverage 20Q and Guess-my-City (GUESS)-style (Abdulhai et al., 2023) settings for the under-specified regime. In both, an answerer holds a secret entity (an entity for 20Q; a city for GUESS). The questioner incrementally seeks information to recover the secret entity. The key difference is that 20Q constrains the questioner to yes/no questions, whereas GUESS permits open-ended questions. This incremental, information-seeking interaction naturally satisfies **C1**. Crucially, the setting also enables **C2**: at every turn the questioner can issue a concrete guess, even when information is still incomplete. This contrasts with math problem settings (e.g., Laban et al. (2025)), where intermediate turns often lack the conditions required to score correctness, impeding per-turn accuracy assessment. However, naively simulating two LLMs to play the roles of questioner and answerer can **violate C3**. Early turns may contain irrelevant or misleading questions, yielding stagnant or even decreasing confidence. To address this, we reformulate the interaction into a *Hinter-Guesser* paradigm that structures the information flow while retaining uncertainty.

**Hinter-Guesser Paradigm.** (1) **QA Stage.** A *Hinter* (LLM) is assigned a secret entity and must provide, at each turn, a helpful but non-trivial hint. A *Guesser* then (i) makes a best-guess answer and (ii) flags whether multiple answers remain plausible (*uniqueness probing*). (2) **Uniqueness Probing.** Even when the Guesser’s answer is correct, the Guesser indicates if other candidates still fit the evidence. This distinguishes a coincidentally right guess from the moment the answer becomes sufficiently supported by the clues, aligning confidence with identifiability rather than chance. (3) **Stopping & Filtering.** The dialogue proceeds until the Guesser both answers correctly and certifies uniqueness. We retain only successful dialogues (eventually solvable) and discard trajectories that



Dataset	Exemplar Prompt at Turn 4
20Q	<p><b>User:</b> Given the following information, provide the title of the Wikipedia page that best answers the last question fragment. If unsure, provide your best guess. The answer should be concise. You have some clues about the answer:</p> <p><b>Assistant:</b> Is it human-made? <b>User:</b> Yes</p> <p><b>Assistant:</b> Is it typically found indoors? <b>User:</b> Yes</p> <p><b>Assistant:</b> Is it commonly encountered in living rooms? <b>User:</b> Yes</p> <p><b>Assistant:</b> Is it larger than a book? <b>User:</b> Yes</p> <p><b>User:</b> Now guess what it is:</p> <p>Keyword: <b>television</b></p>
GUESS	<p><b>User:</b> Given the following information, name the single CITY that best fits them. If unsure, provide your best guess. The answer should be concise. You have some clues about the answer:</p> <p><b>Assistant:</b> What continent is the city in? <b>User:</b> Asia</p> <p><b>Assistant:</b> Is the city coastal or inland? <b>User:</b> Inland</p> <p><b>Assistant:</b> What’s the climate like in the city? <b>User:</b> Tropical</p> <p><b>Assistant:</b> What region within the continent is the city located? <b>User:</b> Southeast Asia</p> <p><b>User:</b> Now guess what it is:</p> <p>Keyword: <b>Bogor, Indonesia</b></p>

Table 1: Examples from the 20Q and GUESS datasets. Ideally, the model’s confidence should increase monotonically throughout the conversation. Note that the conversation history is collected from our Hinter-Guesser pipeline and remains fixed for confidence evaluation (*i.e.*, no strategic information-gathering is involved during evaluation).

fail to converge. In total, we collected 1,848 dialogue turns from 20Q, spanning 226 entities. For GUESS, we collected 1,625 turns spanning 223 entities. Then during confidence evaluation, since the conversations are **fixed**, models **do not** perform any strategic decision-making; instead, they only make guesses and predict confidence scores.

## 4.2 Fully-specified Datasets

We select two incremental, quizbowl-style QA datasets: GRACE (Sung et al., 2025) and TRICKME (Wallace et al., 2019), where clues become increasingly specific and a unique gold answer exists from the outset. GRACE and TRICKME directly supporting C1–C3 as evidence strengthens. We follow their standard protocols without modification, report per-turn accuracy and confidence, and defer details to the Appendix B.

## 5 Experiments

### 5.1 Setup

**Models.** In our experiments, we use Llama3.1 Instruct (8B and 70B) (Meta, 2024), Qwen2.5 Instruct (8B and 72B) (Yang et al., 2024). Temperature is set to 1 for sampling and otherwise 0.

**Confidence Estimation.** For a fair comparison across methods and models, we first let the model answer the question once to obtain an answer  $a$ . We compare  $a$  with the ground truth answer and label it correct or not. For each confidence estimation method, we then estimate the model confidence

in this answer  $a$ ; in parallel, we also estimate the model’s confidence at each turn with respect to the ground-truth answer.

**Controlling for Conversational Length.** A core hypothesis of our work is that a reliable confidence signal should increase monotonically as more task-relevant information becomes available. However, this trend could be a superficial artifact of dialogue length, where models become more confident simply because the turn index  $i$  is higher. To disentangle these factors, we design the following experiment: For a given dialogue  $d$  at turn  $i$ , we create an adversarial condition by replacing the original informative hint with a **placebo QA pair** that adds conversational history without revealing task-relevant information (*e.g.*, Q: “Is this a valid hint?” A: “Yes.”). We then compare the model accuracy and confidence across three states:

1. **Baseline (turn  $i - 1$ ):** The model’s prediction and confidence given history  $h_{d,i-1}$ .
2. **Original (turn  $i$ ):** Prediction and confidence after processing the original informative hint from turn  $i$ .
3. **Placebo (turn  $i'$ ):** Prediction and confidence after processing  $h_{d,i-1}$  followed by the uninformative placebo hint.

If confidence methods are robustly tracking information, we expect a significant increase in accuracy and confidence from the baseline to the original state. In contrast, the transition from the baseline

	Method	20Q		GUESS		GRACE		TRICKME	
		InfoECE	$\tau$	InfoECE	$\tau$	InfoECE	$\tau$	InfoECE	$\tau$
LLama3.1-8b	Accuracy	24.95		14.52		35.73		41.91	
	VANILLA-VERB	67.82	-6.36	74.89	-6.58	51.00	52.21	59.26	54.55
	COT-VERB	63.75	46.97	70.28	37.84	45.21	43.25	87.70	48.63
	SC	18.05	36.73	38.14	9.43	10.57	52.40	18.97	55.37
	P(TRUE)	69.02	42.10	67.08	19.91	50.43	48.48	55.61	52.37
	P(SUFFICIENT)	41.08	38.57	35.17	68.51	23.77	53.94	33.74	58.34
Llama3.1-70b	Accuracy	33.87		18.58		48.27		53.75	
	VANILLA-VERB	59.63	17.60	65.52	16.92	39.06	47.13	47.47	44.49
	COT-VERB	58.39	34.49	70.16	18.24	96.04	61.30	80.97	57.27
	SC	32.99	28.98	56.88	2.59	15.91	41.36	19.90	38.26
	P(TRUE)	67.82	40.82	79.97	3.29	37.04	58.94	35.62	64.25
	P(SUFFICIENT)	13.05	48.43	5.27	81.51	11.52	66.86	23.16	71.38
Qwen2.5-7b	Accuracy	25.22		12.92		27.34		34.06	
	VANILLA-VERB	58.05	61.13	48.68	26.99	50.40	55.55	54.62	56.39
	COT-VERB	64.93	52.60	71.84	56.01	66.20	50.38	65.37	49.89
	SC	45.44	13.85	50.53	36.10	32.78	40.16	33.28	42.54
	P(TRUE)	46.68	47.16	37.86	22.04	35.15	44.04	39.45	51.00
	P(SUFFICIENT)	36.64	55.24	26.63	51.44	28.67	47.79	35.26	52.11
Qwen2.5-72b	Accuracy	32.36		16.12		47.49		53.88	
	VANILLA-VERB	47.92	67.97	67.04	52.81	43.18	72.59	41.62	71.32
	COT-VERB	51.43	72.33	64.63	79.00	46.49	73.04	43.50	70.35
	SC	45.69	28.90	68.93	12.52	32.38	49.17	36.50	48.52
	P(TRUE)	42.12	68.88	57.56	54.87	31.86	64.02	32.87	69.28
	P(SUFFICIENT)	45.86	66.81	28.32	83.76	32.93	66.04	32.41	71.24

Table 2: InfoECE and  $\tau$  across models and datasets. Numbers are in percentages and best results are **bolded**.

to the adversarial state should yield a negligible change, despite the additional turn.

**Multi-turn vs. Single-turn.** Laban et al. (2025) suggest that LLMs can “get lost” in multi-turn conversations, performing worse than when all information is presented in a single turn. We investigate whether this phenomenon holds in our progressive information-seeking setting. For each turn  $i$  in a dialogue  $d$ , we define two conditions:

1. **Multi-turn:** The model is prompted with the full dialogue history up to that point,  $h_{d,i}$ , which includes the sequence of hints  $\{q_{d,1}, a_{d,1}, \dots, q_{d,i-1}, a_{d,i-1}\}$  preceding the current query  $q_{d,i}$ .
2. **Single-turn:** We create a single prompt containing a concise summary  $S_{d,i}$  that synthesizes all information from the hints provided up to turn  $i$ .

We then compare accuracy  $z_{d,i}$  and confidence  $c_{d,i}$  under both conditions. If models perform significantly better in the single-turn condition, it would suggest a cognitive burden in integrating information incrementally. Conversely, comparable or superior performance in the multi-turn setting would indicate that our structured, progressive framework effectively guides the model.

## 5.2 How reliable are confidence estimation methods in multi-turn settings?

We assess reliability along two axes: 1) *per-level calibration* using InfoECE (from 0 to 1, lower is

better); 2) *monotonicity* of confidence over turns using Kendall’s  $\tau$  (from -1 to 1, higher is better). We report  $\tau$  both on the model’s **current** answer and on the **gold** answer at each turn (Table 2). Figure 2 visualize how average confidence and accuracy evolve as information accumulates throughout the question-answering process.

**Calibration is generally poor and sufficiency probes help the most.** Across all models, both verbalized-based confidence (VANILLA-VERB, COT-VERB) and logit-based P(TRUE) are poorly calibrated, with InfoECE values typically between 40 and 80. Self-consistency (SC) is usually the most calibrated on fully-specified incremental QA. In under-specified games, sufficiency probing can be strikingly better-calibrated: for Llama3.1-70B, InfoECE drops to 13.05 on 20Q and 5.27 on GUESS, while maintaining competitive performance on GRACE and TRICKME. Overall, SC is a strong default for calibration. When the answer space is pruned gradually, P(SUFFICIENT) is a more *efficient* alternative. It narrows the gap and sometimes surpasses SC.

**Monotonicity on the current answer: sufficiency is usually best.** Ideally, confidence should rise as clues accumulate. P(SUFFICIENT) most consistently follows this trend: *e.g.*,  $\tau = 83.76$  on GUESS with Qwen2.5-72B and  $\tau = 71.38$  on TRICKME with Llama3.1-70B. In contrast, SC often shows weak monotonicity in under-specified settings (even single digits on GUESS). There are

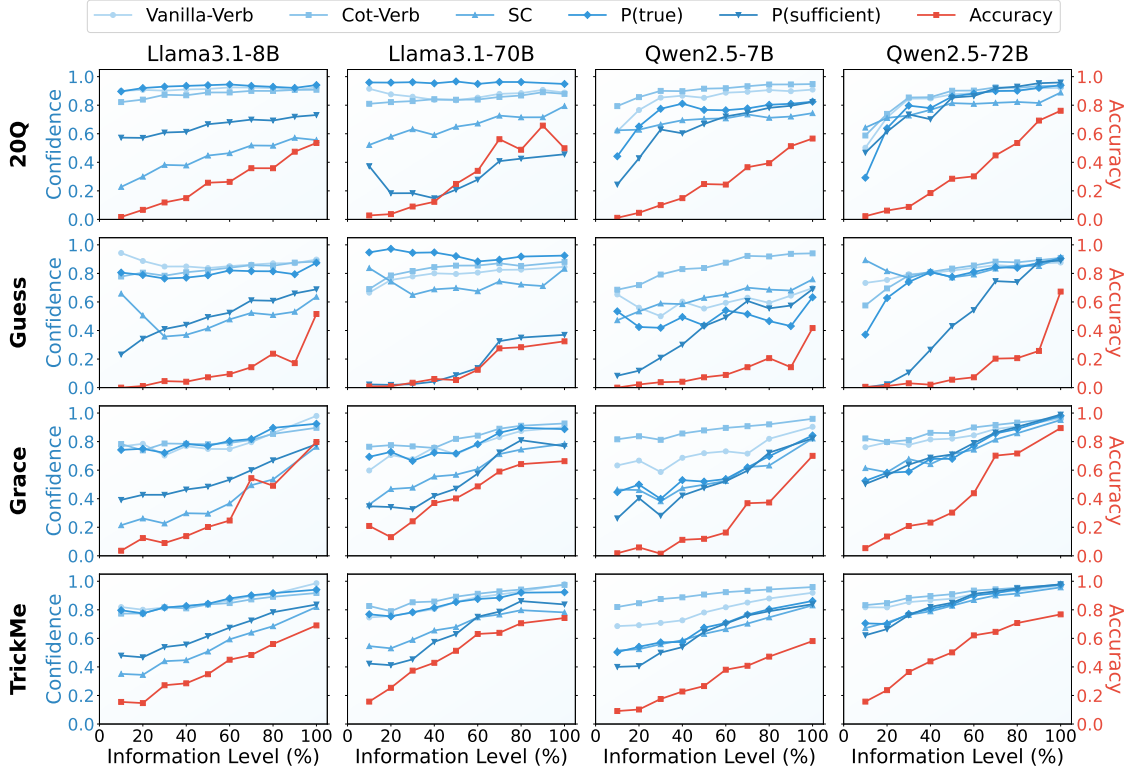


Figure 2: Evolution of average confidence and accuracy across different information levels. While accuracy (right y-axis, red line) generally increases, the confidence metrics (left y-axis, blue line) exhibit varying trends.

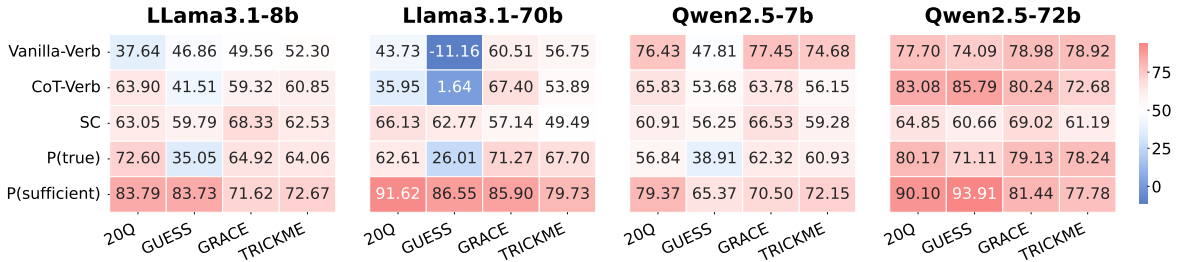


Figure 3: Kendall’s  $\tau$  for *ground truth* answers. Compared to the  $\tau$  for *each turn’s* answers, all methods show substantially better monotonicity. All values are shown as percentages.

model-family specific exceptions: with Qwen2.5 models, verbalized confidence (VANILLA-VERB or COT-VERB) occasionally attains the highest  $\tau$  scores on 20Q and GRACE, despite their generally poor calibration.

**Monotonicity on the ground truth: large gains and clear leaders.** As shown in Figure 3, when confidence is evaluated against the **ground truth** at each turn, all methods show substantial increases in  $\tau$ . Although the ground truth is unavailable in real-world applications, this trend suggests that models can partially recognize when current hints align with the correct answer. P(SUFFICIENT) dominates here, achieving  $\tau = 93.91$  on GUESS with Qwen2.5-72B, and  $\tau = 91.62, 86.55, 85.90$  on

20Q, GUESS, and GRACE with Llama3.1-70B, respectively. VANILLA-VERB can occasionally match or edge out on specific pairs (*e.g.*, Qwen2.5 on 20Q), but it remains poorly calibrated.

**Scaling and Model Family Effects.** As parameters increase, we observe a consistent rise in accuracy and a marked improvement in  $\tau$  (ranking calibration), particularly for the P(SUFFICIENT). For instance, **Qwen 2.5-72B** achieves the highest  $\tau$  score of 83.76% on the GUESS dataset, significantly outperforming its 7B counterpart. However, the effect on INFOECE is more nuanced; while larger models generally provide more reliable rankings, smaller models occasionally exhibit lower absolute calibration errors in specific configurations.

Method	20Q			GUESS		
	Conf <sub>i-1</sub>	Conf <sub>placebo,i</sub>	Conf <sub>i</sub>	Conf <sub>i-1</sub>	Conf <sub>placebo,i</sub>	Conf <sub>i</sub>
Llama3.1-8b	VANILLA-VERB	90.69	88.42 (2.27↓)	92.22 (1.53↑)	87.87	89.54 (1.67↑)
	COT-VERB	85.20	85.18 (0.02↓)	88.10 (2.90↑)	79.26	81.32 (2.06↑)
	SC	39.47	44.29 (4.82↑)	45.80 (6.33↑)	37.22	34.13 (3.09↓)
	P(TRUE)	91.51	89.74 (1.77↓)	94.08 (2.57↑)	73.63	85.38 (11.75↑)
	P(SUFFICIENT)	52.15	49.21 (2.94↓)	66.71 (14.56↑)	44.18	45.88 (1.70↑)
Llama3.1-70b	VANILLA-VERB	87.30	84.03 (3.27↓)	85.84 (1.46↓)	71.30	73.70 (2.40↑)
	COT-VERB	85.11	84.29 (0.82↓)	86.50 (1.39↑)	78.39	78.77 (0.38↑)
	SC	65.49	62.21 (3.28↓)	63.85 (1.64↓)	52.42	53.18 (0.76↑)
	P(TRUE)	95.48	93.43 (2.05↓)	94.56 (0.92↓)	88.16	88.14 (0.02↓)
	P(SUFFICIENT)	19.95	15.27 (4.68↓)	33.29 (13.34↑)	14.27	2.97 (11.30↓)
Qwen2.5-7b	VANILLA-VERB	86.31	80.97 (5.34↓)	86.97 (0.66↑)	68.86	67.80 (1.06↓)
	COT-VERB	89.36	85.13 (4.23↓)	92.21 (2.85↑)	81.48	82.96 (1.48↑)
	SC	66.73	72.52 (5.79↑)	69.69 (2.96↑)	70.49	65.74 (4.75↓)
	P(TRUE)	81.17	68.24 (12.93↓)	77.24 (3.93↓)	39.69	35.91 (3.78↓)
	P(SUFFICIENT)	46.66	39.44 (7.22↓)	67.58 (20.92↑)	41.38	20.04 (21.34↓)
Qwen2.5-72b	VANILLA-VERB	77.90	81.31 (3.41↑)	87.48 (9.58↑)	76.35	73.07 (3.28↓)
	COT-VERB	82.26	84.42 (2.16↑)	88.63 (6.37↑)	80.56	75.74 (4.82↓)
	SC	80.71	76.73 (3.98↓)	79.73 (0.98↓)	70.45	70.31 (0.14↓)
	P(TRUE)	80.65	81.75 (1.10↑)	85.01 (4.36↑)	51.98	66.59 (14.61↑)
	P(SUFFICIENT)	80.46	79.82 (0.64↓)	81.68 (1.22↑)	64.96	53.16 (11.80↓)

Table 3: Average confidence comparison across different models and datasets. Values in parentheses show the change relative to Conf<sub>i-1</sub>. Underlined values indicate statistically significant changes ( $p < 0.05$ ).

### 5.3 Does confidence track information or just turn count?

We test whether confidence increases are driven by accumulating information or are merely an artifact of conversational length. As detailed in our experimental setup and shown in Table 3, we compare a baseline confidence (turn  $i - 1$ ) with the confidence after receiving either an *informative* hint (Original) or a *non-informative* placebo hint. Across all 40 comparisons, informative turns yield more significant changes than placebos (27 vs. 18 with  $p < 0.05$ ). Among them, P(SUFFICIENT) most cleanly disentangles information gain from mere turn accumulation.

**Sufficiency probes actively penalize uninformative turns, tracking evidence over length.** The P(SUFFICIENT) method proves to be the most robust by actively penalizing uninformative turns. It frequently shows a statistically significant *decrease* in confidence after a placebo hint (*e.g.*, a drop from 14.27 to 2.97 for Llama3.1-70B on GUESS). This behavior, where the model lowers its sufficiency assessment in response to a useless hint, confirms it is tracking evidence, not just turn count. As a result, it achieves the clearest separation between conditions: confidence decreases or remains flat with a placebo, but increases with an informative hint. In contrast, other methods, particularly verbalized ones, can be misled into increasing confidence simply because the conversation has progressed.

### Placebo hints reveal important differences between methods.

The adversarial condition with placebo hints shows that models are not simply becoming more confident as a conversation gets longer. For many methods, the change in confidence after a placebo hint is not statistically significant (high  $p$ -value). For example, for Llama3.1-70B on GUESS, COT-VERB, SC, and P(TRUE) show negligible changes ( $p > 0.6$ ). This suggests a degree of robustness against superficial conversational structure.

### P(TRUE) is confounded by turn count (especially in GUESS).

In open-ended GUESS, P(TRUE) often *risks even under placebo*, consistent with a length artifact (mean  $\Delta_{\text{placebo}} = +5.64$ ; significant in 2/4 model pairs). Notably, Llama3.1-8B jumps +11.75 under placebo ( $p < 10^{-12}$ ) and Qwen2.5-72B jumps +14.61 ( $p < 10^{-6}$ ). Although P(TRUE) also increases with genuine information (mean  $\Delta_{\text{info}} = +14.07$  on GUESS), the placebo lift makes it less reliable for disentangling information from dialogue length. On 20Q, placebo effects tend to be negative (mean  $-3.91$ ), suggesting format-dependent behavior.

### Self-consistency (SC) shows moderate robustness.

SC typically exhibits small or negligible placebo movements and sizable gains with informative hints (mean  $\Delta_{\text{info}} = +9.83$  on GUESS). However, it is *not* immune: *e.g.*, Llama3.1-8B on 20Q



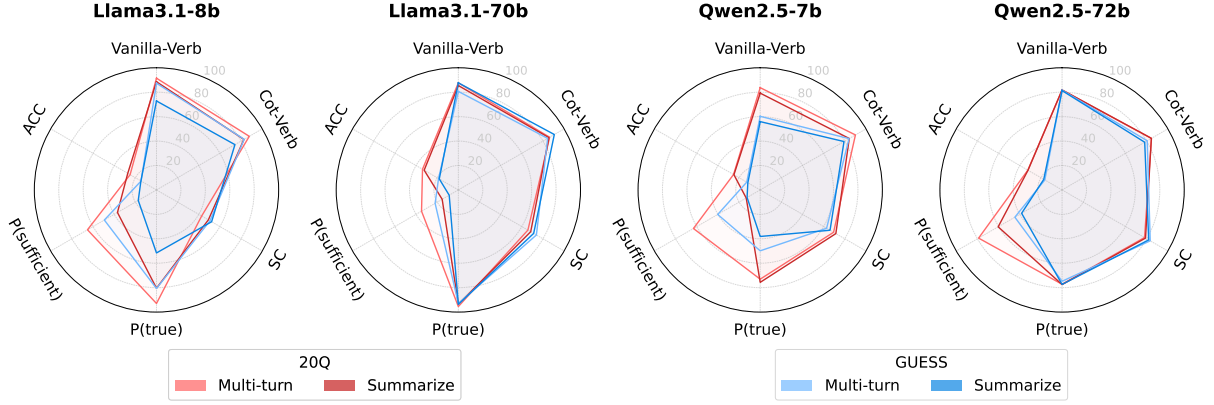


Figure 4: Performance comparison across four language models. Six evaluation dimensions are shown: Vanilla-Verb, Cot-Verb, SC, P(true), P(sufficient), and ACC. Red indicates 20Q benchmarks, blue indicates GUESS benchmarks.

increases under placebo by  $+4.82$  ( $p=0.0025$ ). Thus, SC generally tracks information better than verbalized scores but can still pick up turn-index artifacts in some settings.

**Verbalized confidence is unstable across conditions.** VANILLA-VERB/COT-VERB show small average placebo shifts (often non-significant) and modest informative gains; in some cases they even move counterintuitively (*e.g.*, Qwen2.5-7B on GUESS decreases by  $-9.29$  with an informative hint,  $p=0.005$ ). This instability, together with poor calibration (§ 5.2), limits their utility for turn-by-turn reliability.

#### 5.4 Single-Turn Summary vs. Multi-Turn Interaction

We compare model behavior when consuming clues incrementally (**multi-turn**) versus reading a concise synthesis of the same clues in one prompt (**single-turn summary**). Across models and datasets, accuracy differences between the two settings are small (mean absolute gap  $<1$ ), indicating no systematic advantage for either format (Figure 4). For example, on 20Q Llama3.1-8B slightly improves with summaries ( $24.95 \rightarrow 27.16$ ), whereas Llama3.1-70B slightly drops ( $33.87 \rightarrow 32.31$ ). On GUESS, Qwen2.5-72B gains modestly ( $16.31 \rightarrow 17.29$ ), while Qwen2.5-7B loses ( $12.74 \rightarrow 11.69$ ). In short, unlike the “getting lost” effect reported by Laban et al. (2025), our progressive information-seeking setup yields comparable task accuracy in multi-turn and single-turn conditions (see Appendix C for detailed InfoECE comparisons). One possible reason is that our tasks do not involve complicated arithmetic reasoning compared with Laban et al. (2025).

**Confidence shifts depend strongly on the signal.** While accuracy is stable, confidence responds markedly to prompt format. The sufficiency probe  $P(\text{SUFFICIENT})$  consistently *drops* under single-turn summaries (*e.g.*, on 20Q: Qwen2.5-7B  $63.13 \rightarrow 13.23$ ; Llama3.1-70B  $34.80 \rightarrow 15.30$ ), suggesting that compressing the dialogue into a synopsis removes turn-structure cues that the probe exploits to assess whether evidence is *enough*.  $P(\text{TRUE})$  and verbalized confidence often decrease with summaries for smaller models (*e.g.*, Llama3.1-8B on GUESS:  $P(\text{TRUE})$   $80.66 \rightarrow 51.58$ , VANILLA-VERB  $87.33 \rightarrow 72.82$ ), but can *increase* for larger models in some cases (*e.g.*, Llama3.1-70B on GUESS: VANILLA-VERB  $80.63 \rightarrow 87.65$ , COT-VERB  $84.43 \rightarrow 90.72$ ) without commensurate accuracy gains—an instance of potential miscalibration inflation. By contrast, SC is comparatively stable and sometimes *rises* with summaries on 20Q (*e.g.*, Llama3.1-8B:  $42.70 \rightarrow 49.04$ ), but shows mixed movement on GUESS.

## 6 Conclusion

We present the first systematic study of confidence estimation for LLMs in multi-turn conversations. We establish a formal evaluation framework grounded in two key desiderata with novel metrics and datasets. Our evaluation across various confidence estimation methods reveals that widely-used techniques struggle to maintain calibration and monotonicity in dynamic dialogues. We find that our proposed logit-based probe,  $P(\text{SUFFICIENT})$ , achieves comparatively better performance; however, the task remains significantly under-resolved. Building on our foundation, we advocate for future research into methods that: (1) satisfy both calibration and monotonicity; (2) effectively distinguish

task-relevant information from conversational filler; and (3) remain robust across both single-turn summaries and multi-turn interactions.

## Limitation

Our work, while providing a foundational framework, has several limitations that open avenues for future research. (1) The progressive datasets and Hinter–Guesser protocol simplify real conversations, omitting phenomena such as topic shifts, repairs, and mixed intents, which may limit transfer to messy, open-world dialogue. (2) Our study focuses on specific information-seeking tasks; the dynamics of confidence in more open-ended, creative, or collaborative conversations remain an open question. (3) Our evaluation emphasizes calibration and rank monotonicity; the downstream impact on user utility and human trust requires controlled user studies and field deployments. (4) We study *confidence* rather than *uncertainty*; extending our framework to uncertainty quantification and its relationship to confidence in multi-turn settings is an important next step.

## Ethics Statement

Our research follows standard ethical guidelines. We verified the licenses of all software and datasets used in this study. We introduce new multi-turn evaluation datasets built from existing resources and automated generation. Despite best efforts, the data may contain Western-centric biases, and we did not address multilingual coverage, which may limit generality across languages and cultures. Although some source datasets were created with human-in-the-loop protocols (*e.g.*, GRACE (Sung et al., 2025)), our experiments are fully automated. We recruited no new participants or annotators; thus no compensation or IRB oversight was required. We identified no privacy concerns, as we do not collect, store, or release personally identifiable information. We do not anticipate additional risks. We used an AI assistant only for grammar checking.

## References

Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu, and Sergey Levine. 2023. *Lmrl gym: Benchmarks for multi-turn reinforcement learning with language models*. Preprint, arXiv:2311.18232.

Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024. *MT-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7421–7454, Bangkok, Thailand. Association for Computational Linguistics.

Yijiang River Dong, Tiancheng Hu, and Nigel Collier. 2024. *Can LLM be a personalized judge?* In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 10126–10141, Miami, Florida, USA. Association for Computational Linguistics.

Haodong Duan, Jueqi Wei, Chonghua Wang, Hongwei Liu, Yixiao Fang, Songyang Zhang, Dahua Lin, and Kai Chen. 2024. *BotChat: Evaluating LLMs’ capabilities of having multi-turn dialogues*. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3184–3200, Mexico City, Mexico. Association for Computational Linguistics.

Jiahui Geng, Fengyu Cai, Yuxia Wang, Heinz Koepl, Preslav Nakov, and Iryna Gurevych. 2024. *A survey of confidence estimation and calibration in large language models*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6577–6595, Mexico City, Mexico. Association for Computational Linguistics.

Tiancheng Hu, Benjamin Minixhofer, and Nigel Collier. 2025. Navigating the alignment-calibration trade-off: A pareto-superior frontier via model merging. *arXiv preprint arXiv:2510.17426*.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, and 17 others. 2022. *Language models (mostly) know what they know*. Preprint, arXiv:2207.05221.

Michael Kirchhof, Gjergji Kasneci, and Enkelejd Kasneci. 2025. *Position: Uncertainty quantification needs reassessment for large-language model agents*. Preprint, arXiv:2505.22655.

Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. 2024. *MT-eval: A multi-turn capabilities evaluation benchmark for large language models*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 20153–20177, Miami, Florida, USA. Association for Computational Linguistics.

Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and Jennifer Neville. 2025. *Llms get lost in multi-turn conversation*. Preprint, arXiv:2505.06120.

- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. [Generating with confidence: Uncertainty quantification for black-box large language models](#). *Preprint*, arXiv:2305.19187.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*.
- Meta. 2024. [Llama 3 model card](#).
- Artem Shelmanov, Ekaterina Fadeeva, Akim Tsvigun, Ivan Tsvigun, Zhuohan Xie, Igor Kiselev, Nico Daeheim, Caiqi Zhang, Artem Vazhentsev, Mrinmaya Sachan, Preslav Nakov, and Timothy Baldwin. 2025. [A head to predict and a head to question: Pre-trained uncertainty quantification heads for hallucination detection in LLM outputs](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 35700–35719, Suzhou, China. Association for Computational Linguistics.
- Yoo Yeon Sung, Eve Fleisig, Yu Hou, Ishan Upadhyay, and Jordan Lee Boyd-Graber. 2025. [GRACE: A granular benchmark for evaluating model calibration against human calibration](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 19586–19587, Vienna, Austria. Association for Computational Linguistics.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Eric Wallace, Pedro Rodriguez, Shi Feng, Ikuya Yamada, and Jordan Boyd-Graber. 2019. [Trick me if you can: Human-in-the-loop generation of adversarial examples for question answering](#). *Transactions of the Association for Computational Linguistics*, 7:387–401.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023. [Mint: Evaluating llms in multi-turn interaction with tools and language feedback](#). *Preprint*, arXiv:2309.10691.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. [Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. [Qwen2 technical report](#).
- Ruihan Yang, Caiqi Zhang, Zhisong Zhang, Xinting Huang, Sen Yang, Nigel Collier, Dong Yu, and Deqing Yang. 2025a. [LoGU: Long-form generation with uncertainty expressions](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 18947–18968, Vienna, Austria. Association for Computational Linguistics.
- Ruihan Yang, Caiqi Zhang, Zhisong Zhang, Xinting Huang, Dong Yu, Nigel Collier, and Deqing Yang. 2025b. [UNCLE: Benchmarking uncertainty expressions in long-form generation](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 30328–30344, Suzhou, China. Association for Computational Linguistics.
- Zihao Yi, Jiarui Ouyang, Zhe Xu, Yuwen Liu, Tianhao Liao, Haohao Luo, and Ying Shen. 2024. [A survey on recent advances in llm-based multi-turn dialogue systems](#). *Preprint*, arXiv:2402.18013.
- Boxuan Zhang and Ruqi Zhang. 2025. [CoT-UQ: Improving response-wise uncertainty quantification in LLMs with chain-of-thought](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 26114–26133, Vienna, Austria. Association for Computational Linguistics.
- Caiqi Zhang, Zhijiang Guo, and Andreas Vlachos. 2024a. [Do we need language-specific fact-checking models? the case of Chinese](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1899–1914, Miami, Florida, USA. Association for Computational Linguistics.
- Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024b. [LUQ: Long-text uncertainty quantification for LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5244–5262, Miami, Florida, USA. Association for Computational Linguistics.
- Caiqi Zhang, Chang Shu, Ehsan Shareghi, and Nigel Collier. 2025a. [All roads lead to Rome: Graph-based confidence estimation for large language model reasoning](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 31802–31812, Suzhou, China. Association for Computational Linguistics.
- Caiqi Zhang, Ruihan Yang, Zhisong Zhang, Xinting Huang, Sen Yang, Dong Yu, and Nigel Collier. 2024c. [Atomic calibration of llms in long-form generations](#).
- Caiqi Zhang, Xiaochen Zhu, Chengzu Li, Nigel Collier, and Andreas Vlachos. 2025b. [Reinforcement learning for better verbalized confidence in long-form generation](#). *Preprint*, arXiv:2505.23912.
- Ej Zhou, Caiqi Zhang, Tiancheng Hu, Chengzu Li, Nigel Collier, Ivan Vulić, and Anna Korhonen. 2025. [Beyond the final layer: Intermediate representations for better multilingual calibration in large language models](#). *Preprint*, arXiv:2510.03136.

Xiaochen Zhu, Caiqi Zhang, Tom Stafford, Nigel Collier, and Andreas Vlachos. 2025. [Conformity in large language models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3854–3872, Vienna, Austria. Association for Computational Linguistics.



## Appendix

### A Instruction Prompt Examples.

I: GENERATION TEMPLATE
{Question} Use \boxed{ } to wrap your answer; for example \boxed{{a concise answer}}. Your answer is:
II: VANILLA-VERB TEMPLATE
{Question} Proposed Answer: {Answer} How confident are you that the proposed answer is correct? The confidence score should be a number from 0 (completely unsure) to 100 (absolutely certain). Your response MUST strictly adhere to this format: ### Confidence: \boxed{Your confidence score from 0-100.} Your response is:
III: COT-VERB TEMPLATE
{Question} Proposed Answer: {Answer} How confident are you that the proposed answer is correct? Analyze first and think step by step. The confidence score should be a number from 0 (completely unsure) to 100 (absolutely certain). Your response MUST strictly adhere to this format: ### Analysis: \boxed{Analyze the question and guess the answer. Think step by step and justify your confidence score.} ### Confidence: \boxed{Your confidence score from 0-100.} Your response is:

Table 4: Three TEMPLATE formats: generation and verbalize confidence estimation.

### B Incremental QA datasets

**GRACE** (Sung et al., 2025): a granular benchmark composed of quizbowl-style questions in which clues become increasingly specific. Each clue is self-contained and unambiguous. Models are evaluated on *how early*, *how accurately*, and *how confidently* they answer as clues unfold, providing fine-grained signals for step-wise calibration. This directly supports **C1** and **C2**, and the pyramidal clue structure encourages **C3** as evidence strengthens.

**TrickMe** (Wallace et al., 2019): a human-in-the-loop, adversarially-authored QA dataset built also in a quizbowl interface. Writers iteratively craft incremental clues to elicit confident model mistakes while remaining solvable by humans, yielding challenging, diverse questions that reveal miscalibration under partial information. It uses exactly the same task format as GRACE thus also satisfy our three criteria.

I: P(Sufficient) TEMPLATE
{Question} Based only on the information and hints provided above, does that information sufficiently entail that the correct answer is exactly {Answer}? A. Yes — the information is sufficient to conclude {Answer}. B. No — the information is insufficient, allows alternatives, or contradicts {Answer}. Output format: <b>**A**</b> or <b>**B**</b> only (single uppercase letter; no spaces, punctuation, or explanation):
II: P(Ture) TEMPLATE
{Question} Proposed Answer: {Answer} Is the proposed answer: A. True B. False Output format: <b>**A**</b> or <b>**B**</b> only (single uppercase letter; no spaces, punctuation, or explanation):
SUMMARIZE PROMPT TEMPLATE
You are given a block of hints about an entity (e.g., from a 20 Questions game). Summarize the hints into a concise statement or two, keeping the essential information instead of the Question/Answer format. Do not add new information or remove key attributes. Keep it as short and factual as possible.

Table 5: Two PROMPT\_TEMPLATE formats for binary classification tasks.

Table 6: SUMMARIZE\_PROMPT template for converting Q&A format hints into concise factual statements.

Together, GRACE and TrickMe instantiate the fully-specified setting where a unique gold answer exists from the outset, but models must calibrate *when* to commit as evidence accrues. We use both datasets without altering their underlying incremental-clue protocols, and we report per-turn accuracy and confidence to align with our evaluation framework.

### C Calibration shifts reveal a scaling-dependent format effect.

While accuracy remains comparable across formats (§5.4), calibration quality—measured by InfoECE—responds divergently by model scale (Figure 5). The sufficiency probe P(SUFFICIENT) exhibits opposite trends: for smaller models, summarization *degrades* calibration (e.g., Llama3.1-8B on 20Q: 6.99→24.57; on GUESS: 3.82→9.41), suggesting reliance on turn-by-turn structure. In stark contrast, larger models show substantial *improvements* under summarization (e.g., Llama3.1-70B

on 20Q: 40.29→9.81; on GUESS: 34.70→6.16; Qwen2.5-72B on GUESS: 27.51→2.55), indicating more effective integration of compressed evidence.  $P(\text{TRUE})$  improves markedly for Llama3.1-70B (20Q: 68.00→53.28; GUESS: 66.07→36.87) but shows minimal change for smaller models. SC and verbalized methods remain largely format-invariant (shifts typically  $<5$  InfoECE points) but consistently poorly calibrated ( $>50$  InfoECE). This scaling-dependent divergence suggests that while smaller models depend on conversational structure for reliable calibration, larger models can flexibly exploit either format, sometimes achieving superior calibration from summarization.

## **D Question Examples**

We list some examples of the four datasets we use in the study in Table 7.

## **E Placebo QA Examples**

We list the placebo QA examples we use for 20Q and GUESS datasets in Table 8.

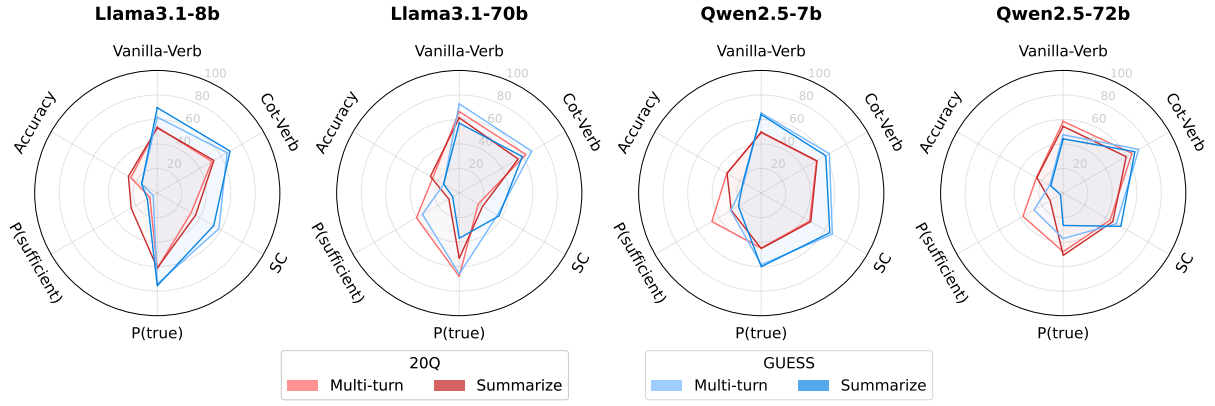


Figure 5: InfoECE of different confidence estimation methods across formats. Five confidence methods (Vanilla-Verb, CoT-Verb, SC, P(true), P(sufficient)) plus Accuracy, compared under multi-turn (darker) vs. summarized (lighter) presentation. Red: 20Q; blue: GUESS. Lower InfoECE indicates better calibration.

Dataset	Prompts
20Q	<p><b>User:</b> Given the following information, provide the title of the Wikipedia page that best answers the last question fragment. If unsure, provide your best guess. The answer should be concise.  You have some clues about the answer:  <b>Assistant:</b> Is it human-made? <b>User:</b> Yes  <b>Assistant:</b> Is it typically found indoors? <b>User:</b> Yes  <b>Assistant:</b> Is it commonly encountered in living rooms? <b>User:</b> Yes  <b>Assistant:</b> Is it larger than a book? <b>User:</b> Yes  <b>User:</b> Now guess what it is:  <b>Keyword:</b> <b>television</b></p>
GUESS	<p><b>User:</b> Given the following information, name the single CITY that best fits them. If unsure, provide your best guess. The answer should be concise.  You have some clues about the answer:  <b>Assistant:</b> What continent is the city in? <b>User:</b> Asia  <b>Assistant:</b> Is the city coastal or inland? <b>User:</b> Inland  <b>Assistant:</b> What's the climate like in the city? <b>User:</b> Tropical  <b>Assistant:</b> What region within the continent is the city located? <b>User:</b> Southeast Asia  <b>User:</b> Now guess what it is:  <b>Keyword:</b> <b>Bogor, Indonesia</b></p>
GRACE	<p><b>User:</b> Given the following information, provide the title of the Wikipedia page that best answers the last question fragment. If unsure, provide your best guess. The answer should be concise.  You have some clues about the answer:  <b>User:</b> It's not Charlie Parker, but a musician with this surname arranged excerpts from Stravinsky's The Firebird and the "Goin' Home" theme from the New World Symphony.  <b>User:</b> Another album by a musician with this surname features chants from the Bhagavad Gita and solos by Pharoah Sanders.  <b>User:</b> A song by a musician with this surname features Elvin Jones on timpani and gong.  <b>User:</b> A set of chord substitutions in a ii-V-I ("two-five-one") progression that proved challenging for pianist Tommy Flanagan are called this surname's "changes."  <b>User:</b> For 10 points, give this surname of harpist Alice and her husband, the saxophonist behind the album Giant Steps.  <b>User:</b> Now guess what it is:  <b>Keyword:</b> <b>Coltrane</b></p>
TRICKME	<p><b>User:</b> Given the following information, provide the title of the Wikipedia page that best answers the last question fragment. If unsure, provide your best guess. The answer should be concise.  You have some clues about the answer:  <b>User:</b> This man was seen driving Desiigner in the music video for Panda.  <b>User:</b> In an interview with Sway, this man yelled "I am Warhol," and compared himself with Shakespeare.  <b>User:</b> This man also controversially said that slavery was a choice.  <b>User:</b> For 10 points, name this songwriter known for his songs "Power," and "Gold Diggers," and more recently, "I Love It" with Lil Pump.  <b>User:</b> Now guess what it is:  <b>Keyword:</b> <b>Kanye_West</b></p>

Table 7: Examples from four datasets, showing question-answer dialogue format.

Guess My City	20Q
<p>Q: Does the city have people living in it? A: Yes</p> <p>Q: Does the city contain buildings? A: Yes</p> <p>Q: Are there roads in the city? A: Yes</p> <p>Q: Does the city have some form of waste disposal, like bins or trash collection? A: Yes</p> <p>Q: Is there access to toilets in the city? A: Yes</p> <p>Q: Does the city have shops or markets? A: Yes</p> <p>Q: Are there schools or educational institutions in the city? A: Yes</p> <p>Q: Does the city have hospitals or clinics? A: Yes</p> <p>Q: Is there some form of public transportation in the city? A: Yes</p> <p>Q: Does the city have restaurants or places to eat? A: Yes</p> <p>Q: Are there offices or workplaces in the city? A: Yes</p> <p>Q: Does the city have places for recreation, such as parks or sports areas? A: Yes</p> <p>Q: Is the city located on land? A: Yes</p> <p>Q: Does the city belong to a country? A: Yes</p> <p>Q: Is there some form of government or administration in the city? A: Yes</p> <p>Q: Does the city have streets or pathways for movement? A: Yes</p> <p>Q: Are there people who work in the city? A: Yes</p> <p>Q: Does the city have some form of shelter or housing? A: Yes</p> <p>Q: Is there electricity available in parts of the city? A: Yes</p> <p>Q: Does the city have some form of water supply or access? A: Yes</p> <p>Q: Are there vehicles that operate in the city? A: Yes</p> <p>Q: Does the city have some form of communication infrastructure? A: Yes</p> <p>Q: Are there businesses operating in the city? A: Yes</p> <p>Q: Does the city have some form of lighting at night? A: Yes</p> <p>Q: Are there emergency services available in the city? A: Yes</p> <p>Q: Does the city have banking or financial services? A: Yes</p> <p>Q: Are there entertainment venues in the city? A: Yes</p> <p>Q: Does the city have postal or delivery services? A: Yes</p> <p>Q: Are there religious or cultural institutions in the city? A: Yes</p> <p>Q: Does the city have some form of law enforcement? A: Yes</p>	<p>Q: Can it be described using words? A: Yes</p> <p>Q: Can people ask questions about it? A: Yes</p> <p>Q: Could it, in principle, be identified or referred to? A: Yes</p> <p>Q: Does it have at least one property? A: Yes</p> <p>Q: Is it what it is? A: Yes</p> <p>Q: Could someone think about it? A: Yes</p> <p>Q: Can it be distinguished from nothing at all? A: Yes</p> <p>Q: Is it possible to talk about it right now? A: Yes</p> <p>Q: Does it have some kind of name or label? A: Yes</p> <p>Q: Would it still count as something even if we know little about it? A: Yes</p> <p>Q: Could it, in theory, be observed or detected? A: Yes</p> <p>Q: Does it interact with its environment in some way? A: Yes</p> <p>Q: Could it be distinguished from absolutely nothing? A: Yes</p> <p>Q: Is it possible to classify it as something rather than nothing? A: Yes</p> <p>Q: Does it occupy some kind of location, even if unknown? A: Yes</p> <p>Q: Is it part of reality? A: Yes</p> <p>Q: Does it have some relation to other things? A: Yes</p> <p>Q: Could one imagine it being measured somehow? A: Yes</p> <p>Q: Is Earth around the Sun? A: Yes</p> <p>Q: Is the Moon larger than the Sun? A: No</p> <p>Q: Can numbers be even? A: Yes</p> <p>Q: Is blue a kind of sound? A: No</p> <p>Q: Can a thought have weight? A: No</p> <p>Q: Is time measured by clocks? A: Yes</p> <p>Q: Do triangles have three sides? A: Yes</p> <p>Q: Is water wetter than fire? A: Yes</p>

Table 8: We construct placebo question–answer pairs for control experiments using Guess My City (30 questions) and 20Q (26 questions). These questions neither provide any useful information to answer the question nor contradict with conversation history.