

---

# Tracing Uncertainty in Language Model “Reasoning”

---

Nils Grünefeld<sup>\*12</sup> Bertram Højer<sup>\*12</sup> Philipp Mondorf<sup>34</sup> Barbara Plank<sup>34</sup> Anna Rogers<sup>12</sup>  
Christian Hardmeier<sup>12</sup> Stefan Heinrich<sup>12</sup> Jes Frellsen<sup>52</sup>

## Abstract

Language model (LM) “reasoning”, commonly described as Chain-of-Thought or test-time scaling, often improves benchmark performance, but the dynamics underlying this process remain poorly understood. We study these dynamics through the lens of uncertainty quantification by treating the “reasoning” traces, the intermediate token sequences generated by LMs, as evolving model states. We summarize each trace by an *uncertainty trace profile*: a small set of features describing the shape of the uncertainty signal over its trace, such as its slope and linearity. We find that across five LMs evaluated on GSM8K and ProntoQA, these profiles predict whether a trace yields a correct final answer with AUROC up to 0.807, improving markedly on recent related work. We reach AUROC 0.801 using only the first few hundred tokens of full traces, suggesting that errors can be detected *early* in the generation. A detailed comparison of correct and incorrect traces further reveals qualitatively distinct uncertainty profiles, with correct traces showing a steeper and less linear decline in uncertainty. Together, the results suggest that our method, grounded in decision-making under uncertainty, provides a principled lens for studying the generative process underlying LM “reasoning”.

## 1. Introduction

The notion of *reasoning* is increasingly being used to describe the output generated by language models (LMs).

---

<sup>\*</sup>Equal contribution <sup>1</sup>Data Science Section, IT University of Copenhagen, Copenhagen, Denmark <sup>2</sup>Pioneer Centre for Artificial Intelligence, Copenhagen, Denmark <sup>3</sup>MaiNLP, Center for Information and Language Processing, LMU Munich, Munich, Germany <sup>4</sup>Munich Center for Machine Learning (MCML), Munich, Germany <sup>5</sup>Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kongens Lyngby, Denmark. Correspondence to: Nils Grünefeld <nilgr@itu.dk>, Bertram Højer <berh@itu.dk>.

Published as a paper at the 1<sup>st</sup> FoGen workshop, ICML 2026, Seoul, South Korea, 2026. Copyright 2026 by the author(s).

Whether referred to as a “scratchpad” (Nye et al., 2021), “Chain-of-Thought” (CoT) (Wei et al., 2022; Kojima et al., 2022), or test-time scaling (Muennighoff et al., 2025), researchers generally refer to the same process: an LM generating a long sequence of intermediate tokens that—ideally—results in a correct final answer. It is currently being debated whether “reasoning” is the right term to describe this process (Stechly et al., 2024; 2025; Kambhampati et al., 2025; Højer, 2025), but irrespective of that debate, open-ended long-form generation has been shown to improve performance on a variety of benchmarks with verifiable labels such as mathematical and logical reasoning (Yang et al., 2025; Guo et al., 2025). The connection between the length of the CoT and task performance is mainly correlational (Liu et al., 2025), although some work suggests that intervening on the “thinking” budget can directly affect performance (Muennighoff et al., 2025).

“Reasoning under uncertainty” is a well-established concept in the psychology of decision-making, which treats reasoning and decision-making as the process of forming judgments and selecting among alternatives, often given incomplete knowledge of outcomes and specific probabilities (Tversky & Kahneman, 1974; 1981; 1986). A definition of reasoning still eludes NLP and AI, and labeling the generative process of LMs as such arguably reflects misguided anthropomorphization. Nevertheless, the notion of reasoning under uncertainty is particularly relevant for this paper. Broadly speaking, definitions of *reasoning* in NLP and AI describe *some* process that is “logical and systematic” (Wei et al., 2022; Kojima et al., 2022; Huang & Chang, 2023), but even this vague “definition” is questionable. Generally, LMs do not enforce the logical validity of generated outputs, and the reinforcement learning (RL) optimization explicitly rewards only the generation of a correct answer, not the validity of the reasoning trace (Guo et al., 2025; Yang et al., 2025).

In this paper, we treat the generated CoT of an LM as an evolving state of the model; a state which can be more or less certain. Figure 1 illustrates this perspective: on the left we see the normalized average of the *distributional uncertainty* (entropy) computed at every step of generation. While correct and incorrect traces overlap and show a similar decrease in uncertainty as a trace progresses, the gap

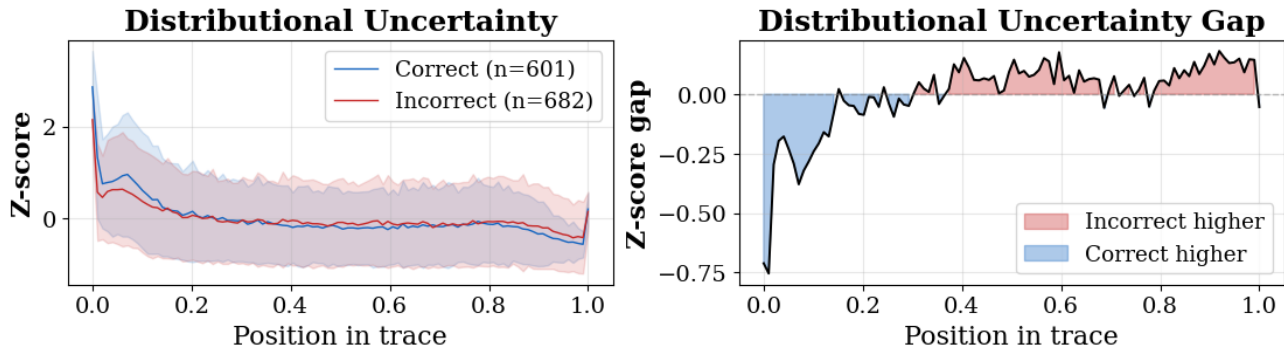


Figure 1. The average uncertainty trace for Qwen 2.5 on GSM8K. Left illustrates the average distributional uncertainty (entropy) computed over all generated sequences. Right shows the average uncertainty gap between correct and incorrect sequences.

plot (right) indicates sections of the averaged traces where correct and incorrect traces clearly deviate; it is this deviation in the uncertainty dynamics we utilize. We apply uncertainty quantification (UQ) to compute three measures of uncertainty to study how the model state develops as a time series, and show that the resulting **uncertainty trace profile** is strongly predictive of final answer correctness for mathematical and logical reasoning tasks.

**Related work.** A growing body of work has explored methods for improving LM performance at inference time via various interventions and uses of test-time compute, often by sampling multiple outputs and selecting among them. Wang et al. (2023) propose *Self-Consistency* which samples CoT reasoning paths, selecting the most frequent answer by majority voting. *Universal self-consistency* (Chen et al., 2023) extends this approach to open-ended tasks by prompting the model itself to identify the most consistent response. To address limitations related to scale, Kang et al. (2025) propose *Self-Certainty*, a metric based on the Kullback-Leibler divergence of each token’s predicted distribution from a uniform distribution, enabling best-of-N selection yielding better scaling as well as generalization to open-ended generation. Similarly, Fu et al. (2025) introduce *DeepConf*, a sampling method that uses model-internal confidence signals to filter low-quality reasoning traces during or after generation. Other approaches rely on hidden state representations to improve performance on verifiable tasks (Højer et al., 2025; Zhang et al., 2025) and Zhao et al. (2026) propose *Circuit-Based Reasoning Verification (CRV)* using transducers and the computation graph of Llama-3.1-8B to distinguish correct from incorrect generations.

**Novelty and contribution.** The aggregation inherent to other methods (*CRV*, *Self-Certainty*) discards most of the structure of the generative process; we treat per-token uncertainty as a time series, enabling richer characterization of how and where reasoning quality changes over the course of generation. Our work is thus complementary to said ap-

proaches: like *Self-Certainty* and *DeepConf*, we leverage uncertainty signals from model outputs and internals, but we are the first to utilize the trace as an evolving state as illustrated in Figure 1. Our primary research questions are thus: **RQ1)** Are uncertainty trace profiles predictive of correctness; and **RQ2)** What are the signals of uncertainty that distinguish incorrect from correct generations?

## 2. Background

### 2.1. “Reasoning” in Language Models

To formalize the mechanics of LM generation, we consider an LM as a parameterized function  $f_\theta : \mathcal{V}^{\leq K} \rightarrow \Delta(\mathcal{V})$ , where  $\mathcal{V}$  is the vocabulary,  $\mathcal{V}^{\leq K}$  denotes a sequence of up to  $K$  tokens, and  $\Delta(\mathcal{V})$  denotes the probability simplex over the vocabulary. The model parameters  $\theta$  define a conditional distribution  $p(v_t | v_{<t}, \theta)$  over the next token  $v_t \in \mathcal{V}$  given the prefix  $v_{<t} = (v_1, \dots, v_{t-1})$ , from which outputs are sampled during generation.

Generating from this model autoregressively is the foundation for producing long sequences of intermediate tokens that simulate reasoning. The finding that such intermediate generations improve performance on some verifiable tasks has resulted in the development of “reasoning models” (RMs) (Yang et al., 2025; Guo et al., 2025). In this paradigm a model is trained with RL to generate an extensive reasoning trace before generating the final answer. Recent discussions emphasize that RMs are post-trained with both a different style of data as well as learning objective, such as RL with verifiable feedback (RLVF). Some argue that models with this additional complexity are therefore necessarily not *just* “next-token predictors” (Downes et al., 2025), while others question how different RLVF is from standard optimization techniques (Shanahan, 2023; Samineni et al., 2025). Irrespective of this debate, the resulting models generate thousands of tokens and achieve improved performance on some verifiable tasks.

At surface-level, both LMs and RMs generate sequences of human-like reasoning. But while longer generated traces are correlated with improved performance on some benchmarks, the logical validity, reliability, and faithfulness of “reasoning” traces have been contested (Lanham et al., 2023; Turpin et al., 2023; Mirzadeh et al., 2025; Mondorf & Plank, 2024; Arcuschin et al., 2025; Chen et al., 2025). This calls into question approaches that analyze RM outputs through their semantic content alone (Ling et al., 2023), especially given evidence that training on semantically invalid CoT traces can still improve performance (Stechly et al., 2025).

## 2.2. Uncertainty in Machine Learning

Most machine learning is inherently probabilistic; even models that produce deterministic point predictions are implicitly embedded in a probabilistic framework (MacKay, 1992; Neal, 2012; Blundell et al., 2015). Within this framework, two kinds of uncertainty are commonly differentiated: aleatoric and epistemic (Hora, 1996; Der Kiureghian & Ditlevsen, 2009; Hüllermeier & Waegeman, 2021).

1. *Aleatoric* uncertainty reflects the inherent, irreducible difficulty of the task itself: variability that no amount of further training could remove.
2. *Epistemic* uncertainty reflects the degree to which a model is informed by its training data about a particular prediction; a decrease over the course of generation can thus be interpreted as “being closer” to regions of the output space encountered during training.

**Uncertainty quantification (UQ).** Quantifying these uncertainties is commonly approached using information-theoretic measures, specifically through an entropy decomposition of the predictive distribution  $H(y | x) = H(y | x, \theta) + I(y; \theta)$ , where  $x$  denotes the input,  $y$  the prediction,  $\theta$  the model parameters, and the two components on the right represent aleatoric and epistemic uncertainty, respectively (Smith & Gal, 2018).

Recently, this decomposition has been subject to criticism (Wimmer et al., 2023; Tomov et al., 2025) and a label-wise variance-based alternative has been proposed (Sale et al., 2023; 2024). This alternative leverages the law of total variance to decompose the total uncertainty for a given label  $k$ , defined as  $\text{Var}[Y_k]$ , into an epistemic component  $\text{Var}_\theta[p(y_k | x, \theta)]$ , capturing the variability of the predictive probability across different parameterizations, and an aleatoric component  $\mathbb{E}_\theta[p(y_k | x, \theta)(1 - p(y_k | x, \theta))]$ , capturing the expected inherent randomness of the prediction under a given parameterization. For efficient estimation of epistemic uncertainty in LMs, we adopt a recently proposed isotropic approximation (Grünefeld et al., 2026), which to our knowledge is the only computationally tractable ap-

proach for variance-based uncertainty quantification at this scale, particularly for the epistemic component.

In this work, we primarily use the variance-based framework for the epistemic component, since its information-theoretic analogue—the mutual information  $I(y; \theta)$ —has no tractable estimator for models on the scale of LMs. On the aleatoric side, we use both frameworks to further distinguish between two sub-types, committal and distributional:

- 1a. *Committal* aleatoric uncertainty,  $\mathbb{E}_\theta[p(\hat{y}_k | x, \theta)(1 - p(\hat{y}_k | x, \theta))]$ , the expected Bernoulli variance over  $\theta$  of the predicted class  $\hat{y}_k$ , taken from the variance-based decomposition.
- 1b. *Distributional* aleatoric uncertainty,  $\mathbb{E}_\theta[H(y | x, \theta)]$ , the expected conditional entropy of the predictive distribution over  $\theta$ , taken from the information-theoretic decomposition.

Committal uncertainty depends only on the probability  $p$  assigned to the selected token, and distinguishes two qualitatively different ways in which a token can be generated. When  $p$  is near 1, the model is committed to the specific token, in what we call a process of **concentration**; when  $p$  is near 0.5, committal uncertainty is at its maximum and the top token is selected not by being likely, but by being marginally less unlikely than its alternatives, in what we call a process of **elimination**. Distributional uncertainty, by contrast, is sensitive to the shape of the entire predictive distribution rather than just the probability of its mode. The two measures can therefore differ: a model may be firmly in the concentration regime, yielding low committal uncertainty, while still exhibiting high distributional uncertainty if the remaining mass is spread diffusely across many alternatives. We provide intuition for these uncertainty types and illustrative examples in Appendix A.

## 3. Methods

We analyze the development of uncertainty over the course of a trace as a novel approach to studying the generative process, in contrast to prior work which discards most of the temporal structure of autoregressive generation.

### 3.1. Uncertainty Estimation

Following the framework of subsection 2.2, we estimate the three uncertainty types at every token position. As described in subsection 2.1, the autoregressive generation is conditioned on the trace prefix  $v_{<t}$  and we let  $y$  be the target. The choice of target distinguishes two settings: for **trace uncertainty** the target is the next token  $v_t$ , while for **answer uncertainty** it is the eventual final answer generated by the model,  $\hat{y}$ .

The epistemic component is estimated using the squared L2 norm of the gradient of the predicted probability of the target with respect to all model parameters, derived via a first-order delta-method expansion of  $\text{Var}_\theta[p(y | v_{<t}, \theta)]$  around a point estimate  $\hat{\theta}$  of the parameters (Grünefeld et al., 2026):

$$\mathbb{U}_E := \text{Var}_\theta[p(y | v_{<t}, \theta)] \quad (1)$$

$$\approx \|\nabla_\theta p(y | v_{<t}, \theta) |_{\theta=\hat{\theta}}\|^2. \quad (2)$$

*Committal aleatoric uncertainty* is estimated as the Bernoulli variance of the target’s probability,

$$\mathbb{U}_C := \mathbb{E}_\theta[p(y | v_{<t}, \theta)(1 - p(y | v_{<t}, \theta))] \quad (3)$$

$$\approx p(y | v_{<t}, \hat{\theta}) \cdot (1 - p(y | v_{<t}, \hat{\theta})). \quad (4)$$

*Distributional aleatoric uncertainty* is computed as the predictive entropy over the full vocabulary, conditioned on the prefix and evaluated at the point estimate  $\hat{\theta}$ :

$$\mathbb{U}_D := \mathbb{E}_\theta[H(\cdot | v_{<t}, \theta)] \approx H(\cdot | v_{<t}, \hat{\theta}) \quad (5)$$

$$= - \sum_{v \in \mathcal{V}} p(v | v_{<t}, \hat{\theta}) \log p(v | v_{<t}, \hat{\theta}). \quad (6)$$

**Trace and answer uncertainty.** At each step  $t$ , trace uncertainty,  $\mathbb{U}^{\text{Tr}}$ , instantiates the three measures above with  $y = v_t$ , characterizing the state of a model with respect to the token it generates next. Answer uncertainty,  $\mathbb{U}^{\text{A}}$ , instantiates them with  $y = \hat{y}$ , characterizing instead the state of a model with respect to the eventual answer; this requires a second forward and backward pass at each step, conditioned on the trace prefix  $v_{<t}$  and applied to  $\hat{y}$  rather than  $v_t$  (see Appendix B) (Grünefeld et al., 2026). Since  $\hat{y}$  is in general a multi-token span, we use the sequence-level extensions of Grünefeld et al. (2026) for epistemic and committal uncertainties, and analogously the average entropy for distributional uncertainty.

We refer to  $\mathbb{U}_E$ ,  $\mathbb{U}_C$ , and  $\mathbb{U}_D$  as the three **types** of uncertainty, and to  $\mathbb{U}^{\text{Tr}}$  and  $\mathbb{U}^{\text{A}}$  as the two uncertainty **channels**. The two are orthogonal: each type has a trace variant and an answer variant, giving six time series per generation that together describe how local and answer-directed uncertainty co-evolve.

### 3.2. Experimental Procedure

**Models & datasets.** We evaluate Llama-3.1 (8B), Llama-3.2 (1B), Qwen-2.5 (0.5B), DeepSeek R1 Distill Qwen (1.5B), and Qwen-3 (0.6B). These models are chosen for comparability with related work (Zhao et al., 2026) and to represent different training techniques. Llama 3 models are trained primarily using supervised fine-tuning and limited

RL with human feedback and thus without explicit “reasoning” optimization (Grattafiori et al., 2024). Qwen-2.5 was trained similarly, but with a focus on “reasoning” data (Yang et al., 2024), whereas DeepSeek R1, and Qwen-3 were trained with the explicit objective of generating “reasoning” traces to improve correctness (Guo et al., 2025; Yang et al., 2025). These models range from standard LMs (the non-“reasoning” type) (Llama) to RMs (the “reasoning”-type) (R1, Qwen3), with Qwen-2.5 landing in-between. This enables an assessment of how uncertainty dynamics differ on this spectrum.

We evaluate on two datasets: GSM8K (Cobbe et al., 2021), chosen for comparability with related work (Zhao et al., 2026), and ProntoQA (Saparov & He, 2023) as a separate benchmark. GSM8K consists of grade school word math problems of varying difficulty while ProntoQA consists of small logical problems using low-likelihood tokens. GSM8K and ProntoQA are examples of benchmarks with verifiable answers on which RMs have shown strong performance improvements compared to LMs. We further study examples from GSM-Symbolic (Mirzadeh et al., 2025); a variant of GSM8K with controlled perturbations that preserve the underlying logic of questions.

**Generating traces.** We evaluate each model on both datasets using greedy decoding. The generated outputs are split into a “reasoning” trace and a final answer using model-specific splitting strategies that are documented in Appendix B. We initially set a max sequence length of 2048 for the trace and filter out responses that do not produce a final answer. For the reasoning models we append the `</think>` token to any trace that did not generate a final answer within the initial generation limit, restart the generation with an additional 512 token budget, and again filter out the traces that did not generate a final answer (details in Appendix D); the filtering removed approximately 10% of samples for DeepSeek R1 and Qwen 3.

### 3.3. Predictive Modeling & Reasoning Dynamics

**Feature engineering & correctness prediction.** We extract a small set of features from all type-channel combinations ( $\mathbb{U}_D^{\text{Tr}}$ ,  $\mathbb{U}_C^{\text{Tr}}$ ,  $\mathbb{U}_E^{\text{Tr}}$ ,  $\mathbb{U}_D^{\text{A}}$ ,  $\mathbb{U}_C^{\text{A}}$ ,  $\mathbb{U}_E^{\text{A}}$ ) to summarize the dynamics of a trace into the uncertainty trace profile. We compute the early mean ( $\mu_{\text{early}}$ ), the middle mean ( $\mu_{\text{mid}}$ ), the late mean ( $\mu_{\text{late}}$ ), the linear slope of the trace ( $\mathbf{m}$ ), and the fit of the slope ( $r^2$ ), resulting in 30 features in total (see Appendix C; Table 2). Using these features, we train two binary classifiers to predict the correctness of the final answer of every given trace, for each model-dataset combination. We choose a logistic regression (LogReg) classifier for its interpretability, and a gradient boosted model (GBoost) for comparability with related work (Zhao et al., 2026). Full experimental details are provided in Appendix D. We use

the area under the receiver operating characteristic curve (AUROC) to measure the predictive performance of each classifier in a five-fold cross-validation setup.

**Early correctness detection.** To assess how early we can predict correctness, we bin traces by number of tokens generated, and extract features based only on the tokens of the trace before the bin limit. We train the same three classifiers based on the constrained traces, to analyze how informative early trace dynamics are for predicting correctness.

## 4. Analysis

Using the uncertainty trace profile to predict correctness improves on the AUROC scores of *CRV* (Zhao et al., 2026) by 11.5% and *Self-Certainty* (Kang et al., 2025) by up to 59%. We report AUROC scores across models and datasets, and perform an in-depth analysis of the uncertainty trace profile for correct and incorrect traces (i.e. traces that result or don’t result in a correct answer).

### 4.1. Correctness Prediction

Table 1 displays the correctness prediction AUROC scores for each model-dataset combination. We report the scores achieved with our uncertainty trace profiles using both LogReg and GBoost classifiers as well as the *Self-Certainty* approach for comparison. Across all models, our uncertainty-based approach matches or outperforms *Self-Certainty*, and most of the predictive power is carried by  $\mathbb{U}^{\text{Tr}}$  (see Appendix E; Table 3). Notably, the uncertainty trace profile is predictive of correctness regardless of model type, in contrast to *Self-Certainty*. Further, we improve on the state-of-the-art *CRV* approach (Zhao et al., 2026), which evaluated a single model, Llama-3.1-8B on GSM8K, and reached an AUROC of 0.702 versus our 0.783.

**Uncertainty type classifiers.** We train classifiers for the trace and answer channels separately and observe similar, although slightly lower AUROC scores on the answer channel (see Appendix E; Table 3). We further train individual LogReg models for each type-channel combination, observing similar scores for  $\mathbb{U}_D$  and  $\mathbb{U}_C$  (AUROC up to 0.78) and slightly lower scores for  $\mathbb{U}_E$  (Appendix E; Table 4). This is likely due to the fact that  $\mathbb{U}_C$  and  $\mathbb{U}_D$  are similar: they both approximate aspects of aleatoric uncertainty and are thus closer to each other than to their epistemic counterpart.

**Early correctness detection.** Figure 2 shows that the UQ lens and our uncertainty trace profiles are informative of correctness relatively early in generation for trace-level uncertainty,  $\mathbb{U}^{\text{Tr}}$ . We observe AUROC up to 0.801 at the 300th token for GSM8K, and a similar trend across all five models, with AUROC rising steadily as larger parts of the full trace

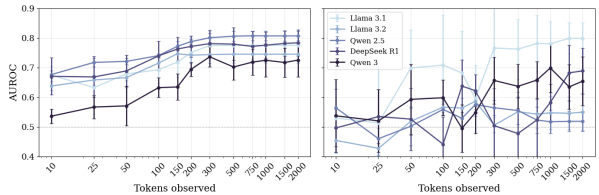


Figure 2. AUROC scores on models trained on features extracted based on increasing shares of the full traces (left: GSM8K, right: ProntoQA).

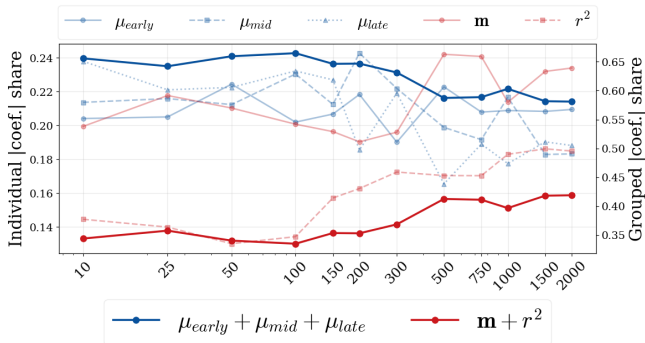


Figure 3. Change in coefficient share over the course of generation.

are included. A similar trend is visible for ProntoQA, although results are more noisy. While we primarily focus on the dynamics of uncertainty traces over the course of generation, this result indicates that our approach could be used to resample traces even before generation halts, improving on the efficacy of resampling approaches.

**Feature importance.** We measure feature importance via the coefficient share of the logistic regression, which is computed as the absolute coefficient of each feature as a fraction of the sum of all absolute coefficients. We aggregate static ( $\mu_{\text{early}}, \mu_{\text{mid}}, \mu_{\text{late}}$ ) and dynamic (slope,  $r^2$ ) features to assess the group-level importance, which is computed by summing the coefficient shares in each group. Figure 3 shows how the coefficient shares change over the course of a trace: at the very beginning of a generated trace the static uncertainty markers have the largest coefficient share, but as the generated sequences become longer, the slope and linearity coefficients of the uncertainty trace profile become more important in predicting correctness.

### 4.2. Feature Analysis

We examine the per-model uncertainty trace profiles shown in Figure 4, which displays features for incorrect and correct traces standardized by their z-score to focus on relative differences between correct and incorrect traces. We standardize within model because the absolute magnitudes of two of our measures are not informative for cross-model comparison. Gradient-based epistemic uncertainty,  $\mathbb{U}_E$ , is

## Tracing Uncertainty in Language Model “Reasoning”

	GSM 8K			ProntoQA		
	LR	GB	SC	LR	GB	SC
Llama 3.1	<b>0.783</b> $\pm 0.031$	0.758 $\pm 0.026$	0.491 $\pm 0.063$	<b>0.799</b> $\pm 0.052$	0.762 $\pm 0.101$	0.566 $\pm 0.080$
Llama 3.2	0.758 $\pm 0.017$	<b>0.767</b> $\pm 0.036$	0.566 $\pm 0.041$	<b>0.550</b> $\pm 0.061$	0.533 $\pm 0.073$	0.549 $\pm 0.065$
Qwen 2.5	<b>0.807</b> $\pm 0.017$	0.787 $\pm 0.016$	0.689 $\pm 0.028$	0.519 $\pm 0.034$	0.476 $\pm 0.050$	<b>0.565</b> $\pm 0.052$
DeepSeek R1	<b>0.786</b> $\pm 0.045$	0.775 $\pm 0.044$	0.703 $\pm 0.028$	<b>0.672</b> $\pm 0.081$	0.639 $\pm 0.106$	0.615 $\pm 0.083$
Qwen 3	0.727 $\pm 0.057$	0.665 $\pm 0.040$	<b>0.728</b> $\pm 0.036$	<b>0.657</b> $\pm 0.078$	0.551 $\pm 0.051$	0.611 $\pm 0.153$

Table 1. AUROC scores (5-fold CV) across five models and classifiers for GSM8K and ProntoQA. Values shown as mean  $\pm$  std. LR = Logistic Regression, GB = Gradient Boosting, SC = Self-Certainty.

not directly comparable across models and entropy,  $\mathbb{U}_D$ , is comparable in principle but dominated in scale by the long tail of effectively irrelevant tokens.

**Levels.** Across all five models,  $\mu_{\text{early}}$ ,  $\mu_{\text{mid}}$ , and  $\mu_{\text{late}}$  of most uncertainty types are higher for incorrect traces than for correct ones: the static level of uncertainty points to an incorrect answer as can be seen by comparing the rows of Figure 4. The gap is more pronounced for trace-level uncertainty,  $\mathbb{U}^{\text{Tr}}$ , than for answer-level uncertainty,  $\mathbb{U}^{\text{A}}$ , consistent with the finding that  $\mathbb{U}^{\text{Tr}}$  features carry most of the predictive power (Appendix E; Table 3). The clear exception is answer-level epistemic uncertainty,  $\mathbb{U}_E^{\text{A}}$ , where we observe the opposite, with incorrect traces showing lower uncertainty. The two epistemic measures therefore point in opposite directions on failure.  $\mathbb{U}_E^{\text{Tr}}$  is elevated on incorrect traces, indicating that the trace-level steps of an incorrect trace are less informed by training data than those of a correct trace;  $\mathbb{U}_E^{\text{A}}$  is depressed on incorrect traces, indicating that the eventual wrong answer is more informed by training data than correct answers tend to be. Our results can be interpreted as follows: when a model generates an incorrect answer, it generates less-supported steps to arrive at a wrong answer that is nonetheless closer to regions encountered during training. This pattern is complementary to the elevated aleatoric uncertainty measures observed on the same incorrect traces.

**Temporal dynamics.** Slope ( $\mathbf{m}$ ) and  $r^2$  characterize how the uncertainty trajectory develops over the course of generation as illustrated by the columns of Figure 4.  $\mathbf{m}$  distinguishes incorrect from correct traces consistently across models: in almost all cases, the slope of incorrect traces is higher relative to correct traces, indicating that uncertainty decreases less in incorrect traces as the generation progresses. The exception is  $\mathbb{U}_E^{\text{Tr}}$ : in Llama 3.1 the slope is still higher for incorrect traces, but the gap shrinks as we move from LMs and inverts in RMs, so that for DeepSeek-R1 and Qwen 3 the  $\mathbb{U}_E^{\text{Tr}}$  slope of incorrect traces is again lower.

The fit quality,  $r^2$ , follows a similar trend. In Llama 3.1, almost every uncertainty type has higher  $r^2$  for correct traces

than for incorrect ones, with correct trajectories following tighter linear trends and incorrect trajectories appearing more irregular. Moving from LMs to RMs this pattern flips such that incorrect traces appear more linear than correct ones.

In Llama 3.1, the local steps of incorrect traces are persistently less informed by training data than those of correct traces, with  $\mathbb{U}_E^{\text{Tr}}$  elevated throughout and not decreasing as the trace progresses. The eventual wrong answer is nonetheless more strongly supported by training data than correct answers tend to be. The trace and the answer it produces thus pull in opposite directions. In the RMs the picture changes:  $\mathbb{U}_E^{\text{Tr}}$  decreases more steeply on incorrect traces than on correct ones, so the local steps move progressively closer to better-supported territory as the model converges on the wrong answer that  $\mathbb{U}_E^{\text{A}}$  already flagged as well-supported. The smoother, more linear  $r^2$  profiles of incorrect RM traces are consistent with this convergence. Together, the level and temporal dynamics describe a qualitative shift in how incorrect traces fail across model classes.

### 4.3. Qualitative Analysis

To illustrate qualitatively how the analysis of uncertainty can be applied to study LM “reasoning”, we select a pair of sequences generated by Qwen 3 when evaluated on a GSM-Symbolic example (Mirzadeh et al., 2025). Recall that GSM-Symbolic examples are perturbed GSM8K questions that preserve the underlying logic. Figure 5 shows an example where the name and the size of the die differ between questions. We show the generated sequences (left) alongside the smoothed trace-level trajectories of distributional uncertainty  $\mathbb{U}_D^{\text{Tr}}$  and epistemic uncertainty  $\mathbb{U}_E^{\text{Tr}}$ , with the uncertainty trace profile features overlaid on the raw trace values (right): horizontal lines show the three means, dashed lines show the linear trend, and  $r^2$  captures the variation around the linear slope.  $\mathbb{U}_D^{\text{Tr}}$  is elevated on the incorrect trace across  $\mu_{\text{early}}$ ,  $\mu_{\text{mid}}$ , and  $\mu_{\text{late}}$  and the incorrect slopes fail to decrease as cleanly as the correct ones, consistent with the analysis of subsection 4.2.

The model correctly computes the relevant probabilities (50% for the first event, 25% for two consecutive odd rolls)

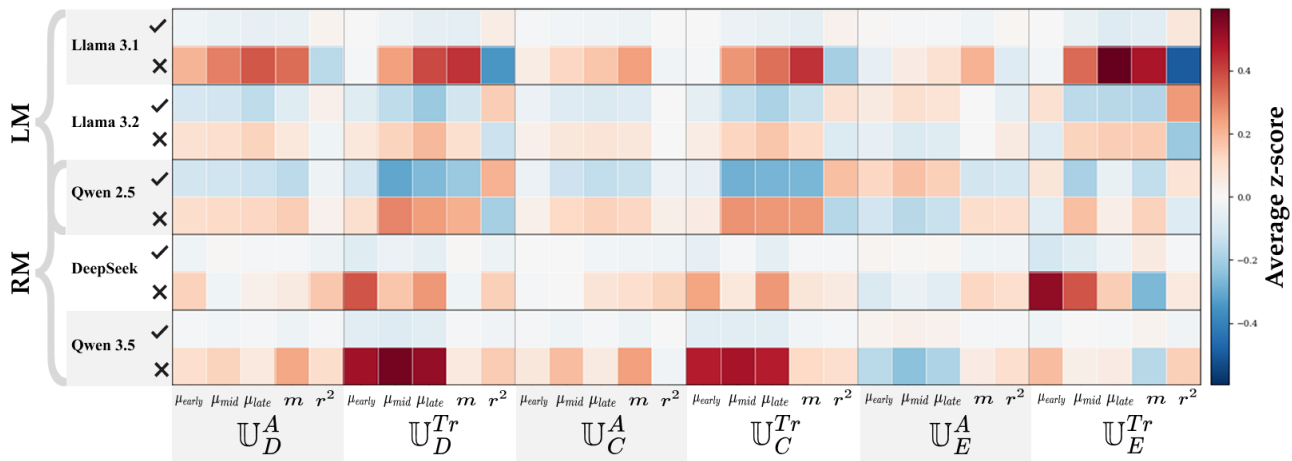


Figure 4. Heatmap of uncertainty features across models and datasets. The heatmap is organized into six feature blocks— $\mathbb{U}_D$ ,  $\mathbb{U}_C$ , and  $\mathbb{U}_E$  for  $\mathbb{U}^A$  and  $\mathbb{U}^{Tr}$  respectively—each containing five features:  $\mu_{early}$ ,  $\mu_{mid}$ ,  $\mu_{late}$ , slope ( $m$ ), and  $r^2$ . Colors indicate relative differences within model between correct and incorrect traces.

for both sequences and correctly considers the same two interpretations of “*how much more likely*”: the *difference* and the *ratio*. In the incorrect sequence, the model commits to using the ratio, shown in the highlighted sentence 31 of Figure 5 which corresponds to the highlighted token number 500; this results in the model incorrectly outputting 200% instead of 25% as the final answer. The correct trace peaks at the multiplicative-probability composition, while the incorrect trace shows a local elevation when committing to the ratio interpretation, with the signal reaching roughly  $2.4\times$  the average of the surrounding 200-token window. For the correct trace, there is no peak at the mention of the ratio, but at the point of committing to the difference. Importantly, focusing on the first semantically observable error in the generated reasoning aligns only with a single peak of the trace, at token 500. Many of the earlier peaks could be relevant for error-detection.

Both trajectories additionally show their largest relative peaks at structural boundaries—the trace opening and the post-think “*Final Answer*” section header—which we read as tokenization-driven gradient elevation rather than reasoning content; against this baseline the smaller “reasoning”-content peaks are the ones that distinguish the two traces and that flag the specific “reasoning” steps at which the incorrect trace shows errors.

## 5. Discussion

The UQ lens reveals a stark difference between generative processes resulting in correct and incorrect final answers, answering RQ1 in the affirmative: uncertainty trace profiles predict correctness with AUROC up to 0.807, outperforming *Self-Certainty* (Kang et al., 2025) and *Circuit-Based Reasoning Verification* (Zhao et al., 2026), and one can pre-

dict correctness with AUROC up to 0.801 using only the first 300 tokens.

UQ provides an interpretative lens through which to understand the process of autoregressive generation that is currently being labeled as “reasoning”. Correct and incorrect traces are distinguishable across the studied models, and the analysis shows that the demarcation is made based on a relative difference in the slope and linearity of the entropy,  $\mathbb{U}_D$ , and Bernoulli variance,  $\mathbb{U}_C$ , with a steeper and more linear slope being indicative of a trace being incorrect. Notably, the epistemic uncertainty presents a complementary behavior, being smaller for incorrect traces than for correct traces. Together, the aleatoric trajectory and the epistemic dissociation answer RQ2: the discriminating signals lie in the joint behavior of the two channels, not in either taken in isolation. Under the classical interpretation of aleatoric and epistemic uncertainty discussed in subsection 2.2, the most intuitive failure mode couples the two measures: an out-of-distribution input should yield elevated epistemic uncertainty, and a difficult input should yield elevated aleatoric uncertainty. Our findings deviate from this interpretation: incorrect traces show *elevated* aleatoric and *depressed* epistemic uncertainty. The model is uncertain during generation, yet the wrong answer it produces is seemingly well-supported by training data; aleatoric and epistemic thus carry separate diagnostic signals: neither alone reveals this failure mode.

**Implications.** First, the predictive power of uncertainty trace profiles offers a practical signal for early correctness estimation, in both LMs and RMs. This has potential applications in selective generation and test-time scaling that aim to terminate or resample traces prior to completion. Second, the observed dissociation between aleatoric and

(Xavier / Yusuf) rolled a (12 / 10)-sided die, how much more likely is it that he rolls a number greater than (6 / 5), expressed as a percentage, compared to rolling two odd numbers in a row?

<p>1 &lt;think&gt; Okay, so Xavier is rolling a 12-sided die.</p> <p>8 Therefore, the probability of rolling a number greater than 6 is 6/12, which simplifies to 1/2.</p> <p>So, the probability of rolling an odd number on the first roll is 1/2, and then the probability of rolling another odd number on the second roll is also 1/2.</p> <p>19 So, we need to find the difference in probabilities and express it as a percentage.</p> <p>38 Therefore, the ratio of A to B is (1/2)/(1/4) = 2, so 2 times more likely.</p> <p>48 Therefore, the answer is 25%.</p>	<p>1 &lt;think&gt; Okay, so Yusuf is rolling a d10 die, and I need to figure out how much more likely it is that he rolls a number greater than 5 compared to rolling two odd numbers in a row.</p> <p>10 So each roll has a <math>5/10 = 1/2</math> probability of being greater than 5.</p> <p>26 So, the probability would be <math>(1/2) * (1/2) = 1/4</math>.</p> <p>30 So, if the probability of greater than 5 is 50%, and the probability of two odd numbers in a row is 25%, then the difference is <math>50\% - 25\% = 25\%</math>.</p> <p>31 Therefore, the ratio is <math>50\%/25\% = 2</math>.</p> <p>148 Therefore, the probability of rolling a number greater than 5 is 200% more likely than rolling two odd numbers in a row.</p>
--	---

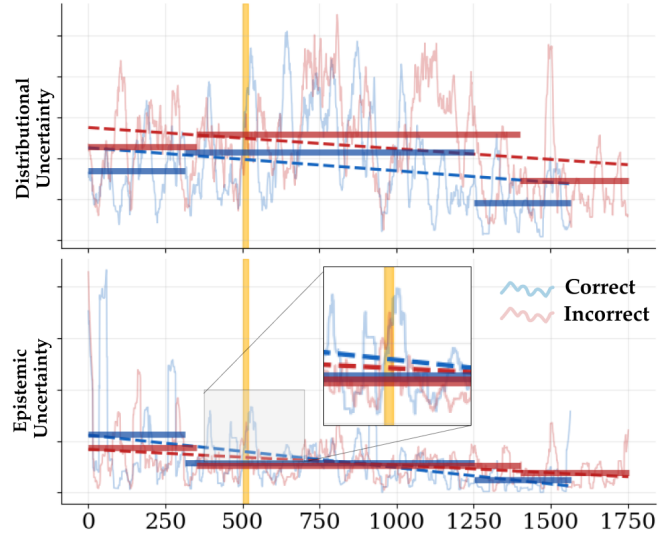


Figure 5. Sample from GSM-Symbolic. Numbers indicate sentence indices.  $\mu_{early}$ ,  $\mu_{mid}$ , and  $\mu_{late}$  are shown as horizontal lines and slope as a dashed linear-fit line, with the corresponding  $r^2$  reflected in the spread of the data around that line. Each plot shows the uncertainty type with the token index on the x-axis, and the location of the first error in the incorrect trace is highlighted in yellow.

epistemic measures of uncertainty points to a characterization of failure modes that simpler measures miss. Third, our case study illustrates failures that uncertainty-based analysis can identify but semantic analysis of the trace cannot, establishing the uncertainty lens as an analytical tool that yields complementary insights to semantic inspection.

**Limitations and future work.** We propose a method for analyzing “reasoning” based on temporal dynamics of the generative process, which have for the most part been ignored by prior work. Our method and analysis rely on interpretations of aleatoric and epistemic uncertainty estimates that are commonly used and well-founded, but nevertheless rest on approximations that are defensible in the overparameterized regime, yet not strictly verifiable. Further, the long sequences generated by RMs sometimes exceed the maximum length of 2048 tokens set for the experiments, potentially distorting results; we mitigate this through the filtering procedure described in subsection 3.2. The results differ in consistency between GSM8K, an older dataset which is likely to be in pre-training data, and ProntoQA, which is more recent and possibly more indicative of out-of-distribution “reasoning” due to the nature of the examples. We emphasize that this paper presents a proof-of-concept for the informativeness of uncertainty measures across “reasoning” traces, which we show on five models (0.5B-8B parameters) spanning the LM to RM spectrum and two commonly used datasets for mathematical and logical reasoning. Future work should expand the evaluation to other types of “reasoning” tasks, different model architectures and sizes, as

well as various out-of-distribution shifts, to further generalize our results.

## 6. Conclusion

We have proposed treating LM “reasoning” traces as time series of uncertainty estimates, summarizing each trace by an uncertainty trace profile that captures the level and dynamics of uncertainty over the course of generation. The uncertainty trace profile is strongly predictive of correctness across five models and two datasets, with AUROC up to 0.807 on full traces and 0.801 using only the first few hundred tokens, a significant improvement over prior approaches. We further show that correct and incorrect traces differ systematically under this lens, and that the epistemic and aleatoric channels carry diagnostic signal in opposite directions, indicating that the two cannot be summarized by a single measure. Taken together, these findings indicate that uncertainty quantification across “reasoning” traces captures substantial information about the dynamics of the process, and offers a lens for studying it that semantic inspection alone does not provide.

## References

Arcuschin, I., Janiak, J., Krzyzanowski, R., Rajamanoharan, S., Nanda, N., and Conmy, A. Chain-of-Thought Reasoning In The Wild Is Not Always Faithful, March 2025.

Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra,

- D. Weight Uncertainty in Neural Networks. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1613–1622. PMLR, June 2015.
- Chen, X., Aksitov, R., Alon, U., Ren, J., Xiao, K., Yin, P., Prakash, S., Sutton, C., Wang, X., and Zhou, D. Universal Self-Consistency for Large Language Model Generation, November 2023.
- Chen, Y., Benton, J., Radhakrishnan, A., Uesato, J., Denison, C., Schulman, J., Somani, A., Hase, P., Wagner, M., Roger, F., Mikulik, V., Bowman, S., Leike, J., Kaplan, J., and Perez, E. Reasoning Models Don’t Always Say What They Think, May 2025.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training Verifiers to Solve Math Word Problems, November 2021.
- Der Kiureghian, A. and Ditlevsen, O. Aleatory or epistemic? Does it matter? *Structural Safety*, 31(2):105–112, 2009. ISSN 0167-4730. doi: 10.1016/j.strusafe.2008.06.020. Risk Acceptance and Risk Communication.
- Downes, S. M., Forber, P., and Grzankowski, A. LLMs are not just next token predictors. *Inquiry: An Interdisciplinary Journal of Philosophy*, January 2025. doi: 10.1080/0020174X.2024.2446240. Published online 12 January 2025.
- Fu, Y., Wang, X., Tian, Y., and Zhao, J. Deep Think with Confidence, August 2025.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al. The Llama 3 Herd of Models, November 2024. Author list truncated; full author list available at the arXiv page.
- Grünefeld, N., Frellsen, J., and Hardmeier, C. An Isotropic Approach to Efficient Uncertainty Quantification with Gradient Norms, March 2026.
- Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., Zhu, Q., Xu, R., Zhang, R., Ma, S., Bi, X., et al. DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. *Nature*, 645(8081):633–638, September 2025. ISSN 1476-4687. doi: 10.1038/s41586-025-09422-z. Author list truncated for brevity; see Nature publication for full list.
- Højer, B. On the Notion that Language Models Reason. In *1st Workshop on Epistemic Intelligence in Machine Learning (EIML), EurIPS 2025*, 2025. doi: 10.48550/arXiv.2511.11810.
- Højer, B., Jarvis, O., and Heinrich, S. Improving Reasoning Performance in Large Language Models via Representation Engineering. In *The Thirteenth International Conference on Learning Representations*. OpenReview, 2025. doi: 10.48550/arXiv.2504.19483.
- Hora, S. C. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety*, 54(2):217–223, 1996. ISSN 0951-8320. doi: 10.1016/S0951-8320(96)00077-4. Treatment of Aleatory and Epistemic Uncertainty.
- Huang, J. and Chang, K. C.-C. Towards Reasoning in Large Language Models: A Survey. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 1049–1065, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.67.
- Hüllermeier, E. and Waegeman, W. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110(3):457–506, March 2021. ISSN 1573-0565. doi: 10.1007/s10994-021-05946-3.
- Kambhampati, S., Valmeekam, K., Bhambri, S., Palod, V., Saldyt, L., Stechly, K., Samineni, S. R., Kalwar, D., and Biswas, U. Position: Stop Anthropomorphizing Intermediate Tokens as Reasoning/Thinking Traces!, April 2025.
- Kang, Z., Zhao, X., and Song, D. Scalable Best-of-N Selection for Large Language Models via Self-Certainty. In *The Thirty-Ninth Annual Conference on Neural Information Processing Systems*, 2025. doi: 10.48550/arXiv.2502.18581.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large Language Models are Zero-Shot Reasoners. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 22199–22213. Curran Associates, Inc., 2022.
- Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., Lukošiušė, K., Nguyen, K., Cheng, N., Joseph, N., Schiefer, N., Rausch, O., Larson, R., McCandlish, S., Kundu, S., Kadavath, S., Yang, S., Henighan, T., Maxwell, T., Telleen-Lawton, T., Hume, T., Hatfield-Dodds, Z., Kaplan, J., Brauner, J., Bowman, S. R., and Perez, E. Measuring Faithfulness in Chain-of-Thought Reasoning, July 2023.

- Ling, Z., Fang, Y., Li, X., Huang, Z., Lee, M., Memisevic, R., and Su, H. Deductive Verification of Chain-of-Thought Reasoning. In *Advances in Neural Information Processing Systems*, volume 36. Curran Associates, Inc., 2023. doi: 10.48550/arXiv.2306.03872.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding R1-Zero-Like Training: A Critical Perspective. In *Conference on Language Modeling (COLM)*, 2025. doi: 10.48550/arXiv.2503.20783.
- MacKay, D. J. C. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3): 448–472, May 1992. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.448.
- Mirzadeh, I., Alizadeh, K., Shahrokhi, H., Tuzel, O., Bengio, S., and Farajtabar, M. GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models. In *The Thirteenth International Conference on Learning Representations*, 2025. doi: 10.48550/arXiv.2410.05229.
- Mondorf, P. and Plank, B. Comparing Inferential Strategies of Humans and Large Language Models in Deductive Reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9370–9402, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.508.
- Muennighoff, N., Yang, Z., Shi, W., Li, X. L., Fei-Fei, L., Hajishirzi, H., Zettlemoyer, L., Liang, P., Candès, E., and Hashimoto, T. S1: Simple test-time scaling. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pp. 20275–20321, Suzhou, China, 2025. Association for Computational Linguistics.
- Neal, R. M. *Bayesian Learning for Neural Networks*, volume 118. Springer Science & Business Media, 2012.
- Nye, M., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., Dohan, D., Lewkowycz, A., Bosma, M., Luan, D., Sutton, C., and Odena, A. Show Your Work: Scratchpads for Intermediate Computation with Language Models, November 2021.
- Sale, Y., Hofman, P., Wimmer, L., Hüllermeier, E., and Nagler, T. Second-Order Uncertainty Quantification: Variance-Based Measures, December 2023.
- Sale, Y., Hofman, P., Löhr, T., Wimmer, L., Nagler, T., and Hüllermeier, E. Label-wise Aleatoric and Epistemic Uncertainty Quantification. In *Proceedings of the Fortieth Conference on Uncertainty in Artificial Intelligence*, volume 244 of *Proceedings of Machine Learning Research*, pp. 3159–3179. PMLR, July 2024.
- Samineni, S. R., Kalwar, D., Valmeekam, K., Stechly, K., and Kambhampati, S. RL in Name Only? Analyzing the Structural Assumptions in RL post-training for LLMs. In *NeurIPS 2025 Workshop on Bridging Language, Agent, and World Models for Reasoning and Planning (LAW)*, 2025. doi: 10.48550/arXiv.2505.13697. Also accepted at NeurIPS 2025 Workshop on Foundations of Reasoning in Language Models (ForLM).
- Saparov, A. and He, H. Language Models Are Greedy Reasoners: A Systematic Formal Analysis of Chain-of-Thought, 2023. URL <https://arxiv.org/abs/2210.01240>. Published at ICLR 2023.
- Shanahan, M. Talking About Large Language Models, February 2023.
- Smith, L. and Gal, Y. Understanding Measures of Uncertainty for Adversarial Example Detection. In Globerson, A. and Silva, R. (eds.), *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, California, USA, August 6-10, 2018*, pp. 560–569. AUAI Press, 2018.
- Stechly, K., Valmeekam, K., and Kambhampati, S. Chain of Thoughtlessness? An Analysis of CoT in Planning. In *Advances in Neural Information Processing Systems*, volume 37. Curran Associates, Inc., 2024. doi: 10.48550/arXiv.2405.04776.
- Stechly, K., Valmeekam, K., Gundawar, A., Palod, V., and Kambhampati, S. Beyond Semantics: The Unreasonable Effectiveness of Reasonless Intermediate Tokens. In *The Thirty-Ninth Annual Conference on Neural Information Processing Systems*, 2025. doi: 10.48550/arXiv.2505.13775.
- Tomov, T., Fuchsgruber, D., Wollschläger, T., and Günemann, S. Entropy Is Not Enough: Uncertainty Quantification for LLMs fails under Aleatoric Uncertainty. In *NeurIPS 2025 Workshop on Structured Probabilistic Inference & Generative Modeling*, 2025.
- Turpin, M., Michael, J., Perez, E., and Bowman, S. R. Language Models Don’t Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting. In *Advances in Neural Information Processing Systems*, volume 36. Curran Associates, Inc., 2023.
- Tversky, A. and Kahneman, D. Judgment under Uncertainty: Heuristics and Biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *Science*, 185 (4157):1124–1131, September 1974. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.185.4157.1124.
- Tversky, A. and Kahneman, D. The Framing of Decisions and the Psychology of Choice. *Science*, 211(4481):453–458, 1981. doi: 10.1126/science.7455683.

- Tversky, A. and Kahneman, D. Rational Choice and the Framing of Decisions. *The Journal of Business*, 59(4): S251–S278, 1986. ISSN 0021-9398.
- Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi, E. H., Narang, S., Chowdhery, A., and Zhou, D. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *The Eleventh International Conference on Learning Representations*, 2023. doi: 10.48550/arXiv.2203.11171.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 24824–24837. Curran Associates, Inc., 2022.
- Wimmer, L., Sale, Y., Hofman, P., Bischl, B., and Hüllermeier, E. Quantifying Aleatoric and Epistemic Uncertainty in Machine Learning: Are Conditional Entropy and Mutual Information Appropriate Measures? In Evans, R. J. and Shpitser, I. (eds.), *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence*, volume 216 of *Proceedings of Machine Learning Research*, pp. 2282–2292. PMLR, July 2023. doi: 10.48550/arXiv.2209.03302.
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., et al. Qwen2 Technical Report, September 2024. Author list truncated; full author list available at the arXiv page.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al. Qwen3 Technical Report, May 2025. Author list truncated; full author list available at the arXiv page.
- Zhang, A., Chen, Y., Pan, J., Zhao, C., Panda, A., Li, J., and He, H. Reasoning Models Know When They’re Right: Probing Hidden States for Self-Verification, April 2025.
- Zhao, Z., Koishekenov, Y., Yang, X., Murray, N., and Cancedda, N. Verifying Chain-of-Thought Reasoning via Its Computational Graph. In *The Fourteenth International Conference on Learning Representations*, 2026. doi: 10.48550/arXiv.2510.09312. Accepted as Oral.

## A. Intuitions for the Uncertainty Types

The three uncertainty types defined in [subsection 2.2](#) can be unfamiliar on first reading. This appendix offers worked intuitions for each, grouped by the aleatoric–epistemic distinction ([subsection A.1](#)) and the committal–distributional distinction within aleatoric uncertainty ([subsection A.2](#)).

### A.1. Aleatoric and Epistemic Uncertainty

Consider trying to predict the outcome of a coin flip for a coin you have just been handed. You have flipped it a few times to estimate whether it is fair, but you remain unsure. Two distinct sources contribute to your uncertainty about the next flip.

**Aleatoric (irreducible) uncertainty.** A coin flip is, in principle, a random event. Even if you knew with perfect certainty that the coin was exactly fair, you still could not predict the outcome of any particular flip with confidence; the randomness is intrinsic to the task. This is aleatoric uncertainty. Flipping the coin more times to learn about its bias does not reduce it. In the language of [subsection 2.2](#), it is the noise in  $y$  given  $x$  that no parameterization  $\theta$  could remove.

**Epistemic (knowledge-based) uncertainty.** You have flipped the coin only a handful of times, and you are not yet sure whether it is fair, slightly biased toward heads, or slightly biased toward tails. This second source of uncertainty is not about the coin flip itself but about your belief regarding the coin’s underlying bias. Someone who had flipped the coin a thousand times would be far more certain about its bias and therefore better calibrated about the next flip, even though the flip itself remains as random as ever. This difference is epistemic uncertainty: it is reducible in principle, by gathering more evidence. In the language of [subsection 2.2](#), it is the variability in the predictive distribution across plausible parameterizations  $\theta$  given the observed data.

### A.2. Committal and Distributional Aleatoric Uncertainty

Within the aleatoric component, our work distinguishes two sub-types that capture different aspects of how a model commits to a prediction. Consider a language model deciding which token to generate next, with a vocabulary of 50,000 tokens. Three illustrative cases bring out the distinction:

**Case 1: One token clearly dominates.** The model assigns probability 0.95 to the token `the`, with the remaining 0.05 spread roughly evenly across the other 49,999 tokens. *Committal* uncertainty is low: the Bernoulli variance  $0.95 \cdot 0.05 = 0.0475$  is close to zero, indicating that the model is firmly committed to `the`. *Distributional* uncertainty is also low: the entropy is small because almost all the mass sits on a single token. This is the regime we call **concentration**: the top token wins by being likely.

**Case 2: One token nominally wins, but mass is broadly distributed.** The model assigns probability 0.15 to `the`, with the remaining 0.85 spread roughly evenly across 1,000 other tokens that each receive roughly 0.00085. *Committal* uncertainty is high: the Bernoulli variance  $0.15 \cdot 0.85 = 0.13$  is substantial relative to its maximum of 0.25, indicating that the model’s top choice carries little of the total probability mass. *Distributional* uncertainty is also high: the entropy is large because mass is spread across many alternatives. This is the regime we call **elimination**: the token `the` wins not by being likely in absolute terms, but by being marginally more probable than each of many alternatives that are individually even less probable.

**Case 3: Two tokens roughly tied.** The model assigns probability 0.5 to `the` and 0.5 to `a`, with negligible mass elsewhere. *Committal* uncertainty is at its maximum: the Bernoulli variance  $0.5 \cdot 0.5 = 0.25$  reflects that the top token wins by an arbitrarily small margin. *Distributional* uncertainty is low: the entropy is just  $\log_2 2 = 1$  bit, because only two tokens carry mass. This case sits at the boundary of the elimination regime, with the elimination occurring among only two alternatives rather than many; with so little of the mass distributed, distributional uncertainty stays low even though committal uncertainty is maximal.

The cases above show that the two measures can move independently: the elimination regime (Case 2) and a near-tied two-way contest (Case 3) both produce high committal uncertainty, but only Case 2 also produces high entropy. The same low-entropy reading can therefore correspond to confident commitment (Case 1) or to a near-tied two-way contest (Case 3), and only committal uncertainty distinguishes them.

## B. Extracting the Final Answer, $\hat{y}$

Computing answer uncertainty  $\mathbb{U}^A$  requires identifying, for each generation, the span  $\hat{y}$  corresponding to the model’s final answer. We split each model’s output into a trace section and an answer section using model-specific rules, since different model families structure their outputs differently. In all cases the split is a position in the generated string, with the concatenation of the two sections reproducing the original output exactly. The strategies were developed iteratively by inspecting splits on samples of correct and incorrect generations from each model-dataset combination, with the additional check that trace and answer sections concatenate back to the original output character-for-character.

**Reasoning models (DeepSeek R1, Qwen 3).** Both models wrap their internal reasoning in explicit `<think>...</think>` tags. We first locate the closing `</think>` tag and then search the post-think content for the latest occurrence of an answer-declaration pattern, such as “Final Answer”, “## Final Answer”, “Answer:”, or “Therefore, the answer is”. The split is placed immediately after this declaration, so the answer section contains only the answer itself and not the surrounding markup. By construction (see subsection 3.2), every retained reasoning model trace contains a closing `</think>` tag and a recognizable answer declaration: traces missing these were either regenerated with an additional 512-token budget or filtered out if the extended budget still did not yield an answer.

**Non-reasoning models (Qwen 2.5, Llama 3.1, Llama 3.2).** These models lack explicit reasoning sections and instead structure their outputs as plain text in which the final answer is conventionally stated in a concluding paragraph. We split at the last paragraph break (an empty line, matched as `\n\s*\n`), taking everything after this break as the answer. If no paragraph break is found, we fall back to the last single newline; if neither is present and cannot reliably split the section we exclude the sample from the  $\mathbb{U}^A$  computation.

## C. Trace Features

Based on the generated traces and the uncertainty measure we compute the following metrics for each type of uncertainty. Recall that the uncertainty types are the entropy of the predictive distribution, an approximation to the aleatoric uncertainty and finally an approximation to the epistemic uncertainty as described in subsection 3.1.

Feature	Symbol	Description
<i>Static Features</i>		
Early mean	$\mu_{\text{early}}$	Arithmetic mean of the first 25% trace.
Middle mean	$\mu_{\text{mid}}$	Arithmetic mean of the middle 50% trace.
Late Mean	$\mu_{\text{late}}$	Arithmetic mean of the last 25% trace.
<i>Dynamic Features</i>		
Slope	$\mathbf{m}$	Slope of a linear trend fitted on a trace.
Linear fit	$r^2$	Goodness-of-fit of the linear trend.

Table 2. Overview of features extracted from every uncertainty trace.

## D. Experimental Setup

This appendix lists implementation details that are needed for reproducibility but kept out of the main text. The model lineup, datasets, generation protocol, uncertainty measures, and feature definitions themselves are described in [subsection 3.2](#) and the trace/answer splitting in the section above.

### D.1. Sample Retention

Two filters are applied to each generated sample at load time, before any analysis:

1. **contains\_answer check.** A regex-based extractor attempts to pull a final answer from the answer portion of each generation. Generations for which extraction fails are dropped.
2. **LLM correctness re-evaluation.** For samples whose extracted answer differs from the reference under the regex matcher, we query an LLM auditor (`gpt-5-2025-08-07`) to re-evaluate correctness under a more lenient semantic match (e.g., 0.5 vs 1/2, mathematically equivalent expressions). When the auditor confirms correctness, we flip `is_correct` to `True`; otherwise we keep the original label.

For Qwen 3 and DeepSeek R1, traces that exhaust the 2048-token budget without emitting a closing `</think>` tag are regenerated with a forced `</think>` injection and an additional 512-token continuation; traces still missing a final answer after the regeneration are dropped.

### D.2. Predictive Modeling

The two classifiers in [subsection 4.1](#) are instantiated as

- `LogisticRegression(max_iter=1000, random_state=42)` with default L2 regularization ( $C = 1$ );
- `GradientBoostingClassifier(n_estimators=100, max_depth=3, random_state=42)`,

both with `StandardScaler` fit on the training fold inside the CV loop. Cross-validation uses `StratifiedKFold(n_splits=5, shuffle=True, random_state=42)`, constructed independently per (model, dataset) combination.

### D.3. Compute and Reproducibility

All classifiers and CV splits are seeded with `random_state=42`. Trace generation is greedy decoding, deterministic given the model weights and tokenizer. All experiments are run on H100 GPUs. Source code is released alongside the paper.

## E. Additional Analysis

### E.1. AUROC by Uncertainty Type

We distinguish between type-channel combinations, and observe that the trace-level uncertainties,  $\mathbb{U}^{\text{Tr}}$ , are generally more predictive than answer-level uncertainties,  $\mathbb{U}^{\text{A}}$ . Additionally, both aleatoric and epistemic uncertainty carry relevant predictive power as seen in columns  $\mathbb{U}_{\text{D}}$  and  $\mathbb{U}_{\text{E}}$ .

		GSM8K		ProntoQA	
		LR	GB	LR	GB
$\mathbb{U}^{\text{Tr}}$	Llama 3.1	0.781 $\pm 0.034$	0.741 $\pm 0.043$	0.777 $\pm 0.057$	0.712 $\pm 0.118$
	Llama 3.2	0.754 $\pm 0.011$	0.747 $\pm 0.033$	0.538 $\pm 0.035$	0.553 $\pm 0.046$
	Qwen 2.5	0.787 $\pm 0.027$	0.777 $\pm 0.023$	0.521 $\pm 0.036$	0.501 $\pm 0.013$
	DeepSeek R1	0.765 $\pm 0.032$	0.753 $\pm 0.043$	0.686 $\pm 0.051$	0.647 $\pm 0.076$
	Qwen 3	0.722 $\pm 0.045$	0.692 $\pm 0.029$	0.674 $\pm 0.073$	0.578 $\pm 0.044$
$\mathbb{U}^{\text{A}}$	Llama 3.1	0.682 $\pm 0.030$	0.670 $\pm 0.019$	0.669 $\pm 0.124$	0.714 $\pm 0.111$
	Llama 3.2	0.675 $\pm 0.017$	0.660 $\pm 0.027$	0.544 $\pm 0.050$	0.529 $\pm 0.047$
	Qwen 2.5	0.755 $\pm 0.030$	0.716 $\pm 0.033$	0.566 $\pm 0.057$	0.481 $\pm 0.048$
	DeepSeek R1	0.642 $\pm 0.055$	0.629 $\pm 0.048$	0.660 $\pm 0.088$	0.638 $\pm 0.069$
	Qwen 3	0.685 $\pm 0.053$	0.616 $\pm 0.043$	0.619 $\pm 0.103$	0.502 $\pm 0.113$

Table 3. AUROC scores (5-fold CV, curated features) across models and classifiers for GSM8K and ProntoQA. Values shown as mean  $\pm$  std. LR = Logistic Regression, GB = Gradient Boosting.

		GSM8K			ProntoQA		
		$\mathbb{U}_{\text{D}}$	$\mathbb{U}_{\text{C}}$	$\mathbb{U}_{\text{E}}$	$\mathbb{U}_{\text{D}}$	$\mathbb{U}_{\text{C}}$	$\mathbb{U}_{\text{E}}$
$\mathbb{U}^{\text{A}}$	Llama 3.1	0.682 $\pm 0.033$	0.630 $\pm 0.030$	0.572 $\pm 0.035$	0.629 $\pm 0.121$	0.603 $\pm 0.150$	0.523 $\pm 0.103$
	Llama 3.2	0.634 $\pm 0.031$	0.582 $\pm 0.033$	0.587 $\pm 0.033$	0.537 $\pm 0.042$	0.516 $\pm 0.063$	0.529 $\pm 0.063$
	Qwen 2.5	0.718 $\pm 0.034$	0.655 $\pm 0.009$	0.678 $\pm 0.026$	0.571 $\pm 0.045$	0.539 $\pm 0.057$	0.496 $\pm 0.025$
	DeepSeek R1	0.613 $\pm 0.048$	0.549 $\pm 0.052$	0.568 $\pm 0.031$	0.693 $\pm 0.103$	0.477 $\pm 0.102$	0.658 $\pm 0.066$
	Qwen 3	0.657 $\pm 0.020$	0.659 $\pm 0.063$	0.636 $\pm 0.045$	0.525 $\pm 0.049$	0.679 $\pm 0.133$	0.555 $\pm 0.091$
$\mathbb{U}^{\text{Tr}}$	Llama 3.1	0.717 $\pm 0.062$	0.705 $\pm 0.060$	0.778 $\pm 0.043$	0.811 $\pm 0.052$	0.776 $\pm 0.047$	0.712 $\pm 0.075$
	Llama 3.2	0.684 $\pm 0.020$	0.679 $\pm 0.020$	0.757 $\pm 0.027$	0.532 $\pm 0.051$	0.520 $\pm 0.057$	0.512 $\pm 0.041$
	Qwen 2.5	0.788 $\pm 0.031$	0.786 $\pm 0.030$	0.738 $\pm 0.020$	0.544 $\pm 0.046$	0.562 $\pm 0.050$	0.517 $\pm 0.047$
	DeepSeek R1	0.703 $\pm 0.026$	0.615 $\pm 0.024$	0.736 $\pm 0.036$	0.664 $\pm 0.060$	0.679 $\pm 0.058$	0.504 $\pm 0.064$
	Qwen 3	0.723 $\pm 0.049$	0.717 $\pm 0.043$	0.549 $\pm 0.084$	0.600 $\pm 0.052$	0.558 $\pm 0.049$	0.528 $\pm 0.087$

Table 4. AUROC scores (Logistic regression with 5-fold CV) by uncertainty type across models for GSM8K and ProntoQA. Values shown as mean  $\pm$  std.