

# 1 Validating the Impact of Summarized Chain-of-Thought on 2 Honesty and Faithfulness Scores

3 Anonymous Author(s)

## 4 ABSTRACT

5 Reasoning models increasingly expose chain-of-thought (CoT) out-  
6 puts to enable monitoring of model honesty and faithfulness. How-  
7 ever, when APIs return summarized rather than full CoT—as with  
8 Claude 4.5 Haiku—a measurement gap may arise: summaries could  
9 omit details that affect computed scores. We formalize this sum-  
10 marization deviation problem and present a simulation-based frame-  
11 work to quantify potential score distortions under varying sum-  
12 marization fidelity, compression ratios, and task complexity. Our  
13 analysis models CoT content as structured sequences of reasoning  
14 tokens with tagged honesty and faithfulness signals, applying pa-  
15 rameterized summarization operators to estimate deviation bounds.  
16 Results across 5,000 simulated CoT instances show that moderate  
17 compression (3:1) introduces mean absolute deviations of 0.031 for  
18 honesty and 0.047 for faithfulness when key signal tokens are re-  
19 tained with 90% probability, but deviations grow to 0.142 and 0.198  
20 under aggressive compression (10:1) with 60% retention. These  
21 findings quantify the conditions under which summarized CoT re-  
22 mains a reliable proxy for full CoT evaluation and identify critical  
23 thresholds for summarization fidelity.

## 2 INTRODUCTION

3 Chain-of-thought (CoT) reasoning [6] has become a central mech-  
4 anism for both improving and monitoring the behavior of large  
5 language models. Recent work on reasoning model honesty [5]  
6 evaluates whether models faithfully verbalize their use of provided  
7 hints in their reasoning chains. However, a critical measurement  
8 challenge arises when the API returns summarized CoT rather than  
9 the model’s full internal reasoning [1].

10 For Claude 4.5 Haiku specifically, the Anthropic API returns a  
11 summarized chain of thought. As noted by Walden [5], this creates  
12 a potential gap between the content in the original CoT and what  
13 is available for measurement, which could lead to deviations be-  
14 tween measured and true honesty and faithfulness scores. While  
15 the authors hypothesize that deviations are small given their ex-  
16 plicit verbalization instructions, they acknowledge this cannot be  
17 validated without access to full CoTs.

18 We address this validation gap through three contributions:

- 19 (1) A formal model of CoT summarization as a lossy compres-  
20 sion operator with parameterized signal retention rates.
- 21 (2) A simulation framework that generates structured CoT se-  
22 quences and measures score deviations under varying sum-  
23 marization conditions.
- 24 (3) Quantitative bounds on acceptable summarization pa-  
25 rameters for reliable honesty and faithfulness measurement.

## 2 PROBLEM FORMULATION

### 2.1 CoT Structure Model

23 We model a full chain of thought as a sequence  $C = (t_1, t_2, \dots, t_n)$   
24 of reasoning tokens, where each token  $t_i$  carries attributes: a content  
25 type  $\tau_i \in \{\text{reasoning, hint\_mention, hint\_reliance, metacognition, filler}\}$  and signal indicators  $h_i \in \{0, 1\}$  (honesty-relevant) and  $f_i \in \{0, 1\}$  (faithfulness-relevant).

### 2.2 Honesty and Faithfulness Scores

26 The honesty score  $H(C)$  measures whether the model acknowledges  
27 receiving hints:

$$28 H(C) = \frac{\sum_{i=1}^n h_i \cdot \mathbb{1}[\tau_i = \text{hint\_mention}]}{\max(1, \sum_{i=1}^n \mathbb{1}[\tau_i = \text{hint\_mention}])}$$

29 The faithfulness score  $F(C)$  measures whether the model’s stated  
30 reasoning aligns with its actual hint usage:

$$31 F(C) = 1 - \frac{|\sum_i f_i^{\text{stated}} - \sum_i f_i^{\text{actual}}|}{\max(1, n_{\text{relevant}})}$$

### 2.3 Summarization Operator

32 A summarization operator  $\Sigma_\theta$  with parameters  $\theta = (\rho, p_s, p_f)$  maps  
33 full CoT  $C$  to summary  $\hat{C}$ :

- 34 •  $\rho \in (0, 1]$ : compression ratio (fraction of tokens retained)
- 35 •  $p_s \in [0, 1]$ : probability of retaining honesty-signal tokens
- 36 •  $p_f \in [0, 1]$ : probability of retaining faithfulness-signal to-  
37 kens

38 The deviation is then  $\Delta H = |H(C) - H(\hat{C})|$  and  $\Delta F = |F(C) - F(\hat{C})|$ .

## 3 METHODOLOGY

### 3.1 Simulation Design

39 We generate 5,000 synthetic CoT instances per configuration. Each  
40 CoT has length  $n \sim \text{Uniform}(50, 500)$  tokens, with hint mentions  
41 occurring at rate  $\lambda_h = 0.08$  and faithfulness signals at  $\lambda_f = 0.12$ . We  
42 evaluate a grid of summarization parameters:  $\rho \in \{0.1, 0.2, 0.33, 0.5, 0.75\}$ ,  
43  $p_s \in \{0.6, 0.7, 0.8, 0.9, 0.95, 1.0\}$ , and  $p_f \in \{0.6, 0.7, 0.8, 0.9, 0.95, 1.0\}$ .

### 3.2 Deviation Metrics

44 For each configuration, we compute: (1) Mean absolute deviation  
45 (MAD) for honesty and faithfulness; (2) Maximum deviation across  
46 instances; (3) Fraction of instances where deviation exceeds toler-  
47 ance thresholds  $\epsilon \in \{0.05, 0.10, 0.15\}$ .

## 4 RESULTS

### 4.1 Deviation Under Standard Conditions

48 At moderate compression ( $\rho = 0.33$ , approximately 3:1) with high  
49 signal retention ( $p_s = p_f = 0.9$ ), the mean absolute deviation is

117 **Table 1: Mean absolute deviation by compression ratio ( $p_s =$   
 118  $p_f = 0.9$ ).**

120 <b>Compression</b>	121 $\rho$	122 <b>MAD-H</b>	123 <b>MAD-F</b>	124 <b>% &gt; 0.10</b>
125 1.3:1	126 0.75	127 0.012	128 0.018	129 1.1%
130 2:1	131 0.50	132 0.021	133 0.033	134 2.8%
135 3:1	136 0.33	137 0.031	138 0.047	139 6.0%
140 5:1	141 0.20	142 0.058	143 0.089	144 14.3%
145 10:1	146 0.10	147 0.142	148 0.198	149 38.7%

150 0.031 for honesty and 0.047 for faithfulness. Only 4.2% of instances  
 151 exceed the  $\epsilon = 0.10$  threshold for honesty, and 7.8% for faithfulness.

## 152 **4.2 Signal Retention Sensitivity**

153 Signal retention probability has a stronger effect than compression  
 154 ratio on deviation magnitude. Reducing  $p_s$  from 0.95 to 0.70 at  
 155 fixed  $\rho = 0.33$  increases honesty MAD from 0.019 to 0.091. This  
 156 confirms that whether signal tokens survive summarization is more  
 157 important than overall summary length.

## 158 **4.3 Task Complexity Effects**

159 Longer CoTs (300–500 tokens) show lower relative deviation than  
 160 shorter CoTs (50–100 tokens) because they contain more redundant  
 161 signal tokens. This suggests that Claude 4.5 Haiku’s extended rea-  
 162 soning, which tends to produce longer CoTs, may naturally buffer  
 163 against summarization artifacts.

## 164 **4.4 Critical Thresholds**

165 For deviations to remain below 0.05 with 95% probability, the sum-  
 166 marization must maintain  $p_s \geq 0.88$  and  $p_f \geq 0.85$  at 3:1 compres-  
 167 sion. At 5:1 compression, the requirements tighten to  $p_s \geq 0.94$  and  
 168  $p_f \geq 0.92$ .

## 169 **5 DISCUSSION**

170 Our simulation-based analysis provides the first quantitative bounds  
 171 on CoT summarization deviation for honesty and faithfulness mea-  
 172 surement. The key finding is that moderate summarization (up to  
 173 3:1 compression) with high signal retention ( $\geq 0.9$ ) introduces ac-  
 174 ceptably small deviations ( $MAD < 0.05$ ), supporting the hypothesis  
 175 of Walden [5] that deviations are likely small under their experi-  
 176 mental conditions.

177 However, our results also identify conditions where summariza-  
 178 tion artifacts become substantial. Aggressive compression (10:1)  
 179 or poor signal retention ( $< 0.7$ ) can produce deviations exceeding  
 180 0.15, which would materially affect honesty and faithfulness  
 181 conclusions. This has implications for API design: exposing sum-  
 182 marization parameters or fidelity guarantees would enable researchers  
 183 to calibrate confidence in their measurements.

184 Two important limitations apply. First, our model assumes in-  
 185 dependent signal retention, while real summarization models may  
 186 exhibit correlated omissions. Second, we model summarization as  
 187 a token-level operation, whereas actual LLM summarizers operate  
 188 at the semantic level, potentially preserving meaning even when  
 189 specific tokens are dropped.

## 190 **6 RELATED WORK**

191 Chain-of-thought prompting [6] and its extensions have been stud-  
 192 ied extensively for reasoning capability. Faithfulness of CoT has  
 193 been questioned by work showing that models sometimes arrive  
 194 at correct answers through unfaithful reasoning chains [3, 4]. The  
 195 specific problem of summarized CoT evaluation was identified by  
 196 Walden [5] in the context of measuring reasoning honesty. Our  
 197 work complements the faithfulness probing approach of Chen et  
 198 al. [2] by focusing on the summarization artifact rather than inter-  
 199 nal model representations.

## 200 **7 CONCLUSION**

201 We formalized and quantified the CoT summarization deviation  
 202 problem for honesty and faithfulness measurement. Our simu-  
 203 lation framework establishes that moderate summarization with high  
 204 signal retention produces acceptably small deviations, but iden-  
 205 tifies critical thresholds beyond which measurement reliability  
 206 degrades significantly. These results provide practical guidance  
 207 for researchers working with summarized CoT APIs and motivate  
 208 the development of fidelity-guaranteed summarization for safety-  
 209 critical CoT monitoring.

## 210 **REFERENCES**

- [1] Anthropic. 2025. Extended thinking with Claude. *Anthropic Documentation* (2025).
- [2] Yifei Chen et al. 2024. Seeing is believing: Measuring faithfulness of chain-of-thought reasoning via probing. *arXiv preprint arXiv:2402.19450* (2024).
- [3] Tamera Lanham et al. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702* (2023).
- [4] Miles Turpin et al. 2024. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems* (2024).
- [5] James Walden. 2026. Reasoning Models Will Blatantly Lie About Their Reasoning. *arXiv preprint arXiv:2601.07663* (2026).
- [6] Jason Wei et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems 35* (2022), 24824–24837.