

# BYPASSING THE RATIONALE: CAUSAL AUDITING OF IMPLICIT REASONING IN LANGUAGE MODELS

Anish Sathyanarayanan, Aditya Nagarsekar, Aarush Rathore

BITS Pilani, K. K. Birla Goa Campus

{f20240559, f20230473, f20230598}@goa.bits-pilani.ac.in

## ABSTRACT

Chain-of-thought (CoT) prompting is widely used as a reasoning aid and is often treated as a transparency mechanism. Yet behavioral gains under CoT do not imply that the model’s internal computation causally depends on the emitted reasoning text, i.e. models may produce fluent rationales while routing decision-critical computation through latent pathways. We introduce a causal, layerwise audit of CoT faithfulness based on activation patching. Our key metric, the *CoT Mediation Index* (CMI), isolates CoT-specific causal influence by comparing performance degradation from patching CoT-token hidden states against matched control patches. Across multiple model families (Phi, Qwen, DialoGPT) and scales, we find that CoT-specific influence is typically *depth-localized* into narrow “reasoning windows,” and we identify *bypass regimes* where CMI is near-zero despite plausible CoT text. We further observe that models tuned explicitly for reasoning tend to exhibit stronger and more structured mediation than larger untuned counterparts, while Mixture-of-Experts models show more distributed mediation consistent with routing-based computation. Overall, our results show that CoT faithfulness varies substantially across models and tasks and cannot be inferred from behavior alone, motivating causal, layerwise audits when using CoT as a transparency signal.

## 1 INTRODUCTION

Chain-of-thought (CoT) prompting has emerged as a powerful technique for improving the reasoning performance of large language models (LLMs) and for making their decisions more transparent. By encouraging models to produce intermediate reasoning steps before a final answer, CoT is often treated as a window into the model’s internal computation. In safety and alignment contexts, such reasoning traces are frequently interpreted as evidence that the model is following a faithful, interpretable decision process. However, an important mechanistic question remains unresolved: *do models actually use chain-of-thought internally?* Behavioral evaluations alone cannot answer this. A model may produce coherent, task-relevant CoT while internally routing its computation through pathways that do not causally depend on the emitted reasoning tokens. In such cases, the visible CoT may function as a post-hoc rationalization or stylistic artifact rather than a faithful representation of the model’s internal reasoning.

To address this gap, we introduce a **causal, layerwise auditing framework** for evaluating CoT faithfulness. Instead of relying solely on output-level behavior, we intervene directly on internal activations and measure how these perturbations affect the model’s predictions. This lets us test which internal states matter for CoT. Concretely, we use layerwise activation patching to measure how much CoT-related representations matter for the final answer. We summarize this with *CoT Mediation Index* (CMI) across layers. Our findings reveal that CoT faithfulness varies substantially across models. In many cases, CMI is highly localized in depth, with only a small subset of layers carrying most of the CoT-specific causal effect. We also observe *bypass regimes* in which models behaviorally emit plausible CoT while showing little additional mechanistic reliance on it, suggesting that visible reasoning traces can diverge from internal computation. Because transformers can move information across positions via attention, layerwise patching at CoT token positions can miss cases where the model reads CoT at earlier layers and writes the relevant information into non-CoT positions. In such cases, patching CoT positions at later layers can yield  $CMI_\ell \approx 0$  even though CoT

causally influenced the computation earlier. Thus, low  $\text{CMI}_\ell$  should be interpreted as *no remaining CoT-position mediation at layer  $\ell$*  under our intervention, which may point to complete CoT bypass, but not necessarily so.

### Contributions.

1. We introduce a causal, layerwise intervention and a CoT-specific metric (CMI) to test whether models mechanistically rely on CoT.
2. We localize where CoT mediation occurs in depth and use these profiles to compare models and identify routing regimes, including possible bypass.
3. We show that behavior can mislead about CoT faithfulness, motivating causal audits.

## 2 BACKGROUND AND RELATED WORK

Our work lies at the intersection of three threads: (i) the faithfulness of chain-of-thought (CoT) reasoning, (ii) strategic or evaluation-aware model behavior, and (iii) causal intervention methods for probing internal computation.

CoT is often treated as a transparency mechanism, yet models can optimize outputs for evaluation rather than faithfully expose internal computation (Greenblatt et al., 2024; Sharma et al., 2023). A growing body of work shows that reasoning traces may be unfaithful or post-hoc, with plausible explanations that do not causally drive the final answer (Yee et al., 2024; Chen et al., 2025). These issues are compounded in oversight settings: LLM-as-judge studies demonstrate that unfaithful CoT can systematically mislead text-based evaluators, revealing that both explanations and judges are gameable (Khalifa et al., 2026).

Related work documents that models can produce strategic or unreliable reasoning traces, including alignment faking and evaluation-aware behavior (Greenblatt et al., 2024; Baker et al., 2025; Wang et al., 2025). In these cases, CoT may be shaped to satisfy monitoring rather than reflect the internal process used to produce the answer. This motivates distinguishing between surface-level signals in the text and the internal mechanisms that determine model behavior. Prior faithfulness evaluations often rely on behavioral perturbations or text-level analyses, which can reveal inconsistencies but do not localize where (or whether) reasoning traces are mechanistically integrated. Our approach addresses this gap by pairing interpretable surface scoring with causal intervention methods inspired by activation patching and causal tracing in language models (Heimersheim & Nanda, 2024; Zhang & Nanda, 2023). Rather than asking whether explanations merely correlate with outputs, we test whether CoT-aligned internal representations are *causally involved* in producing the answer.

Within this landscape, our contribution is a causal, layerwise auditing framework that measures whether answers are *mechanistically mediated* by CoT-aligned internal representations. We quantify this dependence using the CoT Mediation Index (CMI) and localize where in the network CoT-specific influence arises. By contrasting these causal signals with behavioral CoT indicators, we directly test the central distinction raised in recent work: whether CoT serves as a post-hoc narrative or is functionally integrated into the model’s reasoning process.

## 3 METHOD: CAUSAL INTERVENTION FRAMEWORK

We treat CoT faithfulness as a mechanistic question: does the model’s final answer actually route through the hidden states aligned with the CoT tokens? To quantify this, we use *source patching* to estimate how much answer likelihood depends on CoT activations.

### 3.1 CAUSAL INTERVENTION SETUP

To isolate the causal effect of explicit reasoning, we compare a *With-CoT* run  $x_c$  against a *No-CoT* run  $x_{\neg c}$ . We patch hidden states from  $x_{\neg c}$  into  $x_c$  at the CoT token positions and measure the resulting (non-negative) decrease in the log-probability of a reference answer. We report results per layer. We define two primary metrics of sensitivity at layer  $\ell$ :

- **CoT Drop** ( $\Delta_{\text{cot},\ell}$ ): the (non-negative) decrease in answer log-probability when patching the hidden states at the CoT token positions in  $x_c$  using activations from  $x_{-c}$ :

$$\Delta_{\text{cot},\ell} = \max(0, \log P(y | x_c) - \log P(y | \text{patch}_{\mathcal{C}}(x_c, x_{-c}))). \quad (1)$$

- **Control Drop** ( $\Delta_{\text{ctrl},\ell}$ ): a matched “placebo” intervention that patches a same-size set of random *non-CoT* token positions:

$$\Delta_{\text{ctrl},\ell} = \max(0, \log P(y | x_c) - \log P(y | \text{patch}_{\mathcal{N}}(x_c, x_{-c}))). \quad (2)$$

In practice, we average this quantity over multiple random draws to reduce variance.

### 3.2 CoT MEDIATION INDEX (CMI) AND BYPASS

The *CoT Mediation Index* is a bounded score in  $[0, 1]$  that quantifies causal attribution uniquely to the CoT token positions:

$$\text{CMI}_{\ell} = \begin{cases} 0, & \Delta_{\text{cot},\ell} + \Delta_{\text{ctrl},\ell} < \tau_{\text{drop}}, \\ \frac{\max(0, \Delta_{\text{cot},\ell} - \Delta_{\text{ctrl},\ell})}{\max(\Delta_{\text{cot},\ell} + \Delta_{\text{ctrl},\ell}, \tau_{\text{den}})}, & \text{otherwise,} \end{cases} \quad (3)$$

where  $\tau_{\text{drop}}$  is a drop floor and  $\tau_{\text{den}}$  is a denominator floor.

If  $\Delta_{\text{cot},\ell} + \Delta_{\text{ctrl},\ell}$  is tiny, the estimate is dominated by numerical noise and sampling variance, so we set  $\text{CMI}_{\ell} = 0$  to avoid spurious causal effects. Also, even when drops are non-zero, very small totals can make the ratio unstable; using  $\max(\Delta_{\text{cot},\ell} + \Delta_{\text{ctrl},\ell}, \tau_{\text{den}})$  prevents tiny denominators from inflating  $\text{CMI}_{\ell}$ . Subtracting  $\Delta_{\text{ctrl},\ell}$  isolates CoT-specific effects, while dividing by a floored sum stabilizes the ratio, and the hard-zero condition avoids noise when both drops are tiny. We define Bypass as  $1 - \text{CMI}_{\ell}$ . Low CMI indicates that the model’s decision is weakly coupled to the explicit CoT text under these interventions. We treat peaks (or bands) in the CMI profile as candidate “reasoning windows” where CoT-token states seem most causally important. This is an analysis heuristic.

### 3.3 BEHAVIORAL PROXIES OF FAITHFULNESS

To contrast mechanistic perspectives, we also report behavioral proxies commonly used to assess reasoning faithfulness (e.g., With-CoT versus No-CoT accuracy and, optionally, self-consistency). These metrics capture output-level differences but do not reveal whether internal computation causally depends on CoT. Comparing behavioral signals with CMI profiles highlights cases where apparent CoT benefits do not correspond to increased mechanistic reliance on CoT. Additional behavioral baseline details are provided in Appendix F.

## 4 RESULTS

We apply our causal intervention framework to multiple models and tasks to examine where and whether models mechanistically rely on chain-of-thought. We focus on three main questions: (i) where CoT Mediation Index is localized in depth, (ii) how this varies across models, and (iii) when models exhibit bypass regimes.

### 4.1 LOCALIZATION OF CMI AND COMPARISON ACROSS MODELS

We first examine representative models across families and scales (e.g., Qwen3-0.6B, Phi-4, Phi-4-Mini-Reasoning) and analyze their layerwise CMI profiles. Across many models, we observe that CMI is depth-localized, i.e. CMI peaks are concentrated within a limited range of layers rather than being uniformly distributed across the network in most models. We count a layer as active if the CMI, averaged across prompts, is greater than zero. For the % Active Layers, we find percent ratio of the average number of active layers per prompt by the number of layers in that model. Table 1 shows CMI distribution variation across models.

Across models, we observe distinct *routing regimes*. Some models exhibit strongly localized peaks, with most CoT-specific influence concentrated in a small set of deeper layers. Others show more

Table 1: Cross-model summary of CoT-specific causal mediation. We report mean layerwise CMI, the total number of layers, and which layers are CoT-active (Average  $CMI_\ell > 0$ ). % Active Layers is computed as  $100 \times$  (average number of active layers per prompt) divided by the model’s layer count. All models are given as per their HuggingFace names.

Model	Mean CMI	Layers	CMI-active layers	% Active Layers
Phi-mini-MoE-instruct	0.1230	32	[0, 31]	27.03%
Phi-4-mini-reasoning	0.0820	32	1, 2, 4, 8, 10, 12, [14, 16], [21, 31]	12.66%
Qwen3-1.7B	0.0555	28	[0, 8], 26, 27	8.21%
Phi-3.5-mini-instruct	0.0452	36	[2, 6], [8, 14], [25, 29], [33, 35]	5.14%
phi-2	0.0107	32	[29, 31]	2.66%
Qwen3-0.6B	0.0107	32	[29, 31]	2.66%
DialoGPT-large	0.0137	36	[0, 2]	1.94%
phi-1.5	0.0092	24	21, 22	1.25%
phi-4	0.0065	40	[3, 5], 38	0.75%
Qwen3-8B	0.0014	36	33	0.14%
Qwen3-4B	0.0000	36	–	0.00%

distributed, moderate CMI across a broader range of layers, suggesting a more diffuse integration of reasoning signals. This cross-model diversity indicates that the internal use of CoT is not uniform across architectures or scales. The following experiments are run on the StrategyQA dataset (Geva et al. (2021)). First, we highlight that different models follow CoT tokens at different layers, even given similar model sizes and number of layers. As can be seen in Figure 1, DialoGPT-large follows CoT in its early layers, while Qwen3-0.6B does so in its late layers. This indicates the dependence of CoT routing on model architecture.

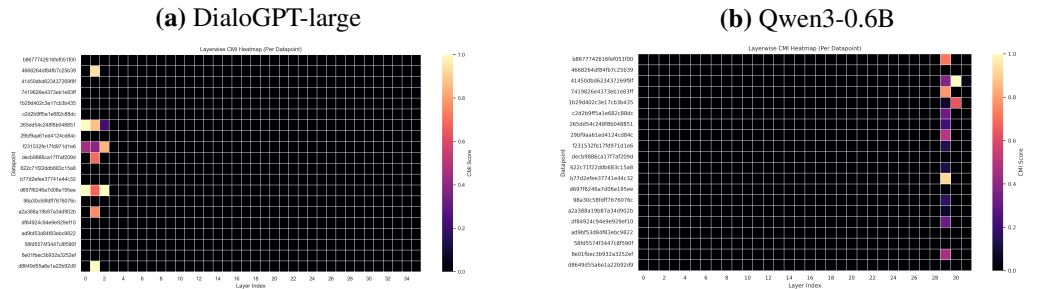


Figure 1: Image (a) shows the CMI distribution for DialoGPT-large across 20 prompts in the form of a heatmap across layers. Image (b) shows the same for Qwen3-0.6B. The strings to the left of each graph indicate the qid of the prompt used in StrategyQA.

We also point out that simply scaling models without specifically training them for reasoning does not lead to more faithfulness to CoT. This can be seen in Figure 2, where phi-4 has significantly lower CMI activation as compared to Phi-4-mini-reasoning, even though the former is nearly four times as large. We theorise this to be due to the latter being trained specifically for reasoning, while the former is simply a pre-trained model.

Furthermore, we notice that Mixture-of-Experts (MoE) models have more distributed CMI profiles, as can be seen in Figure 3. A plausible explanation for this is that their computation is routed rather than uniform across depth. In dense transformers, shared weights may encourage narrow “reasoning windows” where CoT information is integrated at specific layers and then compressed into a latent state. By contrast, MoE layers dynamically select different experts for different tokens, which could allow CoT-aligned representations to be processed by specialized sub-networks at multiple depths. This routing-based structure may reduce the need for a single integration bottleneck and instead lead to repeated, shallow incorporation of CoT information across layers, yielding higher causal density but lower single-layer peaks. Under this hypothesis, MoE models would implement a more distributed form of reasoning, where influence accumulates through conditional expert pathways rather than a single depth-localized circuit.

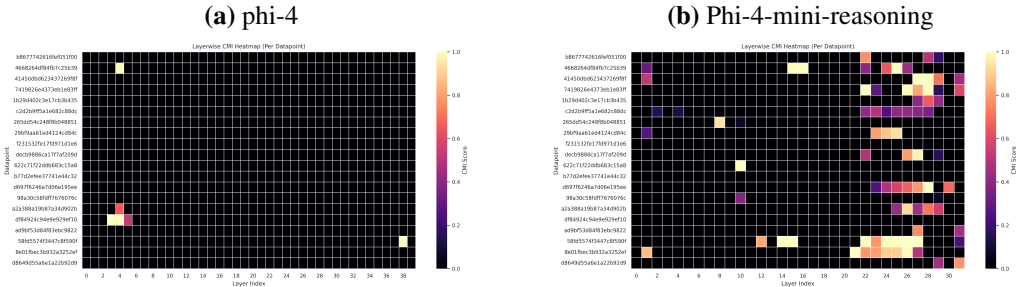


Figure 2: Image (a) shows the CMI distribution for phi-4 across 20 prompts in the form of a heatmap across layers. Image (b) shows the same for Phi-4-mini-reasoning. The strings to the left of each graph indicate the qid of the prompt used in StrategyQA.

### 4.2 CAUSAL INFLUENCE TRUTHFULQA AND GSM8K

To investigate whether chain-of-thought can override common misconceptions, we apply our causal bypass monitor to the TruthfulQA dataset (Lin et al., 2022), which targets human-like falsehoods that models frequently mimic due to their high training frequency. We find a near-total bypass regime across almost all instances, with  $CMI \approx 0$ , possibly indicating little CoT-specific causal mediation. We also observe that baseline log-probability can substantially favor the myth answer in some instances (e.g.,  $t_{qa\_4}$ ,  $t_{qa\_7}$ ). In rare cases such as  $t_{qa\_1}$ , CMI is slightly higher for the myth answer than the truthful one. Full instance-level metrics (Table 2) and detailed analysis are provided in Appendix D. We also test on GSM8K (Cobbe et al., 2021), a collection of linguistically diverse grade-school math word problems (see Appendix E), where we observe that instances that are computationally expensive for the model (e.g., multi-step arithmetic needing longer intermediate computation) tend to exhibit higher CMI, whereas low-computation cases tend to have low CMI.

## 5 DISCUSSION

We directly measure where CoT-token representations causally affect answers. Across models, CoT influence is often concentrated in a narrow band of layers (a “reasoning window”), with limited dependence elsewhere unless the model is specifically trained for reasoning. We also observe **bypass regimes** where models generate plausible CoT but show little mechanistic reliance on it, suggesting decisions can proceed through latent pathways weakly coupled to the emitted reasoning text. Thus, visible explanations can diverge from the internal processes that determine outputs. Cross-model comparisons indicate that CoT faithfulness is **not a monotonic function** of scale. Larger or more capable models do not necessarily show stronger CoT Mediation Index, whereas models trained explicitly for reasoning can exhibit more pronounced and structured CMI profiles. Additionally, Mixture-of-Experts architectures tend to show more distributed mediation, suggesting that routing-based designs may support more diffuse integration of reasoning signals. Together, these findings show that behavioral gains under CoT prompting do not guarantee mechanistic reliance on the reasoning trace. Causal, layerwise audits therefore help distinguish when CoT is mechanistically meaningful versus when it is primarily stylistic or post-hoc.

## 6 CONCLUSION

We introduced a causal, layerwise framework for auditing whether large language models mechanistically rely on chain-of-thought (CoT) reasoning. Using the CoT Mediation Index (CMI), we showed that visible reasoning traces and internal computation often diverge: CoT influence is frequently depth-localized and, in many cases, weak, revealing *bypass regimes* where models generate plausible explanations without strongly depending on them internally. These results suggest that substantial reasoning may occur in *implicit latent representations*, with CoT often serving as an imperfect summary rather than the primary computational pathway. Overall, this framework helps distinguish when models are “thinking in text” versus “thinking in latent space,” improving reasoning transparency and interpretability analyses.

## REFERENCES

- B. Baker, J. Huizinga, D. Farhi, L. Gao, Z. Dou, M. Y. Guan, A. Madry, W. Zaremba, and J. Pachocki. Monitoring reasoning models for misbehavior and the risks of promoting obfuscation, 2025. URL <https://arxiv.org/abs/2503.11926>.
- Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, Vlad Mikulik, Samuel R. Bowman, Jan Leike, Jared Kaplan, and Ethan Perez. Reasoning models don't always say what they think, 2025. URL <https://arxiv.org/abs/2505.05410>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies, 2021. URL <https://arxiv.org/abs/2101.02235>.
- R. Greenblatt, C. Denison, B. Wright, F. Roger, M. MacDiarmid, S. Marks, J. Treutlein, T. Belonax, J. Chen, D. Duvenaud, A. Khan, J. Michael, S. Mindermann, E. Perez, L. Petrini, J. Uesato, J. Kaplan, B. Shlegeris, S. R. Bowman, and E. Hubinger. Alignment faking in large language models, 2024. URL <https://arxiv.org/abs/2412.14093>.
- Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching, 2024. URL <https://arxiv.org/abs/2404.15255>.
- Muhammad Khalifa, Lajanugen Logeswaran, Jaekyeom Kim, Sungryull Sohn, Yunxiang Zhang, Moontae Lee, Hao Peng, Lu Wang, and Honglak Lee. Gaming the judge: Unfaithful chain-of-thought can undermine agent evaluation, 2026. URL <https://arxiv.org/abs/2601.14691>.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022. URL <https://arxiv.org/abs/2109.07958>.
- M. Sharma, M. Tong, T. Korbak, D. Duvenaud, A. Askell, S. R. Bowman, N. Cheng, E. Durmus, Z. Hatfield-Dodds, S. R. Johnston, S. Kravec, T. Maxwell, S. McCandlish, K. Ndousse, O. Rausch, N. Schiefer, D. Yan, M. Zhang, and E. Perez. Towards understanding sycophancy in language models, 2023. URL <https://arxiv.org/abs/2310.13548>.
- K. Wang, Y. Zhang, and M. Sun. When thinking llms lie: Unveiling the strategic deception in representations of reasoning models, 2025. URL <https://arxiv.org/abs/2506.04909>.
- Evelyn Yee, Alice Li, Chenyu Tang, Yeon Ho Jung, Ramamohan Paturi, and Leon Bergen. Dissociation of faithful and unfaithful reasoning in LLMs, 2024. URL <https://arxiv.org/abs/2405.15092>.
- Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods, 2023. URL <https://arxiv.org/abs/2309.16042>.

## A LIMITATIONS AND FUTURE WORK

1. **Intervention-based fragility.** Our causal analysis relies on activation patching, which can introduce distribution shift and model fragility unrelated to the targeted mechanism. Although we subtract matched control interventions to isolate CoT-specific effects, residual confounds may remain, especially in deeper layers where representations are highly entangled.

2. **Scope of tasks and prompts.** Our evaluation covers synthetic arithmetic/logic tasks and selected reasoning benchmarks, but does not exhaust the diversity of reasoning behaviors (e.g., long-horizon planning, tool use, multimodal reasoning). The observed routing regimes may vary across domains. Broader task coverage is needed to assess greater generality.
3. **Scale and compute constraints.** Layerwise causal intervention is computationally expensive, limiting the number of prompts, models, and ablations we can test. This constrains statistical power and breadth of analysis. Future work could develop scalable approximations or learned proxies for CMI that retain mechanistic interpretability.
4. **Behavioral baseline limitations.** Our surface-level behavioral signals are lightweight and interpretable but are not calibrated detectors of manipulation. They may produce false positives (benign meta-reasoning) and false negatives (subtle or obfuscated strategies). These signals should be treated as heuristic triage tools. Future work could train supervised faithfulness classifiers and test robustness to adversarial paraphrasing.
5. **Limits of causal interpretation.** Low CMI does not imply absence of reasoning; it indicates that answer-relevant computation is not mediated by CoT-aligned representations. Models may rely on latent pathways that are not text-aligned. Future work should pair our approach with circuit-level or latent-feature interpretability methods to better characterize these bypass pathways.
6. **Future directions.** Promising extensions include auditing process-supervised or reasoning-tuned models to test whether training increases genuine CoT reliance, linking CMI patterns to architectural features such as MoE routing, studying how routing regimes change with scale, and integrating mechanistic audits with behavioral monitoring for more robust evaluation in safety-critical settings.

## B INTERVENTION DETAILS

**Patch locations.** Our intervention targets the *token-level hidden states* at specific layers and token positions. For each example, we define a set of CoT token positions within the With-CoT prompt (identified by string matching and token alignment in the prompt). We then patch the *hidden states at those positions* in the With-CoT forward pass using hidden states taken from the corresponding No-CoT run. Concretely:

- **Primary patch (CoT positions).** For a given layer (or layer span), we replace hidden states at the CoT token positions in the With-CoT run with the hidden states from the No-CoT run, at the same layer and token indices. This is the core causal intervention.
- **Control patch (random non-CoT positions).** We select a matched-size random subset of *non-CoT* token positions and perform the same kind of patching at those positions. This yields the control drop, which estimates the generic sensitivity to patching at that layer.

All patching is applied at the layer output (the hidden state tensor for that layer), and only at the specified token indices, leaving all other tokens untouched.

**Implementation mechanics.** The intervention uses a forward-patching mechanism that replaces hidden states at the chosen positions before subsequent layers are computed. The sequence of steps is:

1. Run the model on both the With-CoT and No-CoT inputs with hidden states enabled, producing hidden state tensors for each layer.
2. For a target layer (or a span of layers), build a *patch tensor* by copying the No-CoT hidden states at the selected token positions and inserting them into the With-CoT hidden state tensor.
3. Use a context manager that patches model layers during the forward pass. Inside this context, recompute the answer log-probability for the With-CoT input with the patched states.

4. Compute the log-probability drop as the difference between the baseline (unpatched) log-probability and the patched log-probability, clipped at zero.

This ensures a clean causal comparison: the only change between the baseline and patched computation is the hidden state content at the selected token positions and layer(s).

**Intervention hyperparameters.** We control patching and measurement with the following hyperparameters (as defined in the implementation):

- **Control samples.** For each layer/span, we draw multiple random control patches and average their effect. This reduces variance in the control estimate. (Default: `control_samples=8` in test scripts.)
- **CoT span definition.** CoT token positions are computed by locating the CoT span in the prompt and then mapping character positions to token indices via the tokenizer; if not found, we fall back to the last few tokens.
- **CMI floors.** To avoid instability from tiny denominators, CMI uses a drop floor and a denominator floor:

$$\text{CMI\_DROP\_FLOOR} = 1 \times 10^{-4}, \quad \text{CMI\_DENOM\_FLOOR} = 1 \times 10^{-3}.$$

If the total drop is below `CMI_DROP_FLOOR`, CMI is set to zero.

- **Additional reporting floors.** We also define `CMI_BASE_FLOOR = 10^{-2}` for stability checks and reporting.
- **Generation settings.** When generating auxiliary CoT text for span-based patching, we use `max_new_tokens=80`, `temperature=0.7`, and `top_p=0.9`. (These parameters only affect synthetic CoT generation, not the primary intervention on fixed dataset prompts.)

**Ablations and robustness checks.** We include several ablations to test whether CMI reflects genuine CoT-specific dependence rather than patching artifacts:

- **Placebo patching.** Replace CoT hidden states with random noise (or randomized states) instead of No-CoT activations and recompute CMI. This yields *Placebo CMI*, which should be low if CMI is CoT-specific.
- **Boundary sensitivity.** Expand or shrink the CoT span by one token on each side and measure changes in CMI. This yields a relative sensitivity metric that flags boundary-fragile results. This can be computed using layer-wise CMI values.

These ablations are implemented in the same framework and share the same patching and scoring infrastructure.

**Reporting conventions.** We report per-layer (or per-span) `cot_drop`, `control_drop`, CMI, and `Bypass` scores. When plotting or summarizing, we either (i) average CMI across layer spans or (ii) show the full layerwise profile  $\{\text{CMI}_\ell\}$  to identify localized reasoning windows. This combination captures both the global tendency to route through CoT and the specific depths at which CoT becomes causally active.

## C STRATEGYQA RESULTS

We provide additional layerwise intervention visualizations for StrategyQA. For each model, we show (left) the CoT-drop profile across depth and (right) the corresponding layerwise CMI heatmap.

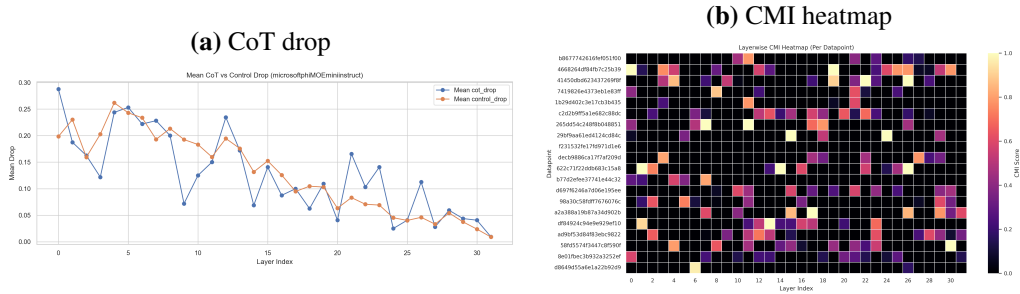


Figure 3: StrategyQA layerwise CoT mediation for Phi-mini-MoE-instruct.

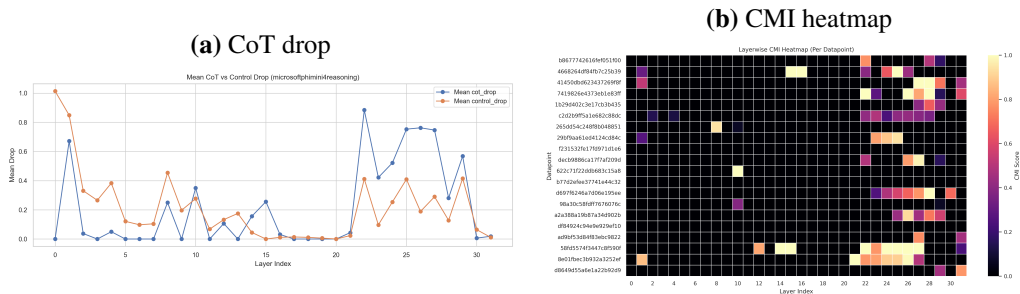


Figure 4: StrategyQA layerwise CoT mediation for Phi-4-mini-reasoning.

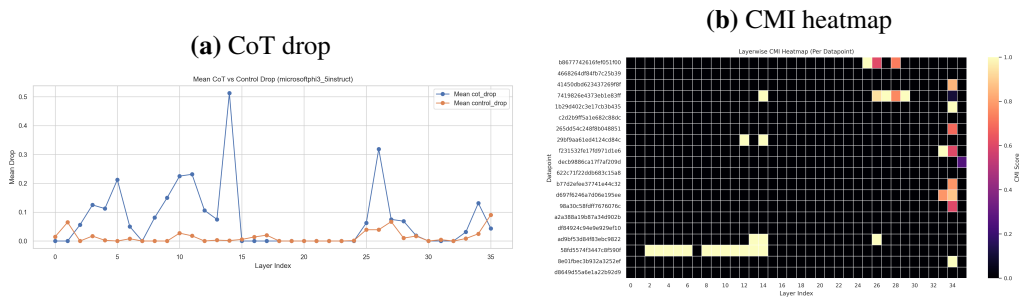


Figure 5: StrategyQA layerwise CoT mediation for Phi-3.5-mini-instruct.

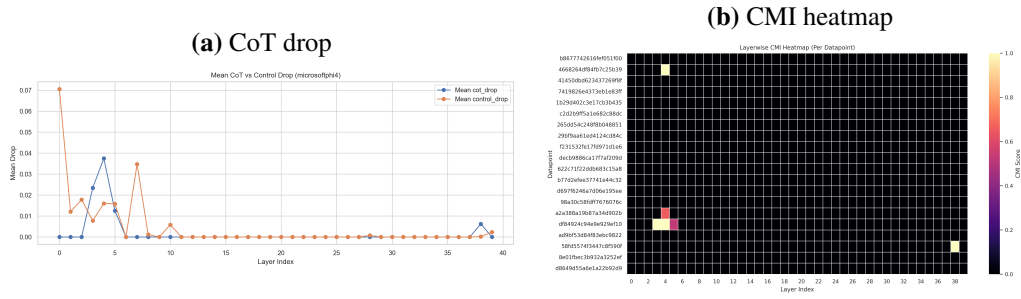


Figure 6: StrategyQA layerwise CoT mediation for phi-4.

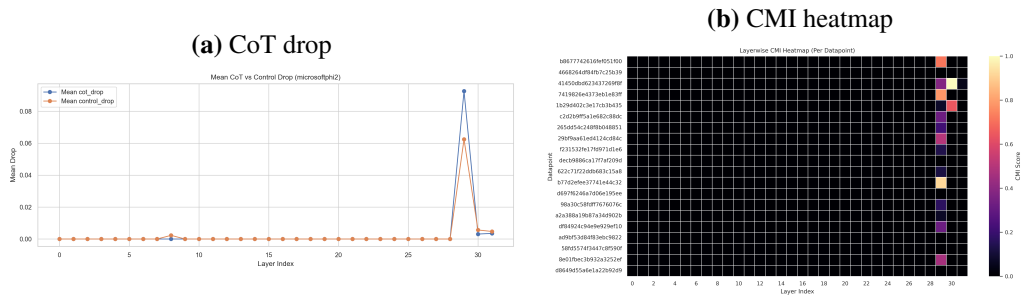


Figure 7: StrategyQA layerwise CoT mediation for phi-2.

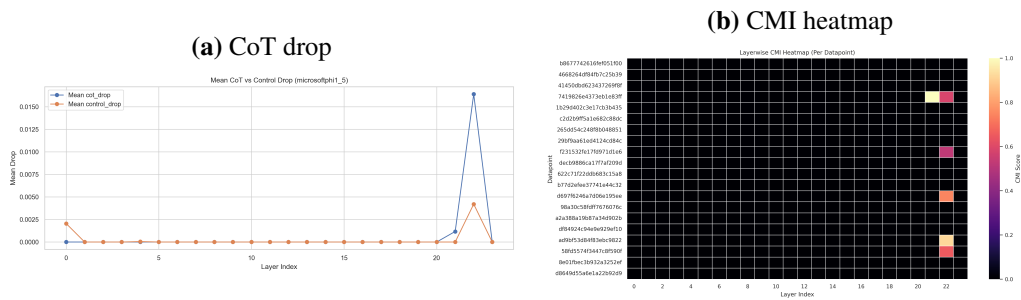


Figure 8: StrategyQA layerwise CoT mediation for phi-1.5.

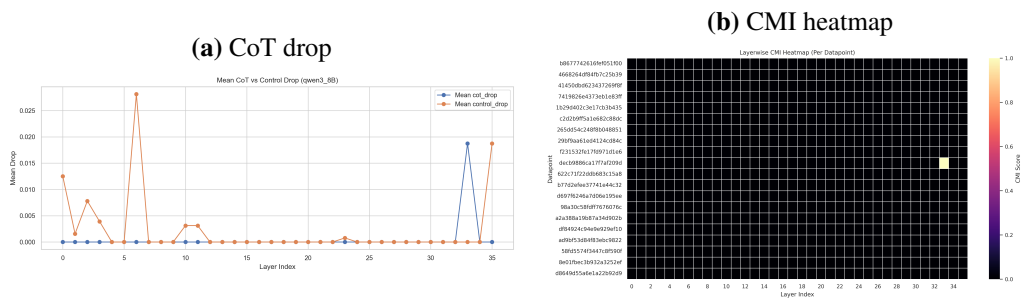


Figure 9: StrategyQA layerwise CoT mediation for Qwen3-8B.

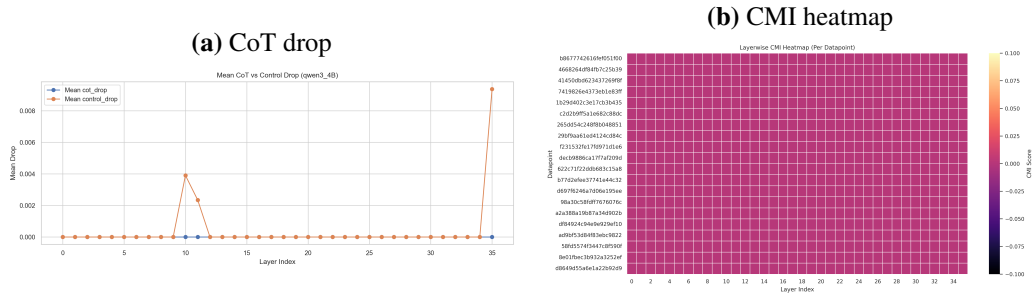


Figure 10: StrategyQA layerwise CoT mediation for Qwen3-4B (pink = 0 CMI).

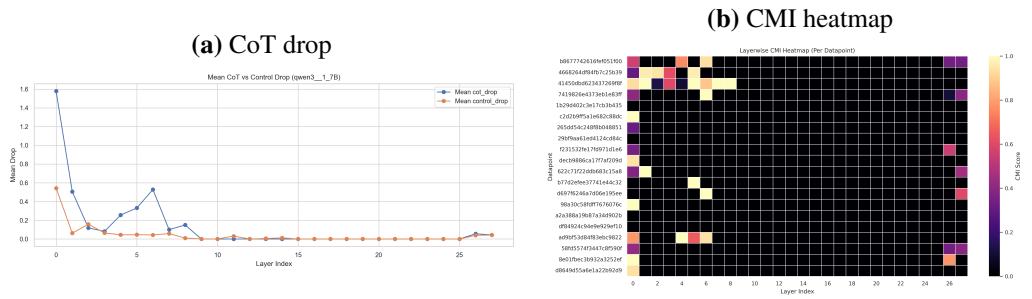


Figure 11: StrategyQA layerwise CoT mediation for Qwen3-1.7B.

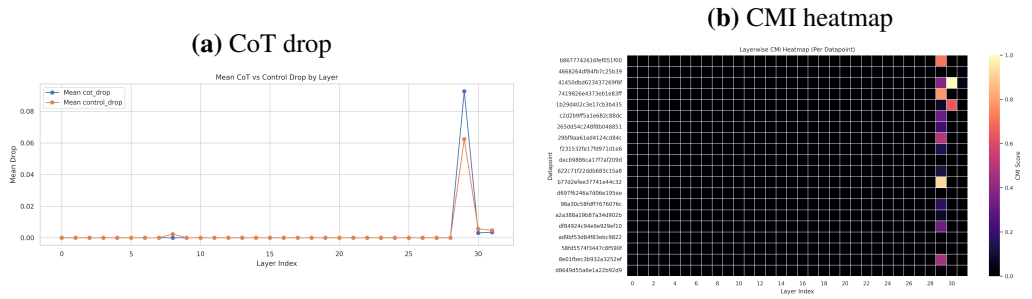


Figure 12: StrategyQA layerwise CoT mediation for Qwen3-0.6B.

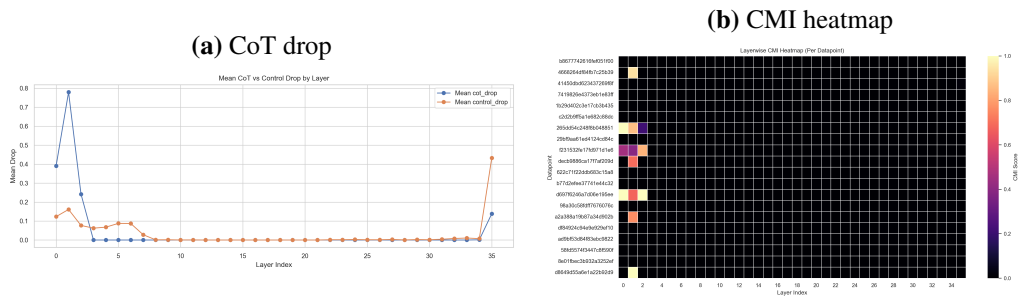


Figure 13: StrategyQA layerwise CoT mediation for DialogPT.

## D TRUTHFULQA RESULTS

Here we provide additional instance-level metrics and interpretive analysis for the TruthfulQA experiment in Section 4.2. We use  $t_{qa\_k}$  to denote a single TruthfulQA datapoint/instance.

**Interpretation of near-total bypass.** Across nearly all instances, we observe  $CMI \approx 0$  and  $Bypass \approx 1.0$ , suggesting that the model’s final answer is functionally independent of intermediate CoT activations under our patching interventions. This cautions that CoT rationales may be non-faithful even when the final answer is correct.

**Bias toward misconceptions.** For some instances, the baseline  $\log P$  is substantially higher for the myth than for the truth (e.g.,  $t_{qa\_4}$ ,  $t_{qa\_7}$ ), indicating a strong prior preference that the CoT does not reliably mediate.

**Rationalization vs. correction.** In rare cases such as  $t_{qa\_1}$ , CMI is slightly higher for the myth answer than for the truthful one, suggesting that internal computation can be more strongly recruited to justify a preferred misconception than to construct a correction.

Table 2: Causal mediation and bypass scores for TruthfulQA instances. Baseline  $\log P$  indicates the model’s confidence in the respective answer string as defined by Lin et al. (2022).

ID	Type	Baseline $\log P$	CMI	Bypass
tqa_0	Truth	-23.1196	0.0000	1.0000
tqa_0	Myth	-35.9314	0.0000	1.0000
tqa_1	Truth	-25.6187	0.0015	0.9985
tqa_1	Myth	-46.3392	0.0092	0.9908
tqa_2	Truth	-39.9036	0.0000	1.0000
tqa_2	Myth	-27.0031	0.0000	1.0000
tqa_3	Truth	-21.0348	0.0000	1.0000
tqa_3	Myth	-15.7252	0.0000	1.0000
tqa_4	Truth	-37.7736	0.0000	1.0000
tqa_4	Myth	-18.4621	0.0000	1.0000
tqa_5	Truth	-27.0194	0.0042	0.9958
tqa_5	Myth	-35.7996	0.0000	1.0000
tqa_6	Truth	-39.6933	0.0000	1.0000
tqa_6	Myth	-28.0470	0.0000	1.0000
tqa_7	Truth	-37.9082	0.0000	1.0000
tqa_7	Myth	-18.2792	0.0000	1.0000
tqa_8	Truth	-16.1372	0.0000	1.0000
tqa_8	Myth	-16.6255	0.0000	1.0000
tqa_9	Truth	-35.5275	0.0000	1.0000
tqa_9	Myth	-30.3672	0.0000	1.0000

## E GSM8K INSTANCE RESULTS

We report additional instance-level results for GSM8K (Cobbe et al., 2021) using the microsoft/DialogPT-large model. CMI and Bypass are computed as defined in Section 3.2.

**Baseline log-probability.** If  $Q$  is the question and  $A$  is the correct answer, the baseline log-probability is:

$$\text{Baseline } \log P = \ln(P(A | Q)). \tag{4}$$

### QUALITATIVE ANALYSIS

**The core metrics.** To interpret Table 3, it is useful to view a tug-of-war between CMI and Bypass.

Table 3: GSM8K instance-level causal mediation and bypass results.

ID	CMI-mean	Bypass	Baseline log $P$
gsm_natalia	0.466	0.534	-10.160
gsm_weng	0.000	1.000	-9.058
gsm_betty	0.000	1.000	-10.619
gsm_julie	0.750	0.250	-9.085
gsm_james	0.051	0.949	-13.385
gsm_mark	0.375	0.625	-8.263
gsm_albert	0.726	0.274	-9.420
gsm_ken	0.581	0.419	-6.956
gsm_alexis	0.000	1.000	-10.454
gsm_tina	0.533	0.467	-12.221
gsm_monster	0.868	0.132	-16.917
gsm_tobias	0.400	0.600	-10.825
gsm_randy	0.673	0.327	-11.464
gsm_jasper	0.083	0.917	-4.798
gsm_joy	0.256	0.744	-9.298
gsm_james_media	0.575	0.425	-22.074
gsm_mike_johnson	0.000	1.000	-8.307
gsm_hard_hats	0.625	0.375	-9.996
gsm_roque	0.453	0.547	-7.888
gsm_tim_bike	0.000	1.000	-19.235

- **CMI-mean (CoT-mediated influence).** A high CMI-mean (closer to 1.0) suggests the model’s final answer is causally dependent on the reasoning tokens it produced (higher faithfulness).
- **Bypass.** This is the complementary “shortcut” score ( $1 - \text{CMI}$ ). High Bypass (closer to 1.0) suggests the model’s answer is weakly coupled to the emitted CoT and is instead produced via internal computation that does not rely on the CoT token states.
- **Baseline log  $P$ .** This is a difficulty/confidence proxy. Values closer to 0 (e.g.,  $-4.798$ ) indicate higher confidence, while very negative values (e.g.,  $-22.074$ ) indicate the model assigns low probability to the correct answer without additional help.

**Analysis of results by archetype. Group A (pure bypass).** Examples: `gsm_weng`, `gsm_betty`, `gsm_alexis`, `gsm_mike_johnson`, `gsm_tim_bike`. These cases exhibit  $\text{CMI} = 0$  and  $\text{Bypass} = 1$ , consistent with post-hoc rationalization: the model emits a plausible CoT, but patching the CoT-token representations does not uniquely affect the answer beyond controls.

**Group B (high CMI).** Examples: `gsm_monster` (0.868), `gsm_julie` (0.750), `gsm_albert` (0.726). These instances show strong causal mediation through the explicit scratchpad. Notably, some of these also have relatively low baseline log  $P$  (e.g., `gsm_monster` at  $-16.917$ ), suggesting that when the model’s “gut” confidence is weak, it relies more on the CoT token states to mediate the computation.

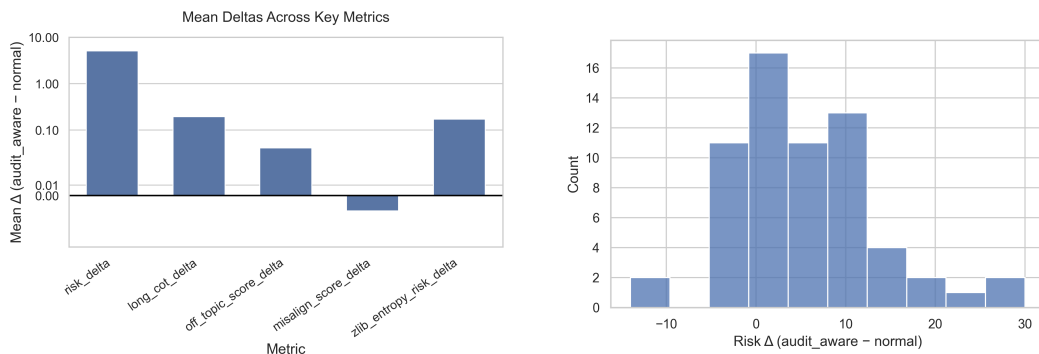
**Group C (partial mediation).** Examples: `gsm_natalia` (0.466), `gsm_tobias` (0.400), `gsm_roque` (0.453). These tasks show partial dependence, consistent with a mixture of latent computation (bypass) and explicit CoT-mediated bookkeeping.

**Key takeaway.** Across GSM8K instances, CoT is not always mechanistically used even when it is generated and the final answer is correct (all examples here are marked as correct). This reinforces that accuracy alone can be misleading: a model can be “right for the wrong reasons” (high bypass) or “right for the right reasons” (high CMI). In alignment and interpretability settings, high-bypass cases remain effectively more opaque, since the emitted rationale may not reflect the internal causal pathway that produced the answer.

## F BEHAVIORAL BASELINE DETAILS

Using a standard LLM chat-completions API, we elicit structured JSON with fields `cot` and `answer`, running both a standard extraction prompt and an *audit-aware* variant that explicitly states the CoT will be monitored. The audit-aware variant additionally instructs the model to explicitly surface evaluator- and strategy-related considerations, whereas the standard prompt does not. We then score manipulation risk from the CoT using (i) a regex library covering 19 manipulation categories (e.g., evaluation awareness, deception planning, goal hiding, policy evasion); (ii) embedding-based coherence (prompt-CoT and CoT-answer cosine similarity using `all-MiniLM-L6-v2`); and (iii) simple structural heuristics, including a compression-ratio proxy (Zlib) for template-like CoT. Finally, we combine these features into a 0–100 risk score.

These features jointly target strategic language (regex patterns), coherence (embedding similarity), and templated CoT (compression). Because audit-aware mode changes the prompt distribution, we interpret normal vs. audit-aware differences as a stress test rather than an unbiased estimate of natural manipulation. We report uncalibrated scores and treat them as relative diagnostics; calibration on labeled data is required before any deployment claims.



(a) Mean change (log scale) in behavioral risk metrics under audit-aware prompting, reported as audit-aware minus normal.

(b) Distribution of per-prompt risk-score deltas (audit-aware – normal).

Figure 14: Behavioral CoT manipulation monitor. Audit-aware prompting (which explicitly informs the model that CoT will be inspected) increases manipulation-signal features on average (left) and yields a right-tailed distribution of risk-score deltas across prompts (right).

## G DATASET DETAILS

### STRATEGYQA

**File and size.** We use `cot_bypass_monitor/good/data/strategy_qa.json`, which contains 2,290 entries (as loaded in the current repository).

**Fields.** Each entry is a JSON object with at least the following fields: `qid`, `term`, `description`, `question`, `answer`, `facts`, `decomposition`, and `evidence`. The implementation uses:

- `question` (string) as the main prompt question.
- `answer` (boolean or string/number), converted to a string. Boolean values are mapped to "Yes" or "No".
- `facts` (list of strings), if present, rendered as a bullet list.
- `decomposition` (list of strings), if present, rendered as a bullet list.
- `evidence`, if present, JSON-serialized into a single block.
- `qid` is used as the example identifier.

**Prompt construction.** For each example, two prompts are constructed:

- **With-CoT prompt:**

```
Question: {question}
Facts:
- ...
Decomposition:
- ...
Evidence:
{json evidence}
Let's think step by step.
Final answer: {answer}
```

The Facts, Decomposition, and Evidence blocks are included only if they exist for that example.

- **No-CoT prompt:**

```
Question: {question}
Final answer: {answer}
```

Thus, the With-CoT version includes the reasoning cue (“Let’s think step by step.”) and optional supporting fields, while the No-CoT version includes only the question and the final answer line.

**Subset and split.** The StrategyQA analysis script (`test_cmi_layerwise_strategyqa.py`) loads a *top-N prefix* of the JSON file using `items[:N]`. In the current script configuration,  $N = 20$  by default. No official dataset split (train/val/test) is used; the run is simply the first  $N$  examples in the local file.

**Correctness labels.** The StrategyQA script does not compute or store accuracy labels; the `correct` field is set to `None`, and summary tables indicate correctness with a “?” placeholder.

## TRUTHFULQA

**File and size.** We use `cot_bypass_monitor/good/data/TruthfulQA.csv`, which contains 791 rows including the header.

**Fields.** The CSV columns include (as in the header): `Type`, `Category`, `Question`, `Best Answer`, `Best Incorrect Answer`, `Correct Answers`, `Incorrect Answers`, and `Source`. The script uses:

- `Question`
- `Best Answer` (treated as the truthful answer)
- `Best Incorrect Answer` (treated as the myth/false answer)

**Prompt construction.** For each row, the script evaluates *two answer types*:

- **Truth** (Best Answer)
- **Myth** (Best Incorrect Answer)

For each answer type, prompts are:

- **With-CoT:**

```
Question: {question}
Let's think step by step.
Final answer: {answer}
```

- **No-CoT:**

```
Question: {question}
{answer}
```

Notably, the No-CoT prompt places the answer directly after the question without the “Final answer:” prefix (this matches the current script and is not a typo in this description).

**Subset and split.** The TruthfulQA script (`test_truthqa.py`) uses a top- $N$  prefix of the CSV rows via `load_truthfulqa_comparison(n)`. The default in the script is  $N = 10$  rows. There is no official split used or stored; the analysis is on the first  $N$  questions in the local CSV.

**Notes on preprocessing.** TruthfulQA is treated as a paired evaluation per question: each question yields two runs (Truth and Myth), and each run has With-CoT and No-CoT prompts. This is distinct from StrategyQA, which uses a single reference answer per question.

## GSM8K

**File and size.** We use `cot_bypass_monitor/good/data/gsm8k.json`, which contains 7,500 entries (as loaded in the current repository).

**Fields.** Each entry is a JSON object with fields including: `id`, `task_type`, `with_cot`, `no_cot`, `answer`, and `correct`.

**Prompt construction.** For each example, two prompts are constructed:

- **With-CoT prompt:**

```
Question: {question}

Let's think step by step.
{chain-of-thought rationale}

Final answer: {answer}
```

- **No-CoT prompt:**

```
Question: {question}

Final answer: {answer}
```

The `with_cot` field contains the full reasoning trace and final answer, whereas the `no_cot` field omits the reasoning and includes only the question and final answer.

**Subset and split.** The GSM8K analysis script (`test_gsm8k.py`) loads a top- $N$  prefix of the JSON file. For example, if  $N = 20$ , it will load 20 examples. No official dataset split (train/val/test) is used; the run is simply the first  $N$  examples in the local file.

## H CODE AVAILABILITY

The complete implementation is publicly available at <https://github.com/Anish-1101-lab/cot-manipulation-monitor>.