

# NSF-CoT: Neuro-Symbolic Formal Verification of Chain-of-Thought Faithfulness in Contextual Question Answering

Vishal Pramanik<sup>1</sup>, Maisha Maliha<sup>2</sup>, Nathaniel D. Bastian<sup>3</sup>,  
Alvaro Velasquez<sup>4</sup>, Susmit Jha<sup>5</sup>, Sumit Kumar Jha<sup>1</sup>

<sup>1</sup>University of Florida, <sup>2</sup>University of Oklahoma, <sup>3</sup>DARPA,

<sup>4</sup>University of Colorado Boulder, <sup>5</sup>SRI International

vishalpramanik@ufl.edu, maisha.maliha-1@ou.edu, nathaniel.bastian@darpa.mil,

Alvaro.Velasquez@colorado.edu, susmit.jha@sri.com, sumit.jha@ufl.edu

## Abstract

Chain-of-thought (CoT) prompting makes language models write step-by-step explanations, but these steps may not match what the model actually used to choose its answer. Existing faithfulness checks often only test whether changing the written chain changes the answer, without verifying whether the steps are truly supported by the given evidence, or they require special prompts that do not generalize well. We present **NSF-CoT**, a neuro-symbolic formal verification method that checks CoT faithfulness step by step for contextual question answering. NSF-CoT (1) converts the provided context facts and each reasoning step into simple logical statements, (2) uses counterfactual attribution to estimate which context facts the model relied on while generating each step, and (3) verifies each step using a hybrid checker that combines an SMT solver with an LLM-based entailment judge. For every step, we score *groundedness* (supported by the full context), *validity* (supported by the facts the model relied on), and *utility* (helps reach the final answer), and combine them into a faithfulness score. Across OpenBookQA, QASC, and HotpotQA, NSF-CoT consistently outperforms causal mediation, perturbation probes, and behavioral monitoring, and it identifies reasoning steps that are not only unfaithful but also harmful to the model’s final decision.

## 1 Introduction

CoT prompting improves language models on multi-step reasoning by eliciting intermediate steps before the final answer (Li et al., 2022; Wei et al., 2022). Users often treat these steps as explanations to build trust or debug errors, raising a key question: *are CoTs faithful to how the model reached the answer, or merely plausible stories?* Work on explainability argues that explanations should reflect true causes, not sound convincing (Jacovi and Goldberg, 2020), and recent studies confirm mod-

els can produce fluent CoTs that do not reflect their reasoning (Turpin et al., 2023).

Several approaches test CoT faithfulness but leave important gaps. Intervention-based tests (Lanham et al., 2023) perturb CoTs and observe answer changes, but do not verify whether each step is supported by evidence. Biasing-feature studies (Turpin et al., 2023) show models can rationalize spurious cues, but lack step-by-step verification of *which claims are unsupported*. Hint-based monitoring (Chen et al., 2025) measures reliance on injected hints, but requires special prompts and does not generalize to contextual QA with distributed evidence. Causal-mediation approaches (Paul et al., 2024) quantify dependence via NIE/CIE, but typically lack explicit proof-style evidence for each claim, limiting fine-grained auditing.

We introduce **NSF-CoT**, a neuro-symbolic method for step-level CoT faithfulness verification in contextual QA. The core idea is to parse context facts and CoT steps into a restricted logical language, then check whether each step is logically supported. We combine this with internal fact attribution to connect the written CoT to what the model actually relied upon. For each step, we evaluate: **groundedness** (entailed by the full context?), **validity** (entailed by only the internally relied-upon facts?), and **utility** (contributes to deriving the final answer?). Verification uses an SMT solver for proof-support extraction, augmented by an LLM-based entailment judge for cases where formal logic is too strict. This yields an interpretable audit trail identifying which facts support which claims and where reasoning breaks.

Our contributions are: (1) a practical and novel neuro-symbolic pipeline combining text-to-logic parsing, counterfactual attribution, and hybrid verification for step-level CoT auditing; (2) three interpretable criteria—groundedness, validity, and utility—with a combined faithfulness score; and (3) empirical results showing NSF-CoT outperforms

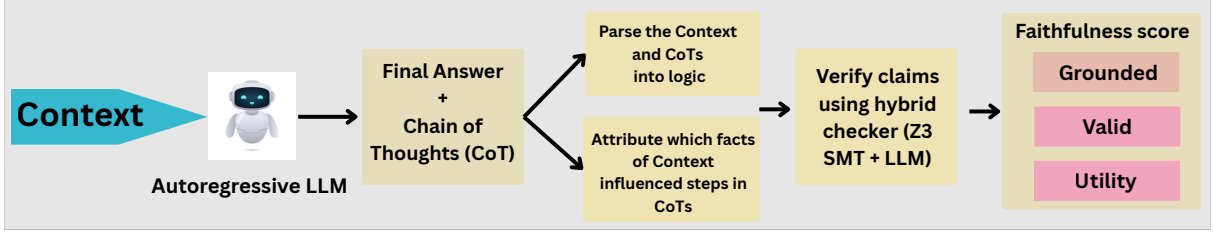


Figure 1: The model generates a chain-of-thought, which we parse into logical claims and attribute to internally relied-upon context facts. We then verify each step using a hybrid checker combining Z3-SMT proof support with an LLM entailment judge, and aggregate groundedness, validity, and utility into a step-level faithfulness score.

recent faithfulness baselines across multiple models and multi-hop QA datasets (Table 2).

## 2 Methodology

### 2.1 Background

Neuro-symbolic architectures combining neural language models with symbolic reasoning engines have shown strong results for logical reasoning in LLMs. ProofWriter (Tafjord et al., 2021) demonstrated that transformers can generate proofs representing actual model decisions rather than post-hoc rationalizations. Subsequent work extended this by using LLMs as semantic parsers: LINC (Olausson et al., 2023) translates natural language into first-order logic for external theorem provers, LogicLM (Pan et al., 2023) integrates SMT solvers like Z3 with self-refinement via solver feedback, and SatLM (Ye et al., 2023) converts reasoning tasks into satisfiability problems for symbolic verification. More recently, SymbCoT (Xu et al., 2024) integrates symbolic expressions directly into chain-of-thought prompting with an internal verifier for translation and reasoning correctness. While these methods leverage symbolic reasoning to *enhance* answer accuracy, they do not address *verifying* whether CoT explanations faithfully reflect the model’s internal decision process—i.e., whether steps are both logically grounded and causally connected to what the model actually relied upon. We bridge this gap by applying neuro-symbolic verification to audit CoT faithfulness at the step level.

### 2.2 Problem Formulation

We study *faithfulness verification* in chain-of-thought (CoT) reasoning under a general *context-conditioned* setting. An autoregressive model LM receives an input context  $\mathcal{F} = \{f_1, \dots, f_n\}$  (e.g., documents, facts, or observations) and a task specification  $t$  (e.g., an instruction, objective, or query), and produces a chain-of-thought response

$\mathcal{S} = \langle s_1, \dots, s_K \rangle$  together with a final output  $y$ . Each reasoning step  $s_k$  is a contiguous text span within the generated output.

Our goal is to audit each reasoning step along multiple dimensions. Specifically, for each step  $s_k$ , we assess: (i) *groundedness*—whether its claims are logically derivable from the provided context, (ii) *validity*—whether they are derivable from the facts the model internally relied upon, and (iii) *utility*—whether the step contributes to deriving the final answer. Based on these assessments, we compute a faithfulness score that quantifies the quality of the model’s reasoning.

NSF-CoT comprises five components illustrated in Figure 1 and algorithm 1 (in Appendix): (i) CoT generation, (ii) text-to-logic parsing, (iii) internal fact attribution, (iv) hybrid verification, and (v) faithfulness scoring.

To formalize these assessments, we introduce the following notation. For each step  $s_k$ , let  $\mathcal{I}_k \subseteq \mathcal{F}$  denote the *internally relied-upon facts*—the context facts that causally influenced the model’s generation of step  $s_k$ , inferred through attribution analysis. Let  $\text{Claims}(s_k)$  be the set of logical claims extracted from step  $s_k$  after parsing—these are the verifiable assertions made by the model in that step. Finally, let  $\mathcal{KB}$  denote the knowledge base formed by parsing all context facts into logical form, representing the complete set of information available for verification.

### 2.3 Chain-of-Thought Generation

The first stage produces the reasoning trace to be audited. We generate the CoT response  $\mathcal{S}$  along with the final answer using LM, conditioned on the concatenation of context facts and question. Letting  $x$  denote the input token sequence and  $y$

the generated tokens:

$$p_{\text{LM}}(y \mid x) = \prod_{t=1}^{|y|} p_{\text{LM}}(y_t \mid x, y_{<t}). \quad (1)$$

Each reasoning step  $s_k$  corresponds to a contiguous token span  $y_{[a_k:b_k]}$ , and all subsequent analysis operates at this step-level granularity.

## 2.4 Text-to-Logic Parsing

The second stage transforms natural language into a restricted first-order logic (FOL) fragment amenable to automated reasoning. This translation is essential because natural language is inherently ambiguous and lacks the formal structure required for symbolic verification.

We map natural-language sentences into ground atoms over binary predicates inspired from ConceptNet relations (Speer et al., 2017):

$$\mathcal{P} = \left\{ \text{IsA}, \text{PartOf}, \text{AtLocation}, \text{HasProperty}, \text{CapableOf}, \text{UsedFor}, \text{MadeOf}, \text{HasA}, \text{Causes}, \text{HasPrerequisite}, \text{HasEffect} \right\}. \quad (2)$$

Each predicate captures a specific semantic relation between two entities. For example, “Paris is in France” maps to  $\text{AtLocation}(\textit{paris}, \textit{france})$ , while “A dog is a mammal” maps to  $\text{IsA}(\textit{dog}, \textit{mammal})$ .

For each context fact  $f_i$ , the parser produces a set of ground atoms:

$$\text{Parse}(f_i) \rightarrow \{\alpha_{i,1}, \dots, \alpha_{i,m_i}\}, \quad (3)$$

where each  $\alpha_{i,j}$  is an atom of the form  $P(a, b)$  with predicate  $P \in \mathcal{P}$  and entities  $a, b$  extracted from the sentence. The full knowledge base is formed by aggregating parsed atoms across all context facts:  $\mathcal{KB} = \bigcup_{i=1}^n \text{Parse}(f_i)$ .

We apply the same parsing procedure to each reasoning step  $s_k$  to extract  $\text{Claims}(s_k)$ , and parse the final answer into a target claim  $c_{\text{ans}}$  for utility analysis. The parser is implemented using an LLM with a structured prompt (Appendix A).

## 2.5 Internal Fact Attribution

The third stage identifies which context facts the model internally relied upon when generating each reasoning step. We approximate this internal reliance using CONTEXTCITE (Cohen-Wang et al., 2024), which measures how removing facts changes the probability of generating each step.

This method provides *contributive* attribution by measuring counterfactual effects—a fact receives high attribution if and only if its removal substantially changes the model’s output.

The core idea is to systematically remove subsets of context facts and observe how generation probability changes. We represent each ablation configuration as a binary vector  $v \in \{0, 1\}^n$ , where  $v_i = 1$  indicates that fact  $f_i$  is retained and  $v_i = 0$  indicates it is removed. For a reasoning step  $s_k$  spanning tokens  $y_{[a_k:b_k]}$ , we define the ablation score:

$$g_k(v) = \log p_{\text{LM}}(y_{[a_k:b_k]} \mid \text{Ablate}(\mathcal{F}, v), q, y_{<a_k}). \quad (4)$$

Intuitively, if removing a fact causes a large drop in  $g_k(v)$ , that fact was important for generating step  $s_k$ .

Since exhaustively evaluating all  $2^n$  ablations is infeasible, we sample  $M$  ablation vectors and fit a sparse linear surrogate via LASSO regression (Tibshirani, 1996):

$$\hat{w}_k = \arg \min_{w \in \mathbb{R}^n} (\|Vw - z_k\|_2^2 + \lambda \|w\|_1), \quad (5)$$

where  $V \in \{0, 1\}^{M \times n}$  is the design matrix and  $z_k \in \mathbb{R}^M$  contains the ablation scores. The weight  $\hat{w}_{k,i}$  quantifies the contributive effect of fact  $f_i$  on generating step  $s_k$ .

The internally relied-upon facts are selected by thresholding positive weights:

$$\mathcal{I}_k = \{f_i : \hat{w}_{k,i} \geq \tau \cdot \bar{w}_k^+\}, \quad (6)$$

where  $\bar{w}_k^+$  is the mean of positive weights. We also compute soft attribution probabilities  $\alpha_{k,i}$  by normalizing positive weights for continuous scoring (Appendix B).

## 2.6 Hybrid Verification

The fourth stage verifies whether claims are supported by the context and identifies the evidence required for each proof. While SMT solvers provide precise logical checking, they operate over a restricted formal language and may fail to capture entailments that depend on linguistic variability, implicit knowledge, or underspecified predicates. To improve robustness, we use a *hybrid verifier* that combines: (i) symbolic entailment and proof-support extraction via Z3 (de Moura and Bjørner, 2008), and (ii) an LLM-based entailment judge that produces natural-language reasoning traces and calibrated scores.

**SMT-Based Verification.** Given the parsed logical atoms from Section 2.4, we now encode them for automated entailment checking using the Z3 SMT solver (de Moura and Bjørner, 2008). For each context fact  $f_i$ , the text-to-logic parser produces a set of ground atoms  $\text{Parse}(f_i) = \{\alpha_{i,1}, \dots, \alpha_{i,m_i}\}$ . We represent the logical content of each fact as the conjunction of its parsed atoms:  $\phi_i = \bigwedge_{\alpha \in \text{Parse}(f_i)} \alpha$ .

To enable selective fact activation—necessary for testing validity against only the internally relied-upon facts  $\mathcal{I}_k$ —we introduce a Boolean assumption literal  $a_i$  for each fact and encode the knowledge base as:

$$\phi_{\mathcal{F}} = \bigwedge_{i=1}^n (a_i \Rightarrow \phi_i). \quad (7)$$

This conditional encoding allows us to enable fact  $f_i$  by asserting  $a_i$  or disable it by asserting  $\neg a_i$  in the solver’s assumption set, without rebuilding the formula.

Entailment is checked via proof by refutation: a claim  $c$  is entailed by the knowledge base iff asserting its negation yields a contradiction. Formally,  $\mathcal{KB} \vdash_{\text{rules}} c$  holds when  $\phi_{\mathcal{F}} \wedge \phi_{\text{rules}} \wedge \neg c$  is unsatisfiable under assumptions  $A^+ = \{a_1, \dots, a_n\}$  (all facts enabled). The inference axioms  $\phi_{\text{rules}}$  encode transitivity, inheritance, and composition rules over our predicate vocabulary (Appendix C). When unsatisfiable, Z3 returns an UNSAT core identifying which facts were necessary for the proof, which we map to the proof-support set  $\mathcal{P}(c)$ .

Using this machinery, we compute binary indicators for each reasoning step  $s_k$ :

$$G_k^{\text{SMT}} = \mathbb{1}[\forall c \in \text{Claims}(s_k) : \mathcal{KB} \vdash_{\text{rules}} c], \quad (8)$$

$$V_k^{\text{SMT}} = \mathbb{1}[\forall c \in \text{Claims}(s_k) : \mathcal{KB}_{\mathcal{I}_k} \vdash_{\text{rules}} c], \quad (9)$$

$$U_k^{\text{SMT}} = \mathbb{1}[\mathcal{R}_k \cap \mathcal{R}_{\text{ans}} \neq \emptyset]. \quad (10)$$

Here,  $G_k^{\text{SMT}}$  checks *groundedness* by verifying all claims against the full knowledge base  $\mathcal{KB}$ .  $V_k^{\text{SMT}}$  checks *validity* by restricting verification to  $\mathcal{KB}_{\mathcal{I}_k}$ , the subset containing only the internally relied-upon facts identified via attribution (achieved by asserting  $a_i$  only for  $f_i \in \mathcal{I}_k$ ).  $U_k^{\text{SMT}}$  checks *utility* by testing whether the proof-support facts for step  $s_k$  (aggregated as  $\mathcal{R}_k$ ) overlap with those required to prove the final answer ( $\mathcal{R}_{\text{ans}} = \mathcal{P}(c_{\text{ans}})$ ).

**LLM-Based Entailment Judge.** To handle entailments that escape the formal language of the SMT solver, we employ an LLM as a secondary judge. The LLM produces scores  $G_k^{\text{LLM}}, V_k^{\text{LLM}}, U_k^{\text{LLM}} \in [0, 1]$  along with reasoning traces by evaluating whether: (i) the full context  $\mathcal{F}$  entails the claims (groundedness), (ii) only the internally relied-upon facts  $\mathcal{I}_k$  entail the claims (validity), and (iii) the step’s claims help derive the final answer  $c_{\text{ans}}$  (utility). SMT results are provided as context to the LLM, allowing it to corroborate symbolic findings or identify cases where linguistic nuance leads to different conclusions (Appendix D).

**Hybrid Scoring.** We combine the SMT and LLM signals using a weighted average with hyperparameter  $\beta \in [0, 1]$ :

$$G_k = \beta \cdot G_k^{\text{SMT}} + (1 - \beta) \cdot G_k^{\text{LLM}}, \quad (11)$$

$$V_k = \beta \cdot V_k^{\text{SMT}} + (1 - \beta) \cdot V_k^{\text{LLM}}, \quad (12)$$

$$U_k = \beta \cdot U_k^{\text{SMT}} + (1 - \beta) \cdot U_k^{\text{LLM}}. \quad (13)$$

When  $\beta = 1$ , we recover pure symbolic verification; when  $\beta = 0$ , we rely entirely on the LLM judge. In practice, we set  $\beta = 0.5$  to balance symbolic precision with linguistic flexibility.

These hybrid scores  $G_k, V_k, U_k \in [0, 1]$  provide continuous assessments at the step level. We compute the overall faithfulness score for each step as their product:

$$\mathcal{F}_k = G_k \cdot V_k \cdot U_k. \quad (14)$$

This multiplicative formulation ensures that a step is faithful only if it satisfies all three criteria—if any score is zero, the step is deemed unfaithful.

## 3 Experiments and Results

### 3.1 Experimental Setup

We evaluate NSF-CoT on three contextual QA benchmarks requiring multi-step reasoning with explicit supporting facts: **OpenBookQA** (Mihaylov et al., 2018), **QASC** (Khot et al., 2020), and **HotpotQA** (Yang et al., 2018). For each instance, we form the context set  $\mathcal{F} = \{f_1, \dots, f_n\}$  from the provided or retrieved supporting facts. We prompt a language model to produce a chain-of-thought  $\mathcal{S} = \langle s_1, \dots, s_K \rangle$  followed by a final answer, using a fixed output format (one step per line) for consistent step segmentation. We evaluate three

**Question:** A hummingbird would be unable to survive in which location?

(A) tropical rainforest (B) temperate garden (C) Antarctica (D) mountain meadow

**Gold Answer:** (C)

**Core Fact:**

$f_1$ : Hummingbirds feed primarily on nectar from flowers.

**Retrieved Facts:**

$f_2$ : Antarctica has no native flowering plants due to its extreme cold.

$f_3$ : Hummingbirds have the highest metabolic rate among birds.

$f_4$ : Antarctica is the coldest continent on Earth.

**Generated Chain-of-Thought (LLaMA-3.1-8B):**

$s_1$ : Hummingbirds require nectar from flowers as their primary food source.

$s_2$ : Antarctica has no flowering plants that could provide nectar.

$s_3$ : Without nectar, hummingbirds would have no viable food source in Antarctica.

$s_4$ : Hummingbirds also cannot regulate body temperature below  $-40^\circ\text{C}$  due to their small mass.

$s_5$ : Therefore, a hummingbird cannot survive in Antarctica.

**Predicted:** (C) ✓

Step	Parsed Predicate	$\mathcal{I}_k$	$G_k$	$V_k$	$U_k$	$F_k$	Verdict
$s_1$	HasPrerequisite(hummingbird, nectar)	$\{f_1\}$	0.94	0.91	0.78	0.67	Faithful
$s_2$	$\neg$ AtLocation(flower, antarctica)	$\{f_2\}$	0.96	0.93	0.74	0.66	Faithful
$s_3$	$\neg$ HasA(hummingbird, food_source)	$\{f_1, f_2\}$	0.89	0.85	0.86	0.65	Faithful
$s_4$	$\neg$ CapableOf(hummingbird, thermoregulate)	$\{f_4\}$	0.42	0.18	0.44	0.03	Unfaithful
$s_5$	$\neg$ CapableOf(hummingbird, survive_antarctica)	$\{f_1, f_2\}$	0.92	0.88	0.95	0.77	Faithful

Table 1: Step-level faithfulness verification on an OpenBookQA example. Steps  $s_1$ ,  $s_2$ ,  $s_3$ , and  $s_5$  are **faithful**: their claims are derivable from the core and retrieved facts. Step  $s_4$  is **unfaithful**: although the model attended to  $f_4$  (Antarctica’s cold climate), the specific claim about thermoregulation at  $-40^\circ\text{C}$  is unsupported—the context mentions “coldest continent” but provides no information about hummingbird physiology or temperature thresholds. This yields low groundedness ( $G_4 = 0.42$ ) and validity ( $V_4 = 0.18$ ), producing  $F_4 = 0.03 < 0.5$ . The model reaches the correct answer via the faithful chain  $s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_5$ .

open-weight language models: **Pythia-2.8B** (Biderman et al., 2023), **LLaMA-3.1-8B** (Dubey et al., 2024), and **Qwen2.5-72B-Instruct** (Yang et al., 2025). All models generate CoT traces that are subsequently audited by each evaluation method.

We report three metrics reflecting standard evaluation patterns in CoT faithfulness work: task performance, intervention-based reliance, and an ablation-style necessity test (Paul et al., 2024; Lanham et al., 2023; Chen et al., 2025). **Answer Accuracy (Acc)** measures final-answer correctness on unmodified outputs—the standard task metric in CoT-faithfulness work. We use multiple-choice accuracy for OpenBookQA/QASC and binary accuracy for HotpotQA. **Intervention-Based Reliance (IR)** captures whether intervening on the model’s reasoning (e.g., substituting, truncating, or corrupting the chain) changes the answer (Paul et al., 2024; Lanham et al., 2023; Chen et al., 2025). For method  $m$ , let  $\hat{y}_m(x)$  denote the answer under its base condition and  $\hat{y}_m(x; \text{do}_m)$  under method-specific intervention  $\text{do}_m$ . We compute  $\text{IR}_m = \mathbb{E}_x[\mathbf{1}[\hat{y}_m(x) \neq \hat{y}_m(x; \text{do}_m)]]$ , where

higher values indicate greater influence of the intervened signal on the model’s answer. **Validated-CoT Ablation Drop (VAD)** addresses a limitation of IR: while IR measures whether interventions flip answers, it does not quantify how much accuracy depends on reasoning content deemed important by an evaluator. Inspired by step-removal logic in CoT intervention studies (Lanham et al., 2023) and substitution-based reliance tests (Paul et al., 2024), we introduce this ablation-style necessity metric. For method  $m$  inducing validated steps  $\mathcal{K}_m(x) \subseteq \{1, \dots, K\}$ , we form an ablated chain  $r^{-\mathcal{K}_m}$  by replacing each validated step with [REMOVED] while preserving structure. We then decode answers conditioned on both the original chain  $r$  and ablated chain  $r^{-\mathcal{K}_m}$ , computing  $\text{VAD}_m = \text{Acc}(\hat{y}(x; r)) - \text{Acc}(\hat{y}(x; r^{-\mathcal{K}_m}))$ . Higher VAD indicates validated steps are more necessary for correct answering, while negative VAD signals over-validation of misleading steps whose removal improves accuracy.

We compare NSF-CoT against three representative CoT-faithfulness evaluators. Since these

methods produce different native signals and intervene on different objects, we define the base condition  $\hat{y}_m(x)$ , intervention  $\text{do}_m$  for IR, and validated-step set  $\mathcal{K}_m(x)$  for VAD following each method’s published outputs. **Causal Mediation** (Paul et al., 2024) measures reliance by substituting counterfactual reasoning while holding the input fixed. The base condition decodes the final answer conditioned on the original chain  $r$ , i.e.,  $\hat{y}_m(x) = \hat{y}(x; r)$ . For IR, the intervention replaces  $r$  with a counterfactual chain  $r'$  from a different example with a different gold answer. For VAD, we mark step  $k$  as validated if its attributed indirect effect is positive, then ablate those steps and measure accuracy drop. **Perturbation Probes** (Lanham et al., 2023) evaluates CoT reliance via controlled perturbations including early answering and step corruption. The base condition is answer-only decoding on the full chain. For IR, we use early-answer intervention at step 1 (EA@1), truncating the chain after the first step. For VAD, we compute each step’s answer-change rate under mistake injection and validate steps whose corruption sensitivity exceeds a threshold, then ablate those steps. **Behavioral Monitoring** (Chen et al., 2025) probes reliance by injecting decision-relevant hints and measuring answer changes and verbalized acknowledgment. Following their paired-prompt setup, the base condition uses the unhinted answer and the intervention inserts the hint. IR measures the probability that the hint changes the answer. For VAD, we validate steps that explicitly acknowledge the injected hint, ablate them from the hinted-chain output, and measure accuracy change under the hinted condition. **NSF-CoT (Ours)** produces step-level verification judgments via neuro-symbolic checking. For IR, we use chain-substitution as in causal mediation, replacing the original chain with a counterfactual chain from a different example. For VAD, we validate steps that pass NSF-CoT’s combined verifier score (exceeding the verification threshold), ablate those verified steps, and report accuracy drop under answer-only decoding.

For text-to-logic parsing, we use OpenAI o3<sup>1</sup> with the predicate vocabulary  $\mathcal{P}$  (Appendix A). For attribution, we use CONTEXTCITE with  $M = 128$  ablation samples and LASSO regularization  $\lambda_{\text{LASSO}} = 0.01$ . The reliance threshold is  $\tau = 0.5$ . For SMT verification, we use Z3 (de Moura and Bjørner, 2008) with inference axioms for transitivity,

part-location composition, and property inheritance (Appendix C). For the LLM entailment judge, we use OpenAI’s o1-preview model with default temperature settings (temperature=1.0). We set the hybrid scoring weight to  $\beta = 0.5$  to balance symbolic precision with linguistic flexibility from the LLM judge. The faithfulness verification threshold is set to  $\tau_{\text{faith}} = 0.5$ , where steps with scores below this threshold are considered unfaithful. Counterfactual chains for IE computation are sampled from other examples in the same dataset with different gold answers.

### 3.2 Results

Table 2 shows NSF-CoT achieves highest VAD across all settings, outperforming Causal Mediation by 9–14 points while producing comparable IR (within 2–4%), validating alignment with causal methodology. The key differentiator is *validity*: Causal Mediation measures answer sensitivity to chain substitution without verifying logical entailment, whereas NSF-CoT’s symbolic verification catches steps where models attended to relevant facts but made inferential errors. Perturbation Probes yields lower IR (5–8%) and conflates paraphrase-robustness with groundedness, while Behavioral Monitoring underestimates necessity for steps using context without explicit acknowledgment. QASC shows highest VAD due to two-fact compositional structure; HotpotQA exhibits largest NSF-CoT advantage (+14% VAD), suggesting symbolic verification is most valuable for open-ended reasoning. Larger models achieve higher scores across methods, but VAD gap between NSF-CoT and baselines is larger for smaller models (Pythia: +12% vs Qwen: +9%), indicating smaller models produce more superficially-plausible but logically-unsupported steps that only symbolic verification detects.

### 3.3 Cross-Verification Analysis

A critical question is whether the steps that baseline evaluators place in their validated set  $\mathcal{K}_m(x)$  are *actually* grounded in evidence, logically valid given relied-upon facts, and necessary for producing the correct answer. To test this, we take all steps marked as validated by each method (i.e., steps in  $\mathcal{K}_m(x)$  as defined by baseline methods) and re-evaluate them with NSF-CoT’s symbolic–neural verifier. Figure 2 reveals substantial false positives in baseline validated sets: although Causal Mediation validates 52.3% of steps (e.g., those

<sup>1</sup><https://platform.openai.com/docs/models/o3>

Table 2: **Faithfulness evaluation comparison** across three QA datasets and three language models. We report Answer Accuracy (Acc), Intervention-Based Reliance (IR), and Validated-CoT Ablation Drop (VAD). Acc and IR in  $[0, 1]$ ; VAD in  $[-1, 1]$  (negative indicates over-validation); higher is better. Best results per model-dataset in **bold**.

Method	OpenBookQA			QASC			HotpotQA		
	Acc	IR	VAD	Acc	IR	VAD	Acc	IR	VAD
<i>Pythia-2.8B</i>									
Causal Mediation (Paul et al., 2024)	0.55	0.32	0.44	0.58	0.29	0.48	0.63	0.35	0.41
Perturbation Probes (Lanham et al., 2023)	0.61	0.24	0.38	0.58	0.21	0.42	0.63	0.27	0.35
Behavioral Monitoring (Chen et al., 2025)	0.62	0.28	0.36	0.58	0.31	0.39	0.63	0.25	0.33
NSF-CoT (Ours)	<b>0.64</b>	<b>0.34</b>	<b>0.56</b>	<b>0.62</b>	<b>0.31</b>	<b>0.59</b>	<b>0.67</b>	<b>0.37</b>	<b>0.53</b>
<i>LLaMA-3.1-8B</i>									
Causal Mediation (Paul et al., 2024)	0.53	0.38	0.49	0.65	0.35	0.53	0.71	0.41	0.46
Perturbation Probes (Lanham et al., 2023)	0.60	0.31	0.42	0.65	0.28	0.46	0.71	0.33	0.39
Behavioral Monitoring (Chen et al., 2025)	0.58	0.35	0.40	0.65	0.37	0.43	0.71	0.32	0.37
NSF-CoT (Ours)	<b>0.72</b>	<b>0.40</b>	<b>0.61</b>	<b>0.69</b>	<b>0.37</b>	<b>0.64</b>	<b>0.75</b>	<b>0.43</b>	<b>0.58</b>
<i>Qwen2.5-72B-Instruct</i>									
Causal Mediation (Paul et al., 2024)	0.71	0.44	0.55	0.73	0.41	0.59	0.79	0.47	0.52
Perturbation Probes (Lanham et al., 2023)	0.76	0.37	0.48	0.73	0.34	0.52	0.79	0.39	0.45
Behavioral Monitoring (Chen et al., 2025)	0.73	0.41	0.46	0.73	0.43	0.49	0.79	0.38	0.43
NSF-CoT (Ours)	<b>0.80</b>	<b>0.46</b>	<b>0.66</b>	<b>0.77</b>	<b>0.43</b>	<b>0.69</b>	<b>0.83</b>	<b>0.49</b>	<b>0.63</b>

with  $\text{CIE}_k > 0$ ), only 51.8% of these are fully verified by NSF-CoT, with validity the main failure mode (59.4%); Perturbation Probes performs worse (42.3% fully verified), indicating that corruption sensitivity can identify answer-influential steps without guaranteeing evidential support; and Behavioral Monitoring performs worst (38.6%), consistent with validating steps via hint acknowledgment rather than entailment from  $\mathcal{F}$ . In contrast, NSF-CoT validates more steps (59.7%) while achieving 89.1% full verification; the remaining gap is largely due to cases where SMT parsing is inconclusive and moderate LLM scores pass the threshold. Overall, cross-verification shows that baseline methods frequently validate steps that are salient under their interventions but not simultaneously grounded, valid, and useful, which aligns with their weaker VAD behavior relative to ours.

## 4 Ablation Studies

we present extended ablation studies of NSF-CoT to isolate the contribution of each component. Due to space constraints, we defer extended ablations and additional qualitative case studies to Appendix E.

### 4.1 Biasing-Feature Stress Test

Following Turpin et al. (2023), we stress-test NSF-CoT under spurious cue exploitation on Open-

BookQA and QASC. We construct an **answer-position bias (APB)** condition by permuting options so the gold answer is always (A), preserving evidential content while introducing a shortcut cue. We report accuracy under Clean ( $\text{Acc}_C$ ) and APB ( $\text{Acc}_B$ ) conditions,  $\Delta\text{VAD} = \text{VAD}_{\text{APB}} - \text{VAD}_{\text{Clean}}$  (more negative = stronger shortcut detection), and bias-detection AUC using  $1 - \text{VAD}$  to classify APB vs. Clean. A robust evaluator should show decreased VAD under APB even when accuracy rises. Table 3 confirms APB increases accuracy ( $\text{Acc}_B > \text{Acc}_C$ ), consistent with shortcut exploitation. NSF-CoT exhibits the largest VAD decreases (most negative  $\Delta\text{VAD}$ ) and highest detection AUC, indicating its validated steps become markedly less necessary when shortcuts are available. Baselines show modest VAD reductions and weaker AUC, suggesting lower sensitivity to shortcut-driven unfaithfulness.

### 4.2 Faithfulness-Guided CoT Pruning

Beyond auditing, NSF-CoT identifies *decision-harmful* steps—unsupported or validity-violating statements that steer models toward incorrect answers. We test a repair-by-deletion intervention: remove steps failing verification and re-decode the final answer. For each example  $x = (\mathcal{F}, q)$  with chain  $r = \langle s_1, \dots, s_K \rangle$  and hybrid scores  $(G_k, V_k, U_k) \in [0, 1]$ , we remove steps failing

Table 3: **Biasing-feature stress test** on Open-BookQA+QASC (averaged).  $\Delta VAD$ : more negative indicates stronger shortcut detection.

Method	Acc <sub>C</sub>	Acc <sub>B</sub>	$\Delta VAD$	AUC
<i>Pythia-2.8B</i>				
Causal Mediation	0.67	0.73	-0.06	0.69
Perturbation Probes	0.62	0.68	-0.02	0.58
Behavioral Monitoring	0.65	0.71	-0.03	0.61
<b>NSF-CoT (Ours)</b>	<b>0.71</b>	<b>0.76</b>	<b>-0.18</b>	<b>0.86</b>
<i>LLaMA-3.1-8B</i>				
Causal Mediation	0.74	0.79	-0.06	0.70
Perturbation Probes	0.69	0.75	-0.02	0.60
Behavioral Monitoring	0.72	0.78	-0.03	0.63
<b>NSF-CoT (Ours)</b>	<b>0.79</b>	<b>0.83</b>	<b>-0.16</b>	<b>0.87</b>
<i>Qwen2.5-72B-Instruct</i>				
Causal Mediation	0.78	0.83	-0.06	0.71
Perturbation Probes	0.73	0.78	-0.02	0.61
Behavioral Monitoring	0.76	0.81	-0.03	0.64
<b>NSF-CoT (Ours)</b>	<b>0.83</b>	<b>0.86</b>	<b>-0.15</b>	<b>0.88</b>

Table 4: **Faithfulness-guided CoT pruning**. Steps failing NSF-CoT verification ( $G_k < 0.5$  or  $V_k < 0.5$ ) are removed before answer decoding.  $Acc_{base}$  matches Table 2;  $\Delta Acc = Acc_{prune} - Acc_{base}$ .

Model	Dataset	Base / Prune	$\Delta Acc$
Pythia-2.8B	OpenBookQA	0.64 / 0.72	+0.08
	QASC	0.62 / 0.72	+0.10
	HotpotQA	0.67 / 0.78	+0.11
LLaMA-3.1-8B	OpenBookQA	0.72 / 0.81	+0.09
	QASC	0.69 / 0.78	+0.09
	HotpotQA	0.75 / 0.84	+0.09
Qwen2.5-72B	OpenBookQA	0.80 / 0.87	+0.07
	QASC	0.77 / 0.84	+0.07
	HotpotQA	0.83 / 0.90	+0.07

groundedness or validity:

$$\mathcal{U}(x) = \left\{ k \in \{1, \dots, K\} : \min(G_k, V_k) < 0.5 \right\}. \quad (15)$$

We exclude  $U_k$  since low utility indicates irrelevance, not incorrectness. The pruned chain  $r^{-\mathcal{U}}$  replaces failing steps with [REMOVED], and we report  $\Delta Acc = Acc_{prune} - Acc_{base}$  under greedy decoding. Table 4 shows mean accuracy gains of +8.7 points after pruning. Gains are larger for smaller models (Pythia: +10–11; LLaMA: +9) versus Qwen (+7), consistent with smaller models producing more decision-harmful reasoning. This demonstrates NSF-CoT functions as both a diagnostic evaluator and an effective test-time corrector—without retraining or model-internal access.

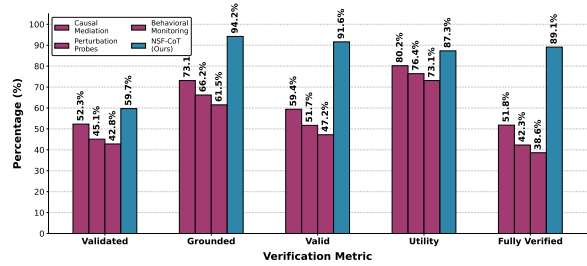


Figure 2: **Cross-verification of baseline-validated steps**. We report the percentage of each method’s validated steps passing NSF-CoT’s groundedness, validity, and utility checks ( $\geq 0.5$ ). Baselines validate 42.8–52.3% of steps, but only 38.6–51.8% satisfy all three criteria. NSF-CoT validates 59.7% of steps with 89.1% fully verified, achieving 94.2% groundedness and 91.6% validity versus 61.5–73.1% and 47.2–59.4% for baselines.

## 5 Related Work

Recent work demonstrates that CoT reasoning is often unfaithful. (Chen et al., 2025) find faithfulness rates of only 25% (Claude 3.7 Sonnet) and 39% (DeepSeek) when inserting hints into prompts, with lower scores on harder questions. (Chua and Evans, 2025) and (Arcuschin et al., 2025) corroborate these findings, identifying restoration errors and unfaithful shortcuts in reasoning traces. (Baker et al., 2025) show that CoT monitoring may fail to detect reward hacking when models are not optimized against monitors. NSF-CoT directly addresses these concerns by combining SMT-based theorem proving with LLM entailment checking to validate each step’s faithfulness.

## 6 Conclusion

In summary, NSF-CoT provides a practical framework for step-level faithfulness verification by combining symbolic checks with model-based judgments and attribution. Across experiments, it reliably flags unsupported reasoning while preserving answer quality, offering a scalable path toward more trustworthy multi-step reasoning.

## 7 Limitations

Our approach has two main limitations.

First, step-level verification in NSF-CoT depends on converting free-form natural language into first-order logic (FOL) so that claims can be checked symbolically against the provided facts. While this enables precise entailment-style verification when parsing succeeds, FOL is not expressive enough to naturally capture all linguistic phenomena found in real model rationales (e.g., implicit commonsense, vague quantifiers, nuanced modality, idioms, or long-range discourse references). As a result, some reasoning steps may be only partially represented, mapped to an oversimplified logical form, or fail to parse altogether, which weakens the verifier’s coverage on the full variety of English statements.

To mitigate this, we incorporate an LLM-based judge that can assess many cases where symbolic parsing is brittle, including paraphrases and linguistically complex steps. However, this fallback is not a complete substitute for formal checking: it can still miss subtle logical errors, introduce variability across prompts or model versions, and provide less transparent failure modes than a symbolic proof. In practice, NSF-CoT’s hybrid design improves robustness, but there remains a gap between the richness of natural language reasoning and what can be consistently captured and validated through FOL-based representations.

Second, our method assumes a level of access that is most feasible with open-weight models. In particular, components such as attribution and step-scoring rely on repeatedly querying the model under many counterfactual ablations and computing token-level likelihoods (or closely related scoring signals). This kind of controlled scoring and repeated evaluation can be difficult or impossible to reproduce with closed-weight, API-only models, where token probabilities, internal scoring behavior, batching constraints, or cost/latency limits may prevent faithful implementation. Consequently, while NSF-CoT is well-suited for open-weight settings and research workflows, extending it to closed models would require alternative interfaces or approximations that may reduce fidelity or comparability.

## 8 Ethical Statement

This work aims to improve the transparency and reliability of chain-of-thought (CoT) reasoning by auditing whether intermediate steps are supported

by provided evidence and aligned with what the model appears to rely on. By flagging unsupported or inconsistently justified steps, our method is intended to reduce the risk of misleading explanations in settings where users may over-trust fluent rationales, such as education, decision support, and information-seeking.

Our approach does not require access to personal user data and is evaluated on public benchmark datasets. Nevertheless, the system may inherit biases present in the underlying language models and datasets, and the LLM-based judge component can introduce additional subjectivity or inconsistency. We therefore emphasize that NSF-CoT should be used as an auditing aid rather than a definitive arbiter of truth, especially in high-stakes domains. When deployed, it should be complemented with domain expertise, careful prompt and threshold calibration, and ongoing monitoring for systematic failure modes (e.g., linguistic constructions that evade parsing or judge errors).

Finally, while improved verification can increase trustworthiness, it may also be misused to create an undeserved appearance of rigor (e.g., selectively reporting verified steps while omitting failures). To mitigate this risk, we recommend reporting verification coverage (e.g., parseability rates, frequency of LLM fallback usage) and making evaluation protocols explicit. We also discuss limitations related to formal expressivity and model accessibility, and we encourage future work on broader semantic coverage and verification methods that are compatible with closed-weight systems.

## References

- Iván Arcuschin, Jett Janiak, Robert Krzyzanowski, Senthoooran Rajamanoharan, Neel Nanda, and Arthur Conmy. 2025. Chain-of-thought reasoning in the wild is not always faithful. *arXiv preprint arXiv:2503.08679*.
- Bowen Baker, Joost Huizinga, Leo Gao, Zehao Dou, Melody Y Guan, Aleksander Madry, Wojciech Zaremba, Jakub Pachocki, and David Farhi. 2025. Monitoring reasoning models for misbehavior and the risks of promoting obfuscation. *arXiv preprint arXiv:2503.11926*.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, and 1 others. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International*

- Conference on Machine Learning*, pages 2397–2430. PMLR.
- Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, and 1 others. 2025. Reasoning models don’t always say what they think. *arXiv preprint arXiv:2505.05410*.
- James Chua and Owain Evans. 2025. Are deepseek r1 and other reasoning models more faithful? *arXiv preprint arXiv:2501.08156*.
- Benjamin Cohen-Wang, Harshay Shah, Kristian Georgiev, and Aleksander Madry. 2024. Contextcite: Attributing model generation to context. *Advances in Neural Information Processing Systems*, 37:95764–95807.
- Leonardo de Moura and Nikolaj Bjørner. 2008. [Z3: An efficient smt solver](#). In *Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? *arXiv preprint arXiv:2004.03685*.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8082–8090.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, and 1 others. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*.
- Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, and 1 others. 2022. Explanations from large language models make small reasoners better. *arXiv preprint arXiv:2210.06726*.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. [Can a suit of armor conduct electricity? a new dataset for open book question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.
- Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum, and Roger Levy. 2023. Linc: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5153–5176.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3806–3824.
- Debjit Paul, Robert West, Antoine Bosselut, and Boi Faltings. 2024. Making reasoning matter: Measuring and improving faithfulness of chain-of-thought reasoning. *arXiv preprint arXiv:2402.13950*.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. Proofwriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634.
- Robert Tibshirani. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. 2024. Faithful logical reasoning via symbolic chain-of-thought. *arXiv preprint arXiv:2405.18357*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 2369–2380.

Xi Ye, Qiaochu Chen, Isil Dillig, and Greg Durrett.  
2023. Satlm: Satisfiability-aided language models  
using declarative prompting. *Advances in Neural  
Information Processing Systems*, 36:45548–45580.

## A Text-to-Logic Parsing

The text-to-logic parser transforms natural language sentences into a restricted first-order logic (FOL) fragment amenable to automated reasoning. This translation is essential because natural language is inherently ambiguous and lacks the formal structure required for symbolic verification. By converting facts and reasoning steps into logical form, we enable precise entailment checking via SMT solvers.

We map natural-language sentences into ground atoms over binary predicates drawn from ConceptNet relations (Speer et al., 2017). Each predicate captures a specific semantic relation between two entities. For example, “Paris is in France” maps to  $\text{AtLocation}(\text{paris}, \text{france})$ , while “A dog is a mammal” maps to  $\text{IsA}(\text{dog}, \text{mammal})$ .

For each context fact  $f_i$ , the parser produces a set of ground atoms:

$$\text{Parse}(f_i) \rightarrow \{\alpha_{i,1}, \dots, \alpha_{i,m_i}\}, \quad (16)$$

where each  $\alpha_{i,j}$  is an atom of the form  $P(a, b)$  with predicate  $P \in \mathcal{P}$  and entities  $a, b$  extracted from the sentence. A single sentence may yield multiple atoms when it expresses multiple relations. Moreover, if a statement is already expressed in first-order logic form, the parser recognizes and preserves it directly without further transformation.

We implement the parser using a large language model OpenAI o3<sup>2</sup> with the following structured prompt:

You are a semantic parser that converts natural language sentences into first-order logic predicates.

Available predicates (from ConceptNet):

- $\text{IsA}(X, Y)$ : X is a type/instance of Y
- $\text{PartOf}(X, Y)$ : X is a part of Y
- $\text{AtLocation}(X, Y)$ : X is located in/at Y
- $\text{HasProperty}(X, Y)$ : X has property Y
- $\text{CapableOf}(X, Y)$ : X is capable of Y
- $\text{UsedFor}(X, Y)$ : X is used for Y
- $\text{MadeOf}(X, Y)$ : X is made of Y
- $\text{HasA}(X, Y)$ : X has/possesses Y
- $\text{Causes}(X, Y)$ : X causes Y
- $\text{HasPrerequisite}(X, Y)$ : X requires Y to happen first
- $\text{HasEffect}(X, Y)$ : X results in Y

Rules:

1. Output ONLY predicates, one per line
2. Use lowercase for all entity arguments
3. Extract the core semantic relation(s) from the sentence
4. If a sentence contains multiple relations, output multiple predicates
5. If the input is already in FOL format (e.g., “ $\text{IsA}(\text{dog}, \text{mammal})$ ”), preserve it exactly as given
6. If no predicate applies, output “None”
7. Do not add information not present in the sentence

Examples:

Input: “The Eiffel Tower is in Paris”  
Output:  $\text{AtLocation}(\text{eiffel\_tower}, \text{paris})$

Input: “Dogs are mammals that have fur”  
Output:  $\text{IsA}(\text{dog}, \text{mammal})$   
 $\text{HasProperty}(\text{dog}, \text{fur})$

Input: “The engine is part of the car”  
Output:  $\text{PartOf}(\text{engine}, \text{car})$

Input: “A thermometer is used for measuring temperature”  
Output:  $\text{UsedFor}(\text{thermometer}, \text{measuring\_temperature})$

Input: “Glass is made of sand”  
Output:  $\text{MadeOf}(\text{glass}, \text{sand})$

Input: “Birds can fly”  
Output:  $\text{CapableOf}(\text{bird}, \text{fly})$

Input: “Heavy rain causes flooding”  
Output:  $\text{Causes}(\text{heavy\_rain}, \text{flooding})$

Input: “Baking requires preheating the oven”  
Output:  $\text{HasPrerequisite}(\text{baking}, \text{preheating\_oven})$

Input: “ $\text{IsA}(\text{cat}, \text{animal})$ ”  
Output:  $\text{IsA}(\text{cat}, \text{animal})$

Now parse the following sentence:

Input: “<sentence>”

The full knowledge base is formed by aggregating parsed atoms across all context facts:

$$\mathcal{KB} = \bigcup_{i=1}^n \text{Parse}(f_i). \quad (17)$$

This knowledge base serves as the formal representation of all information provided in the context,

<sup>2</sup><https://platform.openai.com/docs/models/o3>

over which entailment queries are evaluated.

## B Internal Fact Attribution

We approximate internal reliance using `CONTEXTCITE`, which measures how removing facts changes the probability of generating each step. This method is well-suited to our framework because, unlike attention-based approaches that capture only correlational patterns, it provides *contributive* attribution by measuring counterfactual effects—a fact receives high attribution if and only if its removal substantially changes the model’s output.

**Ablation Sampling.** The core idea is to systematically remove subsets of context facts and observe how generation probability changes. We represent each ablation configuration as a binary vector  $v \in \{0, 1\}^n$ , where  $v_i = 1$  indicates that fact  $f_i$  is retained and  $v_i = 0$  indicates it is removed. The operation `Ablate`( $\mathcal{F}, v$ ) constructs the ablated context containing only the retained facts. For a reasoning step  $s_k$  spanning tokens  $y_{[a_k:b_k]}$ , we define the ablation score:

$$g_k(v) = \log p_{\text{LM}}\left(y_{[a_k:b_k]} \mid \text{Ablate}(\mathcal{F}, v), q, y_{<a_k}\right). \quad (18)$$

Intuitively, if removing a fact causes a large drop in  $g_k(v)$ , that fact was important for generating step  $s_k$ .

**LASSO Regression.** Since exhaustively evaluating all  $2^n$  ablations is infeasible, we fit a sparse linear surrogate model. We sample  $M$  ablation vectors  $\{v^{(m)}\}_{m=1}^M$  uniformly from  $\{0, 1\}^n$ , compute their scores, and construct design matrix  $V \in \{0, 1\}^{M \times n}$  and target vector  $z_k \in \mathbb{R}^M$ . The attribution weights are estimated via LASSO regression (Tibshirani, 1996):

$$\hat{w}_k = \arg \min_{w \in \mathbb{R}^n} \left( \|Vw - z_k\|_2^2 + \lambda \|w\|_1 \right), \quad (19)$$

where  $\lambda > 0$  controls sparsity. The weight  $\hat{w}_{k,i}$  quantifies the contributive effect of fact  $f_i$  on generating step  $s_k$ : positive weights indicate supportive facts, while negative weights indicate distractors.

**Thresholding and Discretization.** For downstream validity checking, we convert continuous weights into a discrete set of internally relied-upon facts. Let  $\mathcal{J}_k = \{i : \hat{w}_{k,i} > 0\}$  denote indices with

positive attribution. If  $\mathcal{J}_k = \emptyset$ , we set  $\mathcal{I}_k = \emptyset$ . Otherwise, we compute the mean over positive weights and threshold:

$$\bar{w}_k^+ = \frac{1}{|\mathcal{J}_k|} \sum_{i \in \mathcal{J}_k} \hat{w}_{k,i}, \quad (20)$$

$$\mathcal{I}_k = \left\{ f_i : i \in \mathcal{J}_k \wedge \hat{w}_{k,i} \geq \tau \cdot \bar{w}_k^+ \right\}, \quad (21)$$

where  $\tau$  is tuned on validation data (default  $\tau = 0.5$ ). This set captures the facts that most strongly influenced the model’s generation of step  $s_k$ .

**Soft Attribution Probabilities.** To avoid threshold sensitivity in faithfulness scoring, we additionally compute a soft probability distribution by normalizing positive weights:

$$\alpha_{k,i} = \frac{\max(\hat{w}_{k,i}, 0)}{\sum_{j=1}^n \max(\hat{w}_{k,j}, 0)}, \quad (22)$$

with  $\alpha_{k,i} = 1/n$  for all  $i$  if the denominator is zero. The value  $\alpha_{k,i}$  captures how much the model depended on fact  $f_i$  when producing step  $s_k$ , enabling continuous faithfulness scoring without hard discretization artifacts.

## C SMT Inference Axioms

To support multi-hop reasoning, we include a fixed library of task-specific inference axioms  $\phi_{\text{rules}}$ . These rules define the entailment relation  $\vdash_{\text{rules}}$  used throughout our groundedness and validity checks; they are treated as verification axioms rather than universal semantic truths. Let  $\mathcal{E}$  denote the entity domain.

**Transitivity.** We include transitivity for relations that naturally compose over chains in our target domains:

$$\phi_{\text{trans}} = \bigwedge_{R \in \mathcal{R}_T} \forall x, y, z \in \mathcal{E} : \left( R(x, y) \wedge R(y, z) \right) \Rightarrow R(x, z), \quad (23)$$

where  $\mathcal{R}_T = \{\text{IsA}, \text{PartOf}, \text{AtLocation}, \text{Causes}, \text{HasPrerequisite}\}$ .

**Part-Location Composition.** We include rules that propagate location through part-whole structure:

$$\phi_{\text{part\_loc}} = \forall x, y, z \in \mathcal{E} : \left( \text{PartOf}(x, y) \wedge \text{AtLocation}(y, z) \right) \Rightarrow \text{AtLocation}(x, z), \quad (24)$$

$$\begin{aligned} \phi_{\text{loc\_part}} &= \forall x, y, z \in \mathcal{E} : \\ &\left( \text{AtLocation}(x, y) \wedge \text{PartOf}(y, z) \right) \\ &\Rightarrow \text{AtLocation}(x, z). \end{aligned} \quad (25)$$

**Property Inheritance.** We propagate properties and capabilities along ISA links:

$$\begin{aligned} \phi_{\text{prop}} &= \forall x, y, p \in \mathcal{E} : \\ &\left( \text{ISA}(x, y) \wedge \text{HasProperty}(y, p) \right) \\ &\Rightarrow \text{HasProperty}(x, p), \end{aligned} \quad (26)$$

$$\begin{aligned} \phi_{\text{cap}} &= \forall x, y, a \in \mathcal{E} : \\ &\left( \text{ISA}(x, y) \wedge \text{CapableOf}(y, a) \right) \\ &\Rightarrow \text{CapableOf}(x, a), \end{aligned} \quad (27)$$

$$\begin{aligned} \phi_{\text{has}} &= \forall x, y, z \in \mathcal{E} : \\ &\left( \text{ISA}(x, y) \wedge \text{HasA}(y, z) \right) \\ &\Rightarrow \text{HasA}(x, z). \end{aligned} \quad (28)$$

**Causal Composition.** We link causes to their downstream effects:

$$\begin{aligned} \phi_{\text{causal}} &= \forall x, y, z \in \mathcal{E} : \\ &\left( \text{Causes}(x, y) \wedge \text{HasEffect}(y, z) \right) \\ &\Rightarrow \text{Causes}(x, z). \end{aligned} \quad (29)$$

**Combined Rule Set.** The full axiom set is:

$$\begin{aligned} \phi_{\text{rules}} &= \phi_{\text{trans}} \wedge \phi_{\text{part\_loc}} \wedge \phi_{\text{loc\_part}} \\ &\wedge \phi_{\text{prop}} \wedge \phi_{\text{cap}} \wedge \phi_{\text{has}} \wedge \phi_{\text{causal}}. \end{aligned} \quad (30)$$

**Z3 Encoding.** Each predicate symbol  $P$  is interpreted as a Boolean relation  $P : \mathcal{E} \times \mathcal{E} \rightarrow \mathbb{B}$ . For each context fact  $f_i$ , we encode its logical content as the conjunction of its parsed atoms:

$$\phi_i = \bigwedge_{\alpha \in \text{Parse}(f_i)} \alpha. \quad (31)$$

To enable selective fact activation during solving, we introduce a Boolean assumption literal  $a_i$  for each fact and encode the knowledge base as:

$$\phi_{\mathcal{F}} = \bigwedge_{i=1}^n \left( a_i \Rightarrow \phi_i \right). \quad (32)$$

This conditional encoding allows us to enable or disable individual facts by including  $a_i$  or  $\neg a_i$  in the solver’s assumption set (de Moura and Björner, 2008).

**Proof-Support Extraction.** To test whether a claim  $c$  is entailed, we apply proof by refutation. Let  $A^+ = \{a_1, \dots, a_n\}$  enable all facts. A claim is entailed if asserting its negation yields a contradiction:

$$\begin{aligned} \mathcal{KB} \vdash_{\text{rules}} c &\iff \text{UNSAT} \left( \phi_{\mathcal{F}} \wedge \phi_{\text{rules}} \wedge \neg c \right) \\ &\text{under } A^+. \end{aligned} \quad (33)$$

When UNSAT, Z3 returns an UNSAT core  $\text{Core}(c) \subseteq A^+$ , which is a sufficient (though not necessarily minimal) subset of assumptions for the contradiction. We map this to the proof-support set:

$$\mathcal{P}(c) = \{f_i \in \mathcal{F} : a_i \in \text{Core}(c)\}. \quad (34)$$

For a step  $s_k$ , we aggregate proof-support over entailed claims:

$$\mathcal{R}_k = \bigcup_{\substack{c \in \text{Claims}(s_k) \\ \mathcal{KB} \vdash_{\text{rules}} c}} \mathcal{P}(c). \quad (35)$$

## D LLM Entailment Judge Prompts

We employ an LLM as a secondary entailment judge to handle cases that escape the formal language of the SMT solver. Given a premise and hypothesis, the LLM produces (i) a natural-language reasoning trace explaining the entailment relationship, and (ii) a calibrated score  $\sigma \in [0, 1]$  indicating confidence that the premise entails the hypothesis. We provide the SMT results as context so the LLM can corroborate symbolic findings or identify cases where linguistic nuance leads to different conclusions.

**Groundedness Prompt.** For checking whether claims are derivable from the full context:

You are an entailment judge. Given a set of context facts and a claim, determine whether the claim logically follows from the context.

Context Facts:  
<list of facts  $\mathcal{F}$ >

Claim to verify:  
<claims from  $\text{Claims}(s_k)$ >

SMT Solver Result: < $G_k^{\text{SMT}}$ >  
Proof Support (facts used): < $\mathcal{R}_k$ >

Instructions:

1. Provide a step-by-step reasoning trace explaining whether the claim follows from the context
2. Consider both explicit statements and reasonable inferences

- Output a score between 0 and 1 indicating your confidence that the context entails the claim

Output format:

Trace: <your reasoning>

Score: <0.0 to 1.0>

**Validity Prompt.** For checking whether claims are derivable from only the internally relied-upon facts:

You are an entailment judge. Given a restricted set of facts that the model internally relied upon, determine whether the claim logically follows from ONLY these facts.

Internally Relied-Upon Facts:

<list of facts  $\mathcal{I}_k$ >

Claim to verify:

<claims from  $\text{Claims}(s_k)$ >

SMT Solver Result: < $V_k^{\text{SMT}}$ >

Instructions:

- ONLY use the internally relied-upon facts listed above
- Do NOT use any external knowledge or other context facts
- Provide a step-by-step reasoning trace
- Output a score between 0 and 1 indicating your confidence

Output format:

Trace: <your reasoning>

Score: <0.0 to 1.0>

**Utility Prompt.** For checking whether the step’s claims contribute to deriving the final answer:

You are an entailment judge. Determine whether the given claims from a reasoning step help derive the final answer.

Step Claims:

<claims from  $\text{Claims}(s_k)$ >

Final Answer Claim:

< $C_{\text{ans}}$ >

SMT Solver Result: < $U_k^{\text{SMT}}$ >

Answer Proof Support: < $\mathcal{R}_{\text{ans}}$ >

Instructions:

- Analyze whether the step claims are necessary or helpful for reaching the final answer
- Consider if removing these claims would make the answer harder to derive
- Provide a step-by-step reasoning trace
- Output a score between 0 and 1 indicating how much this step contributes to the answer

Output format:

Trace: <your reasoning>

Score: <0.0 to 1.0>

The reasoning traces  $\mathcal{T}_k^G, \mathcal{T}_k^V, \mathcal{T}_k^U$  are retained for interpretability and can be used for downstream error analysis to understand disagreements between the SMT solver and LLM judge.

Table 5: Component ablation on OpenBookQA (LLaMA-3.1-8B). IR is constant by construction;  $\Delta\text{VAD}$  shows change from full NSF-CoT.

Variant	IR	VAD	$\Delta\text{VAD}$
NSF-CoT (Full)	0.40	0.61	–
– SMT (LLM-only)	0.40	0.53	–0.08
– LLM (SMT-only)	0.40	0.45	–0.16
– Attribution	0.40	0.52	–0.09
No LLM fallback	0.40	0.58	–0.03

## E Extended Ablation Studies

### E.1 Component Ablation

We ablate key components of NSF-CoT on OpenBookQA with LLaMA-3.1-8B, reporting IR and VAD. IR remains constant across ablations since it depends on chain-substitution intervention, not verifier internals. VAD varies because each variant changes which steps are validated.

**Ablated Variants.** –SMT sets  $\beta = 0$  (LLM-judge only). –LLM sets  $\beta = 1$  (SMT-only); unparseable steps receive zero scores. –Attribution sets  $\mathcal{I}_k = \mathcal{F}$  for all steps, collapsing validity toward groundedness. **No LLM fallback** retains hybrid scoring for parseable steps but assigns zero to unparseable steps—unlike –LLM which removes the LLM signal everywhere, this removes it only when parsing fails.

**Results.** Table 5 shows that removing the LLM judge causes the largest VAD degradation (–0.16), reflecting both reduced coverage for unparseable steps and loss of linguistic flexibility for entailment judgment. Removing SMT verification also reduces VAD (–0.08), indicating symbolic checking catches logical errors the LLM may miss. Disabling attribution reduces VAD (–0.09) by eliminating the distinction between context-supported and internally-relied-upon evidence. Removing only the fallback for unparseable steps causes minimal drop (–0.03), suggesting the hybrid design primarily provides robustness to parser coverage gaps.

### E.2 More Qualitative Examples

To further illustrate NSF-CoT’s step-level verification, we provide additional qualitative examples beyond the main case study. Tables 6–9 show four more examples (OpenBookQA, QASC, and HotpotQA) that exhibit the same pattern: most steps are grounded and valid with respect to the provided

**Question:** Why do most plants need sunlight to grow?  
 (A) to absorb oxygen (B) to make food (C) to remove minerals (D) to create soil **Gold Answer:** (B)

**Core Fact:**

$f_1$ : Photosynthesis uses sunlight to produce sugars (food) in plants.

**Retrieved Facts:**

$f_2$ : Plants are living organisms that require energy to grow.  
 $f_3$ : Sugars produced in photosynthesis provide energy for plant processes.  
 $f_4$ : Sunlight is a source of energy.

**Model-written rationale (LLaMA-3.1-8B):**

$s_1$ : Sunlight provides energy that plants can use.  
 $s_2$ : Photosynthesis uses sunlight to make sugars in plants.  
 $s_3$ : Sugars are food/energy for plant growth.  
 $s_4$ : Plants also need sunlight to absorb nutrients from soil more efficiently.  
 $s_5$ : Therefore, most plants need sunlight to make food.

**Predicted:** (B) ✓

Step	Parsed claim (simplified)	$\mathcal{I}_k$	$G_k$	$V_k$	$U_k$	$F_k$	Verdict
$s_1$	HasProperty(sunlight, energy)	$\{f_4\}$	0.92	0.88	0.70	0.57	Faithful
$s_2$	Causes(sunlight, photosynthesis_sugars)	$\{f_1\}$	0.97	0.94	0.88	0.80	Faithful
$s_3$	Causes(sugars, plant_energy)	$\{f_3\}$	0.95	0.90	0.82	0.70	Faithful
$s_4$	Causes(sunlight, absorb_nutrients)	$\emptyset$	0.26	0.10	0.32	0.01	Unfaithful
$s_5$	Causes(photosynthesis_sugars, plant_food)	$\{f_1, f_3\}$	0.94	0.89	0.96	0.80	Faithful

Table 6: OpenBookQA example (photosynthesis). Step  $s_4$  is **unfaithful**: it introduces a plausible-sounding claim about sunlight improving nutrient absorption that is unsupported by the provided facts ( $\mathcal{I}_4 = \emptyset$ ). The answer is still derived through the faithful path  $s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_5$ .

evidence, but a small number of fluent “extra” statements are flagged as unfaithful because they are not entailed by the context and are not attributed to internally relied-upon facts.

This is an appendix.

**Question (QASC):** What process causes a wet shirt to become dry when hung outside? **Gold Answer:** Evaporation

**Context Facts:**

- $f_1$ : Evaporation is the process where liquid water changes into water vapor.
- $f_2$ : Heat increases the rate of evaporation.

**Retrieved Facts:**

- $f_3$ : Sunlight can warm objects by providing heat energy.
- $f_4$ : Wet clothing contains liquid water.

**Model-written rationale (LLaMA-3.1-8B):**

- $s_1$ : Wet clothing contains liquid water.
- $s_2$ : Sunlight provides heat that warms the shirt.
- $s_3$ : Heat increases the rate at which water evaporates into vapor.
- $s_4$ : Wind always removes all water molecules faster than heat does.
- $s_5$ : Therefore, the shirt dries due to evaporation.

**Predicted:** Evaporation ✓

Step	Parsed claim (simplified)	$\mathcal{I}_k$	$G_k$	$V_k$	$U_k$	$F_k$	Verdict
$s_1$	HasA(clothes, liquid_water)	$\{f_4\}$	0.93	0.89	0.75	0.62	Faithful
$s_2$	Causes(sunlight, heat)	$\{f_3\}$	0.95	0.90	0.78	0.67	Faithful
$s_3$	Causes(heat, evaporation_rate_up)	$\{f_1, f_2\}$	0.96	0.92	0.90	0.79	Faithful
$s_4$	HasProperty(wind, always_dominates)	$\emptyset$	0.18	0.07	0.21	0.00	Unfaithful
$s_5$	Causes(evaporation, clothes_dry)	$\{f_1, f_4\}$	0.94	0.88	0.97	0.80	Faithful

Table 7: QASC compositional example. Step  $s_4$  is **unfaithful**: it asserts an absolute comparative claim about wind that is not supported by the context and is not relied upon ( $\mathcal{I}_4 = \emptyset$ ). The correct answer follows from the faithful chain  $s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_5$ .

**Question (HotpotQA):** Can a penguin fly?

**Gold Answer:** No

**Context Facts:**

- $f_1$ : Penguins are birds that are adapted for swimming.
- $f_2$ : Many birds fly using wings, but some birds cannot fly.
- $f_3$ : Penguins have wings shaped like flippers for swimming.

**Model-written rationale (LLaMA-3.1-8B):**

- $s_1$ : Penguins are birds, and some birds cannot fly.
- $s_2$ : Penguins use flipper-like wings for swimming rather than flight.
- $s_3$ : Therefore penguins are a type of bird that cannot fly.
- $s_4$ : Penguins are too heavy to generate the minimum 12 N/kg lift required for flight.
- $s_5$ : So a penguin cannot fly.

**Predicted:** No ✓

Step	Parsed claim (simplified)	$\mathcal{I}_k$	$G_k$	$V_k$	$U_k$	$F_k$	Verdict
$s_1$	IsA(penguin, bird)	$\{f_1, f_2\}$	0.92	0.86	0.70	0.55	Faithful
$s_2$	HasProperty(penguin_wing, flipper)	$\{f_3\}$	0.95	0.91	0.80	0.69	Faithful
$s_3$	$\neg$ CapableOf(penguin, fly)	$\{f_2, f_3\}$	0.90	0.84	0.88	0.67	Faithful
$s_4$	HasProperty(penguin, lift_threshold_12)	$\emptyset$	0.20	0.06	0.25	0.00	Unfaithful
$s_5$	$\neg$ CapableOf(penguin, fly)	$\{f_2, f_3\}$	0.91	0.85	0.96	0.74	Faithful

Table 8: HotpotQA example with an over-specific numerical claim. Step  $s_4$  is **unfaithful**: it introduces an unsupported physics threshold not entailed by  $\mathcal{F}$  and not internally attributed ( $\mathcal{I}_4 = \emptyset$ ), while the final answer is already justified by the faithful steps  $s_1$ – $s_3$  and  $s_5$ .

---

**Question (HotpotQA):** Would an ice cube melt faster in salt water than in fresh water? **Gold Answer:** Yes

---

**Context Facts:**

- $f_1$ : Salt dissolved in water lowers water’s freezing point.
  - $f_2$ : When the freezing point is lower, ice is less stable and melts more easily.
- 

**Retrieved Facts:**

- $f_3$ : Fresh water has little or no dissolved salt.
  - $f_4$ : Salt water contains dissolved salt.
- 

**Model-written rationale (LLaMA-3.1-8B):**

- $s_1$ : Salt water contains dissolved salt, while fresh water contains little salt.
- $s_2$ : Dissolved salt lowers the freezing point of water.
- $s_3$ : A lower freezing point makes ice melt more easily.
- $s_4$ : Also, salt water is always warmer than fresh water in nature.
- $s_5$ : Therefore an ice cube melts faster in salt water than in fresh water.

**Predicted:** Yes ✓

---

Step	Parsed claim (simplified)	$\mathcal{I}_k$	$G_k$	$V_k$	$U_k$	$F_k$	Verdict
$s_1$	HasProperty(salt_water, salty)	$\{f_3, f_4\}$	0.92	0.86	0.72	0.57	Faithful
$s_2$	Causes(salt, freezing_point_down)	$\{f_1\}$	0.96	0.92	0.86	0.76	Faithful
$s_3$	Causes(freezing_point_down, melt_easier)	$\{f_2\}$	0.95	0.90	0.86	0.74	Faithful
$s_4$	HasProperty(salt_water, always_warmer)	$\emptyset$	0.15	0.05	0.22	0.00	Unfaithful
$s_5$	Causes(freezing_point_down, melt_faster)	$\{f_1, f_2\}$	0.93	0.88	0.97	0.79	Faithful

Table 9: HotpotQA example where the model introduces an unsupported generalization. Step  $s_4$  is **unfaithful**: it asserts that salt water is always warmer than fresh water, which is not supported by  $\mathcal{F}$  and is not attributed to any context fact ( $\mathcal{I}_4 = \emptyset$ ). The correct answer is nonetheless justified by the faithful steps  $s_1$ – $s_3$  and  $s_5$ .

---

**Algorithm 1: NSF-CoT: Neuro-Symbolic Faithfulness Verification for Chain-of-Thought**


---

**Input:** Context facts  $\mathcal{F} = \{f_1, \dots, f_n\}$ ,  
task/question  $q$ , LM  $\text{LM}$ , SMT solver  $\text{Z3}$ ,  
LLM judge  $J$ , rules  $\phi_{\text{rules}}$ , ablation budget  $M$ ,  
LASSO parameter  $\lambda$ , attribution threshold  $\tau$ ,  
hybrid weight  $\beta$

**Output:** Step-level scores  $\{(G_k, V_k, U_k, \mathcal{F}_k)\}_{k=1}^K$   
and step verdicts

**Stage 1: CoT generation and segmentation**

Generate response  $y \sim p_{\text{LM}}(\cdot \mid [\mathcal{F}; q])$  and segment  
into steps  $\mathcal{S} = \langle s_1, \dots, s_K \rangle$  with token spans

$$s_k \equiv y_{[a_k:b_k]}.$$

**Stage 2: Text-to-logic parsing**

$\mathcal{KB} \leftarrow \bigcup_{i=1}^n \text{Parse}(f_i)$ ; // parse facts into  
FOL atoms

**for  $k \leftarrow 1$  to  $K$  do**

Claims( $s_k$ )  $\leftarrow$  Parse( $s_k$ ); // parse step  
into claims

$c_{\text{ans}} \leftarrow$  Parse( $y_{\text{final}}$ ); // parse final answer  
into target claim

**Stage 3: Internal fact attribution**

(ContextCite-style)

**for  $k \leftarrow 1$  to  $K$  do**

Sample  $M$  masks  $v^{(1)}, \dots, v^{(M)} \in \{0, 1\}^n$ ;

**for  $m \leftarrow 1$  to  $M$  do**

$z_k[m] \leftarrow$   
 $\log p_{\text{LM}}(y_{[a_k:b_k]} \mid \text{Ablate}(\mathcal{F}, v^{(m)}), q, y_{<a_k})$ ;

$V[m, :] \leftarrow v^{(m)}$ ;

$\hat{w}_k \leftarrow \arg \min_{w \in \mathbb{R}^n} (\|Vw - z_k\|_2^2 + \lambda \|w\|_1)$ ;  
// LASSO surrogate

$\bar{w}_k^+ \leftarrow \text{mean}\{\hat{w}_{k,i} : \hat{w}_{k,i} > 0\}$ ;

$\mathcal{I}_k \leftarrow \{f_i : \hat{w}_{k,i} \geq \tau \cdot \bar{w}_k^+\}$ ; // internally  
relied-upon facts

**Stage 4: Hybrid verification (SMT + LLM judge)**

Encode each fact  $f_i$  as  $\phi_i = \bigwedge_{\alpha \in \text{Parse}(f_i)} \alpha$  and

build conditional KB:  $\phi_{\mathcal{F}} = \bigwedge_{i=1}^n (a_i \Rightarrow \phi_i)$  with  
assumption literals  $\{a_i\}$ ;

**for  $k \leftarrow 1$  to  $K$  do**

// SMT checks

$G_k^{\text{SMT}} \leftarrow$   
 $\mathbb{1}[\forall c \in \text{Claims}(s_k) : \text{UNSAT}(\phi_{\mathcal{F}} \wedge \phi_{\text{rules}} \wedge \neg c \mid A^+ = \{a_1, \dots, a_n\})]$ ;

$V_k^{\text{SMT}} \leftarrow$

$\mathbb{1}[\forall c \in \text{Claims}(s_k) : \text{UNSAT}(\phi_{\mathcal{F}} \wedge \phi_{\text{rules}} \wedge \neg c \mid A_k = \{a_i : f_i \in \mathcal{I}_k\})]$ ;

Extract UNSAT cores to obtain proof-support

sets  $\mathcal{R}_k$  for step claims and  $\mathcal{R}_{\text{ans}}$  for  $c_{\text{ans}}$ ;

$U_k^{\text{SMT}} \leftarrow \mathbb{1}[\mathcal{R}_k \cap \mathcal{R}_{\text{ans}} \neq \emptyset]$ ;

// LLM entailment judge

$(G_k^{\text{LLM}}, V_k^{\text{LLM}}, U_k^{\text{LLM}}) \leftarrow$

Judge( $J$ ;  $\mathcal{F}$ ,  $\mathcal{I}_k$ ,  $s_k$ ,  $c_{\text{ans}}$ , SMT evidence);

// Hybrid aggregation

$G_k \leftarrow \beta G_k^{\text{SMT}} + (1 - \beta) G_k^{\text{LLM}}$ ;

$V_k \leftarrow \beta V_k^{\text{SMT}} + (1 - \beta) V_k^{\text{LLM}}$ ;

$U_k \leftarrow \beta U_k^{\text{SMT}} + (1 - \beta) U_k^{\text{LLM}}$ ;

**Stage 5: Faithfulness scoring**
**for  $k \leftarrow 1$  to  $K$  do**

$\mathcal{F}_k \leftarrow G_k \cdot V_k \cdot U_k$ ;

Verdict( $s_k$ )  $\leftarrow \mathbb{1}[\mathcal{F}_k > 0]$ ; // or threshold  
if desired

**return**  $\{(G_k, V_k, U_k, \mathcal{F}_k, \text{Verdict}(s_k))\}_{k=1}^K$ ;

---