

# Decompose-and-Formalise: Recursively Verifiable Natural Language Inference

Anonymous ACL submission

## Abstract

Recent work has shown that integrating large language models (LLMs) with theorem provers (TPs) in neuro-symbolic pipelines helps with autoformalisation, entailment verification, and proof-guided refinement of explanations for natural language inference (NLI). However, scaling such refinement to naturalistic NLI remains difficult: long, syntactically rich inputs and deep multi-step arguments amplify autoformalisation errors, where a single local mismatch can invalidate the proof. Moreover, current methods often handle failures via costly global regeneration due to the difficulty of localising the responsible span or step from prover diagnostics. Aiming to address these problems, we propose a decompose-and-formalise framework that (i) decomposes premise-hypothesis pairs into an entailment tree of atomic steps, (ii) verifies the tree bottom-up to isolate failures to specific nodes, and (iii) performs local diagnostic-guided refinement instead of regenerating the whole explanation. Moreover, to improve faithfulness of autoformalisation, we introduce  $\theta$ -substitution in an event-based logical form to enforce consistent argument–role bindings. Across a range of reasoning tasks using five LLM backbones, our method achieves the highest verified explanation rates, improving over the state-of-the-art by 26.2%, 21.7%, 21.6% and 48.9%, while reducing refinement iterations and runtime while preserving strong end-NLI accuracy.

## 1 Introduction

Natural language explanations are the workhorse of human communication and rationality, yet they are frequently incomplete, imprecise, and logically fragile (Valentino and Freitas, 2024; Valentino et al., 2021). This motivates verifiable Natural Language Inference (NLI) systems that can make the implicit commitments of everyday explanations

explicit, and refine them into solver-checkable stepwise certificates (Quan et al., 2024; Sadeddine and Suchanek, 2025). Recent works in neuro-symbolic pipelines offer a principled route to this goal by integrating Large Language Models (LLMs) with external theorem provers (TPs): LLMs autoformalise natural language statements into formal representations, and TPs then verify whether premises and explanations entail a hypothesis in the target logic (Pan et al., 2023; Olausson et al., 2023; Quan et al., 2025a; Wu et al., 2022; Xu et al., 2024, 2025c; Zhang et al., 2025; Jiang et al., 2024). TP feedback can further support iterative refinement of both formalisation and explanations (Quan et al., 2024, 2025b; Xu et al., 2025b).

Despite the pace of progress, scaling formal refinement to naturalistic and complex NLI remains an open problem. Real explanations often involve long, syntactically rich sentences and deep multi-step arguments that compress many coupled semantic commitments into a small surface form. A single local mismatch (such as a role swap, a negation-scope error, or a quantifier mistake) can derail an entire reasoning chain. Moreover, because these errors propagate across steps, end-to-end regeneration and re-proving becomes increasingly unstable and costly as chains become longer and more branched (Pan et al., 2023; Quan et al., 2025b). At the same time, autoformalisation must be both prover-compatible and semantically faithful (preserving scope, roles, quantification, and lexical commitments) because even minor semantic drift can yield proofs that are syntactically valid but misaligned with the original text (Jiang et al., 2024; Xu et al., 2025b; Quan et al., 2025b). These bottlenecks foreground three questions: *RQ1: How can we obtain autoformalisation that is both prover-compatible and semantically faithful?* *RQ2: What kind of verification strategy provides transparent, stepwise verification that preserves multi-step structure and reduces refinement*

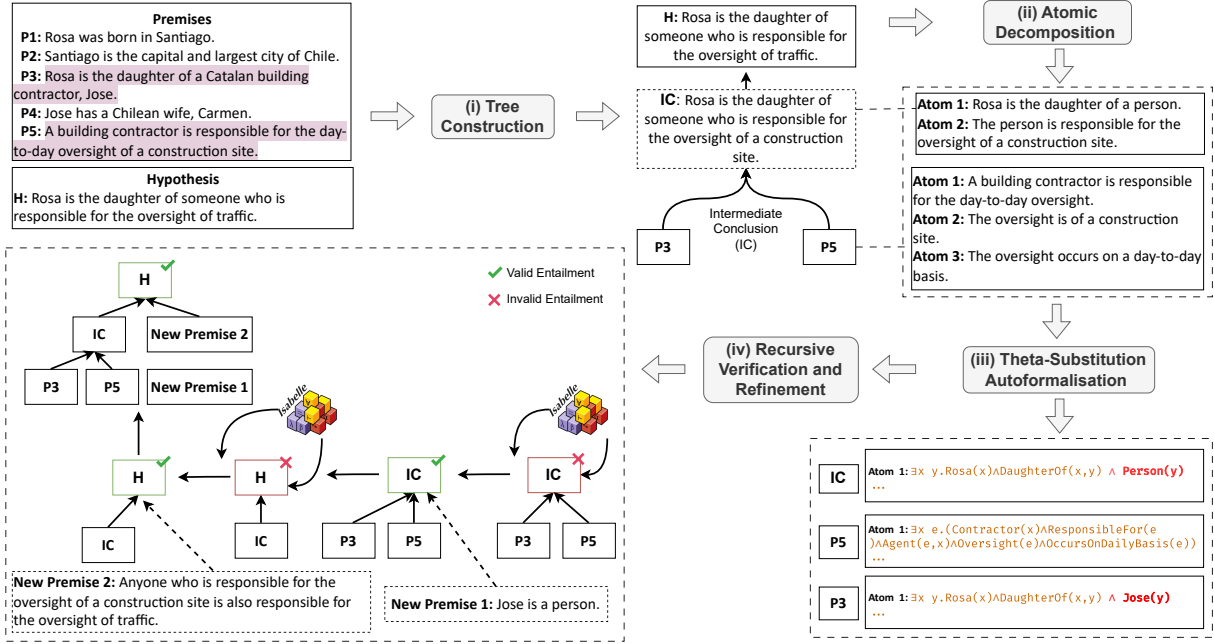


Figure 1: An illustration of the proposed framework for recursively verifiable natural language inference (NLI) via entailment trees, atomic decomposition, and  $\theta$ -substitution autoformalisation. The entailment tree is initially constructed from the given premises and hypothesis, including intermediate conclusions. Each sentence is then decomposed into atomic propositions and autoformalised into logical forms. By utilising an external theorem prover, we progressively verify and refine each subtree until the final hypothesis node is reached. The reasoning chain is considered logically sound and coherent once the entire entailment tree has been fully verified and refined.

083 *drift? RQ3: To what extent can local diagnostic-*  
 084 *guided refinement improve efficiency and robust-*  
 085 *ness without sacrificing end-task accuracy?*

086 To address these challenges, we propose *LLM-*  
 087 *TP Tree*, a scalable decompose-and-formalise  
 088 framework built around an entailment-tree view  
 089 of explanations. The core idea is to replace holistic  
 090 refinement with structured decomposition and  
 091 local repair. Concretely, LLM-TP Tree (i) decom-  
 092 poses each premise-hypothesis instance into an  
 093 entailment tree of atomic steps, (ii) verifies the  
 094 tree bottom-up to isolate proof failures to specific  
 095 nodes, and (iii) performs local diagnostic-guided  
 096 refinement on the failing nodes rather than regen-  
 097 erating the entire explanation. This design makes  
 098 semantic preservation controllable at the smallest  
 099 meaningful granularity, stabilises refinement under  
 100 longer-horizon reasoning, and improves efficiency  
 101 by focusing computation where proofs actually  
 102 fail. To improve faithfulness in autoformalisation,  
 103 we introduce a  $\theta$ -substitution procedure within an  
 104 event-based logical form that enforces consistent  
 105 argument-role bindings during autoformalisation,  
 106 making role commitments explicit and systemati-  
 107 cally repairable. Together, atomic decomposition,  
 108 recursive bottom-up verification, and  $\theta$ -substitution  
 109 yield prover-compatible theories that better track

the intended semantics of the original text while  
 enabling targeted refinement.

We evaluate LLM-TP Tree on FOLIO, ProofWriter, PrOntoQA, and EntailmentBank using five LLM backbones (GPT-4o, GPT-5 nano, Grok-4-fast, Deepseek-V3.1 and Qwen3-Max). We compare against state-of-the-art baselines including Explanation-Refiner (Quan et al., 2024), Faithful-Refiner (Quan et al., 2025b), LINC (Olausson et al., 2023), and Logic-LM (Pan et al., 2023). Across datasets, our approach achieves the strongest refinement performance, improving the average verification rate by 26.2%, 21.7%, 21.6%, and 48.9% points on FOLIO, ProofWriter, PrOntoQA, and EntailmentBank, respectively, while reducing refinement iterations and runtime and maintaining strong NLI accuracy. In summary, our contributions are:

1. We introduce LLM-TP Tree, a recursive entailment-tree framework that verifies and refines multi-step inference chains for verifiable NLI via bottom-up checking and local repair.
2. We propose atomic decomposition, which decomposes complex sentences into atomic propositions and reduces inference-time by

16.5%, 13.7%, 21.1%, and 10.1% across datasets compared to baseline neuro-symbolic frameworks.

3. We propose a  $\theta$ -substitution autoformalisation procedure that enforces consistent argument-role bindings and improves autoformalisation faithfulness by 0.09, 0.106, 0.128, and 0.21 on average over baseline models.
4. We conduct comprehensive automatic and human evaluations with ablations, showing that LLM-TP Tree substantially improves the robustness and efficiency of LLM-TP explanation refinement while preserving strong logical reasoning performance.

## 2 Recursively Verifiable NLI

In this work, we define an NLI instance  $i$  with a premise set  $P_i = \{p_1, p_2, \dots, p_n\}$  and a hypothesis  $h_i$  to be formally verifiable if the entailment condition can be verified by a theorem prover. Let  $\Phi : \text{NL} \rightarrow \mathcal{L}$  translate natural language sentences into a target logic  $\mathcal{L}$ , and let  $\mathcal{S}$  be a solver for  $\mathcal{L}$ . The NLI instance  $i$  is labelled as entailment if  $\mathcal{S}(\Phi(P_i) \vdash \Phi(h_i))$  is provable, and as contradiction if  $\mathcal{S}(\Phi(P_i) \vdash \neg\Phi(h_i))$  is provable.

Formal verifiability certifies only the endpoint relation between  $P_i$  and  $h_i$ . We call an instance recursively verifiable if there exists an explicit inference chain  $\Pi = \langle \delta_1, \delta_2, \dots, \delta_k \rangle$  whose intermediate steps are individually solver-checkable and compose to  $h_i$ . Let  $\Pi^{<j} = \langle \delta_1, \dots, \delta_{j-1} \rangle$  denote the prefix before step  $j$ . Then  $\Pi$  is a recursively verifiable witness if, for all  $j \in \{1, \dots, k\}$ ,  $\mathcal{S}(\Phi(P_i \cup \Pi^{<j}) \vdash \Phi(\delta_j))$  is provable, and  $\mathcal{S}(\Phi(P_i \cup \Pi) \vdash \Phi(h_i))$  is provable.

## 3 Methodology

To enable faithful formal checking of multi-step reasoning chains and to refine them when proof obligations are not discharged, we propose LLM-TP Tree, a neuro symbolic framework for recursively verifiable NLI based on an entailment reasoning tree  $T_i$ . Figure 1 provides an overview.

### 3.1 Entailment Tree Construction

Given a premise set  $P_i$  and a hypothesis  $h_i$ , we first prompt an LLM to construct an entailment tree and a set of intermediate conclusions  $C = \{\delta_1, \delta_2, \dots, \delta_m\}$ , optionally selecting a subset of premises  $P'_i \subseteq P_i$ , such that  $P'_i \cup C \models h_i$  holds.

The tree  $T_i$  consists of nodes labelled by propositions (premises or intermediate conclusions) and directed edges encoding entailment dependencies. For each internal node  $v \in T_i$  with proposition  $\varphi_v$  and descendant nodes  $\mathcal{D}(v)$ , the corresponding local entailment obligation requires that the conjunction of descendants entails the parent, namely  $\bigwedge_{u \in \mathcal{D}(v)} \varphi_u \models \varphi_v$ . The prompt of the entailment tree construction is in Appendix J.

### 3.2 Atomic Decomposition

Sentence complexity is a major contributor to both the latency and the brittleness of autoformalisation (Quan et al., 2025b). Long and structurally complex sentences tend to increase the search space of the translation process, which not only slows down inference but also amplifies semantic inconsistencies in the resulting logical forms. To mitigate these issues, we then decompose each natural language proposition into a set of atomic propositions, so that the subsequent autoformalisation operates on shorter, semantically focused units.

After constructing the entailment reasoning tree, we apply atomic decomposition to the propositions (first subtree) that participate in proof obligations. For the hypothesis  $h_i$ , we generate a set of atoms

$$D(h_i) = \{a_1^h, a_2^h, \dots, a_k^h\}, \quad h_i \equiv \bigwedge_{\ell=1}^k a_\ell^h. \quad (1)$$

Similarly, for the premise set  $P_i = \{p_1, \dots, p_n\}$ , each premise  $p_j$  is decomposed into

$$D(p_j) = \{a_1^{p_j}, a_2^{p_j}, \dots, a_{t_j}^{p_j}\}, \quad p_j \equiv \bigwedge_{\ell=1}^{t_j} a_\ell^{p_j}. \quad (2)$$

We define the global atom set induced by the premises as

$$D(P_i) := \bigcup_{j=1}^n D(p_j), \quad P_i \equiv \bigwedge_{a \in D(P_i)} a. \quad (3)$$

**Entailment-preserving requirement.** We require the decomposition operator  $D(\cdot)$  to be entailment-preserving: for any sentence  $\varphi$  and any atom  $a \in D(\varphi)$ , it must hold that  $\varphi \models a$ . Intuitively, each atom is a logically entailed consequence of the original sentence, so decomposition does not introduce new information that is not supported by the source. We obtain candidate

atoms by prompting an LLM (see Appendix I). To enforce the entailment-preserving requirement in practice, we apply a pretrained NLI classifier (roberta-large-mnli (Liu et al., 2019)) to verify that each atom is entailed by its source sentence within a predefined threshold (0.9). This filtering step reduces invalid decompositions and improves the stability of downstream autoformalisation.

### 3.3 $\theta$ -substitution Autoformalisation

Existing autoformalisation pipelines (Pan et al., 2023; Olausson et al., 2023) often rely on complex, monolithic generation procedures to map natural language sentences into fully specified logical formulas. This design is computationally inefficient and, more importantly, prone to systematic errors that frequently prevent theorem provers from certifying entailment. These observations motivate a more structured autoformalisation strategy that decomposes the translation into smaller, verifiable steps and reduces sensitivity to surface linguistic variation.

We then adopt the Neo-Davidsonian event semantics formulation used in Quan et al. (2024), where events and semantic roles (e.g., agent, patient) are represented explicitly in first-order logic to autoformalise the decomposed atomic propositions in previous step to logical forms. Building on this representation, we propose a multi-step  $\theta$ -substitution procedure that constructs the final logical form by progressively instantiating an abstract template (see Example 2).

A substitution  $\theta$  is a mapping from placeholders to terms:

$$\theta = \{x_1 \mapsto t_1, x_2 \mapsto t_2, \dots, x_n \mapsto t_n\}. \quad (9)$$

Applying  $\theta$  to a formula  $\varphi$ , denoted  $\varphi[\theta]$ , replaces each placeholder  $x_i$  with its corresponding term  $t_i$ .

For a natural language sentence, we start from an abstract logical template and apply a sequence of substitutions:

$$\varphi_0 \xrightarrow{\theta_1} \varphi_1 \xrightarrow{\theta_2} \varphi_2 \xrightarrow{\theta_3} \varphi_{\text{final}}. \quad (10)$$

Each substitution handles a specific aspect: 1) *Entity substitution* ( $\theta_1$ ): maps generic predicates to specific entities 2) *Event substitution* ( $\theta_2$ ): introduces event predicates and event structure. 3) *Role substitution* ( $\theta_3$ ): adds Neo-Davidsonian semantic roles and their arguments.

The final formula is obtained by composing the substitutions:

$$\varphi_{\text{final}} = \varphi_0[\theta_1 \circ \theta_2 \circ \theta_3]. \quad (11)$$

Example 2 illustrates the resulting stepwise derivation in detail.

### 3.4 Entailment tree verification and refinement for recursively verifiable NLI

Finally, given an entailment tree  $\mathcal{T}$ , we verify it recursively by turning each subtree into a set of local proof obligations that are checkable by a TP. Each node  $n \in \mathcal{T}$  stores a natural language statement  $s_n$ , its atomic decomposition  $D(n) = \{a_{n1}, a_{n2}, \dots, a_{n|D(n)|}\}$ , and a role indicator  $\rho(n) \in \{\text{EXPLANATION}, \text{HYPOTHESIS}\}$  that is defined *with respect to the current subtree under verification*. For a subtree  $\mathcal{T}_{\text{sub}} \subseteq \mathcal{T}$  with root  $h = \text{root}(\mathcal{T}_{\text{sub}})$  and leaf set  $E = \text{leaves}(\mathcal{T}_{\text{sub}})$ , we set

$$\rho(n) = \begin{cases} \text{EXPLANATION} & \text{if } n \in E, \\ \text{HYPOTHESIS} & \text{if } n = h. \end{cases} \quad (12)$$

**From atoms to axioms and lemmas.** Let  $\Phi$  denote autoformalisation from natural language atoms to the target logic. For every explanation node  $e \in E$ , each atom  $a \in D(e)$  is translated into an Isabelle/HOL (Nipkow et al., 2002) axiom (example in Appendix D):

$$\text{axiom}(e, a) : \Phi(a). \quad (13)$$

Collecting all explanation atoms yields the axiom set

$$A(E) := \bigcup_{e \in E} \{\Phi(a) \mid a \in D(e)\}. \quad (14)$$

For the hypothesis node  $h$ , each atom  $a \in D(h)$  induces a local proof obligation:

$$\text{lemma}(h, a) : A(E) \vdash \Phi(a). \quad (15)$$

We accept  $h$  as verified for the current subtree if *TP* certifies  $\text{lemma}(h, a)$  for all atoms  $a \in D(h)$ .

**Recursive verification with localized refinement.** The verification proceeds bottom-up. Intuitively, once a subtree root is verified, it can be treated as a support statement for its parent level. Algorithm 1 summarises the procedure.

When *TP* cannot discharge a lemma  $\text{lemma}(h, a)$  for some atom  $a \in D(h)$ , we

### Example 3.1: $\theta$ -substitution Autoformalisation

**Sentence:** A forest fire would cause deer to die or leave a woodland.

**Logic Template Generation:**

$$\varphi_0 = \forall x y z. P(x) \wedge Q(y) \wedge R(z) \rightarrow S \quad (4)$$

**Step 1: Entity Substitution** (Apply  $\theta_1 = \{P/\text{ForestFire}, Q/\text{Deer}, R/\text{Woodland}\}$ ):

$$\varphi_1 = \varphi_0[\theta_1] = \forall x y z. \text{ForestFire}(x) \wedge \text{Deer}(y) \wedge \text{Woodland}(z) \rightarrow S \quad (5)$$

**Step 2: Event Substitution** (Apply  $\theta_2 = \{S/(\text{Die}(e_1) \vee \text{Leave}(e_2))\}$ ):

$$\varphi_2 = \varphi_1[\theta_2] = \forall x y z. \text{ForestFire}(x) \wedge \text{Deer}(y) \wedge \text{Woodland}(z) \rightarrow (\text{Die}(e_1) \vee \text{Leave}(e_2)) \quad (6)$$

**Step 3: Semantic Role Substitution** (Apply  $\theta_3$ ):

$$\theta_3 = \begin{cases} \text{Die}(e_1) \mapsto (\text{Die}(e_1) \wedge \text{Agent}(e_1, y)) \\ \text{Leave}(e_2) \mapsto (\text{Leave}(e_2) \wedge \text{Agent}(e_2, y) \wedge \text{Patient}(e_2, z)) \end{cases} \quad (7)$$

**Final Result:**

$$\varphi_{\text{final}} = \varphi_2[\theta_3] = \forall x y z e_1 e_2. \text{ForestFire}(x) \wedge \text{Deer}(y) \wedge \text{Woodland}(z) \rightarrow (\text{Die}(e_1) \wedge \text{Agent}(e_1, y)) \vee (\text{Leave}(e_2) \wedge \text{Agent}(e_2, y) \wedge \text{Patient}(e_2, z)) \quad (8)$$

Figure 2:  $\theta$ -substitution autoformalisation example.

extract diagnostics from the prover (e.g., the failed goal and relevant constraints) and use them to localise the failure to a small set of implicated explanation nodes  $\hat{E}$ . We then prompt an LLM to refine only these implicated nodes, update the subtree accordingly, and re-check the same local obligations. This loop continues until all atoms in  $D(h)$  are certified for the subtree root  $h$ , after which the verified root is promoted and the algorithm proceeds to higher levels. Repeating this process bottom-up yields a recursively verifiable proof certificate for the entire entailment tree.

## 4 Experimental Setup

We evaluate our neuro-symbolic framework from two complementary perspectives: (i) explanation refinement apply theorem-prover checking and refine the explanation; and (ii) logical reasoning, which measures standard end-task accuracy to ensure that enhanced verification and refinement do not come at the cost of predictive performance.

**Datasets** We use four datasets ranging from synthetically generated deductive reasoning benchmarks to expert-written real-world corpora: ProofWriter (Tafjord et al., 2021), PrOntoQA (Saparov and He, 2023), FOLIO (Han et al., 2024), and EntailmentBank (Dalvi et al., 2021). Additional dataset details (i.e. number of samples) are provided in Appendix A.

**Baselines and Models** For explanation refinement, we compare against two state-of-the-

art explanation refinement models, Explanation-Refiner (Quan et al., 2024) and Faithful-Refiner (Quan et al., 2025b). For logical reasoning, we consider two prompting-based baselines, including a direct standard prompting and chain-of-thought (CoT) prompting. In addition, we compare against two representative neuro-symbolic frameworks, Logic-LM (Pan et al., 2023) and LINC (Olausson et al., 2023). We evaluate five distinct LLM backbones, spanning a general-purpose chat model (GPT-4o (OpenAI, 2024)), models evaluated in non-thinking mode (GPT-5 Nano (OpenAI, 2025), Grok-4 Fast (xAI, 2025)), and models evaluated in thinking mode (DeepSeek-V3.1 (Deepseek-AI, 2025), Qwen3-Max (Qwen Team, 2025)) and additionally include GPT-3.5 in logical reasoning tasks to match the model coverage reported in prior work. LLM implementation details is described in Appendix D.

## 5 Empirical Results and Evaluation

**LLM-TP Tree effectively to construct and refines verifiable explanation-based NLI.** Table 1 reports the main results on the explanation refinement task. Across all datasets and all LLM backbones, LLM-TP Tree achieves the best *Fin.* performance, indicating that it is consistently more effective at producing a checkable witness for explanation-based NLI. Averaged over backbones, LLM-TP Tree improves the final verified rate by 26.23%, 21.73%, 21.60%, and 48.93% on each dataset compared to Explanation-Refiner. The ef-

Dataset	Approach	GPT-4o		GPT-5 nano		Grok-4 fast		Deepseek-V3.1		Qwen3-max	
		Init.	Fin.	Init.	Fin.	Init.	Fin.	Init.	Fin.	Init.	Fin.
FOLIO	Explanation-Refiner	53.28	63.93	46.72	59.84	45.08	52.46	50.82	62.30	55.74	68.85
	Faithful-Refiner	68.03	78.69	53.28	68.03	54.92	65.57	71.31	81.15	69.67	84.43
	<b>LLM-TP Tree</b>	81.15	<b>90.16</b>	66.39	<b>81.15</b>	62.30	<b>79.51</b>	<b>92.62</b>	<b>95.08</b>	85.25	<b>95.08</b>
ProofWriter	Explanation-Refiner	62.67	76.67	59.33	73.33	53.33	69.33	52.67	66.67	64.00	78.00
	Faithful-Refiner	80.00	86.67	72.00	84.67	60.00	78.67	83.33	90.67	84.00	91.33
	<b>LLM-TP Tree</b>	91.33	<b>95.33</b>	85.33	<b>91.33</b>	72.67	<b>90.00</b>	<b>98.00</b>	<b>98.00</b>	92.00	<b>98.00</b>
PrOntoQA	Explanation-Refiner	64.00	78.67	60.00	74.67	61.33	73.33	62.00	74.00	64.67	78.00
	Faithful-Refiner	82.67	90.00	73.33	80.00	72.00	84.00	85.33	91.33	82.00	91.33
	<b>LLM-TP Tree</b>	92.00	<b>97.33</b>	93.33	<b>98.00</b>	86.00	<b>91.33</b>	<b>100.00</b>	<b>100.00</b>	98.00	<b>100.00</b>
EntailmentBank	Explanation-Refiner	15.33	26.67	12.00	23.33	8.67	13.33	15.33	26.67	16.67	28.00
	Faithful-Refiner	23.33	44.00	17.33	47.33	14.00	46.00	25.33	48.00	22.00	52.00
	<b>LLM-TP Tree</b>	21.33	<b>70.00</b>	18.00	<b>74.00</b>	16.67	<b>68.67</b>	<b>71.33</b>	<b>78.67</b>	22.67	<b>78.67</b>

Table 1: Comparison results on the explanation refinement tasks. *Init.* shows the number of explanation that is initially logically valid. *Fin.* shows the number of explanation that is finally logically valid after refinement. **Bold** values indicate the best performance. Arrows indicate absolute performance gain of LLM-TP Tree over the baseline.

Dataset	ER	FR	LT	LT w/o AD
FOLIO	4.27	3.75	3.04	3.57
ProofWriter	3.47	3.08	2.42	2.88
PrOntoQA	3.36	2.73	2.04	2.44
EntailmentBank	4.97	4.46	4.04	4.35
<b>Avg.</b>	4.02	3.51	2.88	3.31

Table 2: Average refinement iterations required averaged over five LLM backbones. Comparison between ER (Explanation-Refiner), FR (Faithful-Refiner), LT (LLM-TP Tree) and LT w/o AD (LT without atomic decomposition)

fect is most pronounced on EntailmentBank: while holistic refiners remain below 52% final validity, LLM-TP Tree reaches 68.67–78.67% across backbones. Beyond final refinement quality, the *Init.* column probes whether a backbone can produce an initially provable explanatory chain before any refinement. On synthetical benchmarks such as ProofWriter and PrOntoQA, LLM-TP Tree yields substantially higher initial validity than holistic refiners; for instance, on ProofWriter with Qwen3-max, LLM-TP Tree attains 92.00% *Init.* versus 64.00% for Explanation-Refiner. Notably, EntailmentBank exhibits low initial validity across methods, suggesting that constructing globally coherent verifiable chains is difficult under noisier multi-hop explanations; in this regime, our proposed mechanism yields the largest gains after refinement.

**Atomic decomposition improves the efficiency of verification-driven explanation refinement.**

We report the average number of refinement iterations, averaged across different LLM backbones, required to refine an explanation for each dataset, as shown in Table 2. Across all four datasets (FOLIO, ProofWriter, PrOntoQA, and EntailmentBank), LT converges in fewer iterations than baseline refiners, reducing the required iterations by an average of 1.14 relative to ER and by an average of 0.63 relative to Faithful-Refiner (FR). To isolate the advantages of atomic decomposition, we run an ablation that disables atomic sentence splitting. Relative to the full LT system, removing atomic decomposition increases the number of required refinement iterations by between 0.31 and 0.53. This ablation indicates that atomic decomposition contributes materially to LT’s efficiency gains: by splitting premises and intermediate conclusions into atomic subclaims, verification failures can be localised to a small set of fine-grained hypotheses, enabling targeted repairs of the minimal logical mismatch and thereby reducing both the number of refinement iterations and the cost of each verification/refinement pass. We also report the mean running time per refinement iteration (including theorem-proving time) under different LLM backbones, and provide the full per-backbone iteration results in Appendix F.

**Multi-step  $\theta$ -substitution autoformalisation improves semantic faithfulness.** A faithful autoformalisation should preserve the semantic content of the original sentence as much as possible. Fol-

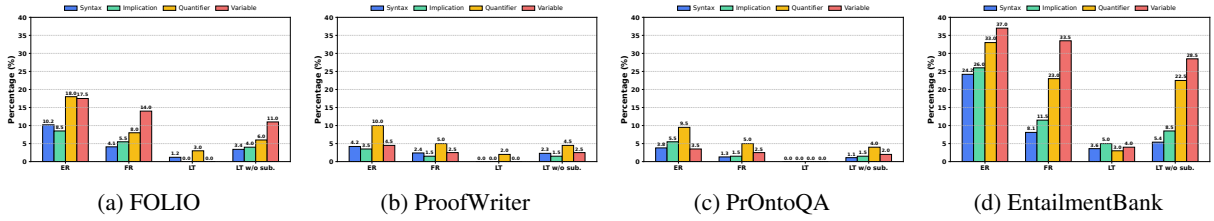


Figure 3: Distribution of the autoformalisation errors (syntax, implication, quantifier, variable in Isabelle/HOL Theories).

Dataset	ER	FR	LT	LT w/o sub.
FOLIO	0.766	0.822	0.856	0.828
ProofWriter	0.834	0.892	0.940	0.914
PrOntoQA	0.824	0.918	0.952	0.924
EntailmentBank	0.614	0.726	0.824	0.770
<b>Avg.</b>	0.760	0.840	0.893	0.859

Table 3: Average autoformalisation faithfulness (cosine similarity) averaged over five LLM backbones. Comparison between ER (Explanation-Refiner), FR (Faithful-Refiner), LT (LLM-TP Tree) and LT w/o sub. (LT without  $\theta$ -substitution autoformalisation).

lowing the rule-based informalisation procedure by Quan et al. (2025b), we quantify faithfulness by converting each predicted logical form back into natural language and computing the cosine similarity between the informalised sentence and the original sentence. Table 3 summarises the results. Our approach achieves the highest average faithfulness across all datasets, outperforming Explanation-Refiner by an average of 0.13 and Faithful-Refiner by an average of 0.053. As expected, faithfulness is generally higher on synthetic benchmarks, whose language more closely follows rule-based templates, whereas EntailmentBank is more challenging due to its more naturalistic multi-hop explanations. Notably, our method yields the largest margin on EntailmentBank, indicating stronger robustness of  $\theta$ -substitution autoformalisation under distributional and linguistic variability. We further perform an ablation by removing  $\theta$ -substitution (LT w/o sub.). This variant shows a consistent drop in similarity of an average of 0.034 across the same datasets compared to the full model. We also shows the full ablation study results in Appendix E

### 5.1 Failure and Error Analysis

We analyse refinement drift using proof-depth alignment. Figure 4 measure proof-depth alignment by comparing the TP-extracted used proof depth of each explanation with the dataset-provided gold depth (ideal  $y=x$ ). In refined cases, all frame-

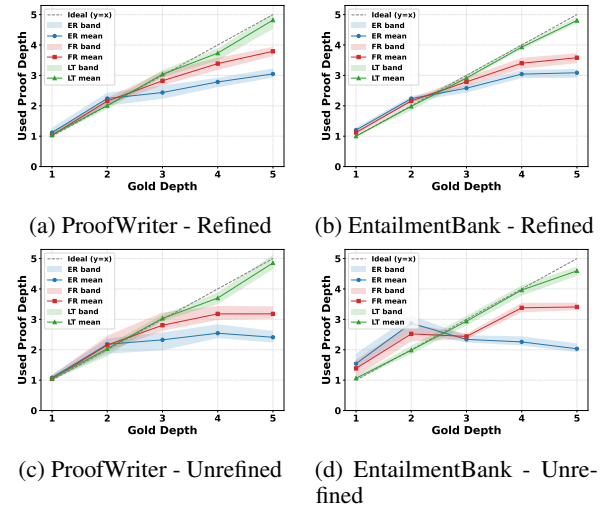


Figure 4: Proof depth alignment between gold and actual constructed proof depths across frameworks. The solid curve shows the mean used depth averaged over five LLMs, while the shaded band shows the range across backbones. Top: Trends in refined cases. Bottom: Trends in unrefined cases.

works show a positive trend that is close to the ideal  $y=x$  line. By contrast, in unrefined cases, ER and FR exhibit a clear refinement drift in the form of depth compression: the used depth grows slowly, saturates, or even decreases, increasingly undershooting the gold depth. This indicates that holistic refiners tend to collapse longer multi-step explanations into shorter proofs. In contrast, LLM-TP Tree tracks the gold depth much more closely suggesting that its refinement process better preserves proof-step granularity in longer-horizon inference. Detailed distributions are provided in Appendix H. We further analyse autoformalisation failures by categorising errors in the generated Isabelle/HOL theories into four types (details in Appendix B) across 200 sampled theories. Figure 3 shows that LLM-TP Tree consistently yields the lowest error rates across datasets, with the largest margin on the more naturalistic EntailmentBank. Removing  $\theta$ -substitution (LT w/o sub.) leads to a clear increase in semantic and structural errors, particularly quan-

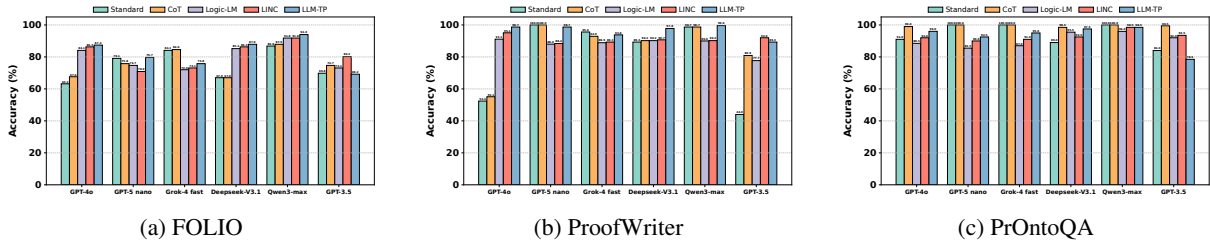


Figure 5: Comparison on the logical reasoning accuracy tasks.

tifier and variable errors on EntailmentBank, highlighting the importance of staged substitution autoformalisation for maintaining a prover-compatible logical structure.

## 5.2 Logical Reasoning

Figure 5 reports accuracy on three logical reasoning benchmarks. Across all the datasets and LLMs, LLM-TP Tree is the strongest neuro-symbolic method: compared against Logic-LM and LINC, it achieves higher accuracy. The improvements are substantial on the more challenging datasets. These results indicate that explicit theorem-proving with a tree-structured search provides a reliable verification signal and reduces errors relative to existing neuro-symbolic pipelines, as long as the backbone model can produce sufficiently correct autoformalisation. However, Neuro-symbolic approaches sometimes perform worse than standard prompting. Despite its advantage over other baselines, LLM-TP Tree is not universally better than direct prompting. In particular, on ProntoQA Standard/CoT prompting are consistently strongest across all backbones, leaving limited headroom for theorem-proving. This behaviour is consistent with two factors: (i) some datasets are already near-saturated by strong classifiers and CoT prompting, so additional symbolic steps cannot provide much benefit; (ii) neuro-symbolic methods introduce an additional failure mode (autoformalisation errors). When autoformalisation produces an incomplete or incorrect theory, the prover cannot establish either label and the system defaults to Unknown, directly harming accuracy. This is most evident with GPT-3.5, where LLM-TP Tree drops below LINC on all three datasets.

## 6 Related Work

Neuro-symbolic pipelines have shown strong effectiveness on multi-hop reasoning, as they combine learned language understanding with symbolic reasoning to support complex, multi-step inference.

Early work such as Weber et al. (2019) integrates a symbolic reasoner (Prolog) with a rule-learning component to solve multi-hop reasoning problems. More recently, researchers have increasingly combined LLMs with symbolic solvers across a range of logical reasoning tasks. For example, Pan et al. (2023) and Olausson et al. (2023) autoformalise natural-language contexts into symbolic representations and then apply external solvers to perform deductive reasoning over the resulting formal structures. Along a related direction, Xu et al. (2025a) proposes a decompose-search-resolve framework that breaks a problem into smaller logical substructures and performs structured search to resolve each reasoning step. Prover feedback has also been used to improve reliability: Quan et al. (2025a) extracts erroneous proof steps from solver feedback and iteratively refines reasoning in a loop. In contrast to endpoint-focused verification, our work targets recursively verifiable NLI by enforcing chain-level checking and localised refinement over structured entailment trees.

## 7 Conclusion

We presented LLM-TP Tree, a neuro-symbolic framework for recursively verifiable NLI that follows a decompose-and-formalise paradigm. By combining entailment-tree verification with entailment-preserving atomic decomposition and  $\theta$ -substitution autoformalisation, our method enables localised, refinement that targets the implicated part of an explanation rather than regenerating it holistically. Experiments on four datasets with multiple LLM backbones show improved TP-verified explanation rates, more faithful autoformalisation, and reduced refinement drift while maintaining strong end-task reasoning accuracy. In future work, we will extend this framework to broader naturalistic settings, improve the robustness of entailment-tree construction and diagnostic localisation, and further reduce verification cost through better reuse of verified subtrees.

## 553 Limitations

554 Our framework produces solver-checkable certifi-  
555 cates, but the guarantee is necessarily conditional  
556 on the chosen autoformalisation function  $\Phi$ , the  
557 target formalism, and the capabilities of the under-  
558 lying theorem prover. In particular, a verification  
559 can be valid in the induced formal theory while still  
560 being partially misaligned with the original natural  
561 language intent if  $\Phi$  introduces subtle scope, po-  
562 larity, or role mismatches that are not exposed by  
563 our current faithfulness checks. This dependence is  
564 amplified in naturalistic NLI, where missing back-  
565 ground knowledge, implicit commonsense assump-  
566 tions, or discourse phenomena (e.g., coreference,  
567 modality, temporality) may not be representable  
568 in the event-based logical form, making some in-  
569 stances difficult to certify without additional ax-  
570 ioms or richer semantics. The effectiveness of lo-  
571 calised refinement also depends on the quality of  
572 the intermediate structures proposed by the LLM  
573 and on the ability to map prover diagnostics back to  
574 the responsible region. When the initial entailment  
575 tree omits critical intermediate conclusions or intro-  
576 duces a misleading proof plan, bottom-up checking  
577 can localise which obligations are not discharged  
578 but may not reliably recover the globally correct  
579 chain without substantial restructuring. Finally,  
580 although our approach reduces expensive global  
581 regeneration, it still incurs non-trivial compute due  
582 to repeated LLM calls, repeated autoformalisation,  
583 and repeated prover invocations. Prover timeouts or  
584 inconclusive outcomes can remain a practical bot-  
585 tleneck on harder instances, and our current system  
586 does not provide a formal convergence guarantee  
587 to the minimal or unique refined explanation.

## 588 References

589 Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan  
590 Xie, Hannah Smith, Leighanna Pipatanangkura, and  
591 Peter Clark. 2021. Explaining answers with entail-  
592 ment trees. *EMNLP*.

593 Deepseek-AI. 2025. [Deepseek-v3 technical report](#).  
594 *Preprint*, arXiv:2412.19437.

595 Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhent-  
596 ing Qi, Martin Riddell, Wenfei Zhou, James Coady,  
597 David Peng, Yujie Qiao, Luke Benson, Lucy Sun,  
598 Alexander Wardle-Solano, Hannah Szabó, Ekaterina  
599 Zubova, Matthew Burtell, Jonathan Fan, Yixin Liu,  
600 Brian Wong, Malcolm Sailor, and 16 others. 2024.  
601 [FOLIO: Natural language reasoning with first-order  
602 logic](#). In *Proceedings of the 2024 Conference on  
603 Empirical Methods in Natural Language Processing*,

pages 22017–22031, Miami, Florida, USA. Associa-  
tion for Computational Linguistics.

Dongwei Jiang, Marcio Fonseca, and Shay Cohen. 2024.  
[LeanReasoner: Boosting complex logical reasoning  
with lean](#). In *Proceedings of the 2024 Conference of  
the North American Chapter of the Association for  
Computational Linguistics: Human Language Tech-  
nologies (Volume 1: Long Papers)*, pages 7497–7510,  
Mexico City, Mexico. Association for Computational  
Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-  
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,  
Luke Zettlemoyer, and Veselin Stoyanov. 2019.  
Roberta: A robustly optimized bert pretraining ap-  
proach. *arXiv preprint arXiv:1907.11692*.

Tobias Nipkow, Markus Wenzel, and Lawrence C Paul-  
son. 2002. *Isabelle/HOL: a proof assistant for  
higher-order logic*. Springer.

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, Singapore.  
Association for Computational Linguistics.

OpenAI. 2024. [Gpt-4 technical report](#). *Preprint*,  
arXiv:2303.08774.

OpenAI. 2025. Introducing gpt-5. <https://openai.com/index/introducing-gpt-5/>.

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, Singapore. Association for Computa-  
tional Linguistics.

Xin Quan, Marco Valentino, Danilo Carvalho, Dhairya  
Dalal, and Andre Freitas. 2025a. [PEIRCE: Unify-  
ing material and formal reasoning via LLM-driven  
neuro-symbolic refinement](#). In *Proceedings of the  
63rd Annual Meeting of the Association for Compu-  
tational Linguistics (Volume 3: System Demonstra-  
tions)*, pages 11–21, Vienna, Austria. Association for  
Computational Linguistics.

Xin Quan, Marco Valentino, Louise A. Dennis, and An-  
dre Freitas. 2024. [Verification and refinement of nat-  
ural language explanations through LLM-symbolic  
theorem proving](#). In *Proceedings of the 2024 Con-  
ference on Empirical Methods in Natural Language  
Processing*, pages 2933–2958, Miami, Florida, USA.  
Association for Computational Linguistics.

Xin Quan, Marco Valentino, Louise A. Dennis, and An-  
dre Freitas. 2025b. [Faithful and robust LLM-driven  
theorem proving for NLI explanations](#). In *Proceed-  
ings of the 63rd Annual Meeting of the Association*

604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659



765	challenging five-hop subset (fictional-entity reason-	513
766	ing). For explanation refinement, we randomly sam-	514
767	ple 150 instances from this subset. Since this	515
768	subset is binary in nature, we keep boolean-labelled	516
769	instances and cast each instance into the NLI-style	
770	format $(P, h)$ .	
771	<b>FOLIO.</b> FOLIO (Han et al., 2024) is a human-	518
772	curated first-order logic reasoning benchmark, pair-	519
773	ing natural-language statements with formal logical	520
774	structure to support theorem-prover-based verifi-	521
775	cation. For explanation refinement, we use the	522
776	evaluation split commonly adopted by prior neuro-	523
777	symbolic work after preprocessing, and remove	524
778	indeterminate labels to keep a boolean setting. Con-	525
779	cretely, we start from 182 evaluation instances after	526
780	filtering out 22 erroneous examples identified by	527
781	Olausson et al. (2023) and exclude those with un-	528
782	certain labels, yielding 122 instances. We then	529
783	treat the remaining instances as NLI pairs $(P, h)$	530
784	for verification and refinement.	531
785	<b>EntailmentBank.</b> EntailmentBank (Dalvi et al.,	532
786	2021) is a multi-hop entailment benchmark that pro-	533
787	vides naturalistic premises and hypotheses together	534
788	with structured explanation annotations (interme-	535
789	diate conclusions and entailment structure). For	536
790	explanation refinement, we sample 150 instances	537
791	and select examples spanning different gold proof	538
792	depths to ensure coverage of both shallow and	539
793	longer-horizon multi-step chains. Each instance	540
794	is cast into $(P, h)$ , where $P$ is the set of supporting	541
795	statements (facts) and $h$ is the hypothesis.	542
796	<b>A.2 Logical Reasoning</b>	543
797	<b>ProofWriter.</b> For logical reasoning, we evalu-	544
798	ate on the ProofWriter D5 split under the open-	545
799	world assumption to retain the three-way label	546
800	space (True/False/Unknown). We sample 225	547
801	instances stratified by gold reasoning depth $d \in$	548
802	$\{1, 2, 3, 4, 5\}$ , with 45 instances per depth. Within	549
803	each depth bucket, we balance labels by sampling	550
804	15 instances per class.	551
805	<b>FOLIO.</b> For logical reasoning, the original vali-	552
806	dation split contains 204 instances, while we select	553
807	the 182 evaluation instances after filtering out 22	554
808	erroneous examples identified by Olausson et al.	555
809	(2023) for evaluation. We evaluate on this 182-	556
810	instance set under the dataset’s original label space.	557
811	<b>PrOntoQA.</b> For logical reasoning, we follow	558
812	prior work (Pan et al., 2023) and evaluate on the	559
	five-hop PrOntoQA subset. We randomly sam-	560
	ple 200 instances from this subset (binary labels),	561
	which represents the most challenging multi-hop	
	configuration in the benchmark.	
	<b>B Autoformalisation Errors</b>	
	Following the error taxonomy in Quan et al.	562
	(2025b), we categorise autoformalisation errors	563
	into four broad types: <i>syntax</i> , <i>implication</i> , <i>quanti-</i>	564
	<i>fier</i> , and <i>variable</i> errors. We define an autoformali-	565
	sation error as any mismatch between the intended	566
	logical meaning of the input sentence(s) (as cap-	567
	tured by a reference logic template or gold formalis-	568
	ation) and the produced formal statement. Notably,	569
	these categories differ in how easily they can be	570
	detected: syntax errors are often surfaced immedi-	571
	ately by the theorem prover, whereas the remaining	572
	three categories typically produce well-typed, well-	573
	formed formulas that may still be semantically un-	574
	faithful and therefore require semantic-level check-	575
	ing.	576
	<b>Syntax errors.</b> Syntax errors are those that pre-	577
	vent the proof assistant from parsing or type-	578
	checking the generated statement, and are there-	579
	fore usually the easiest to detect. Typical causes	580
	include missing brackets, malformed binder syntax,	581
	incorrect use of infix operators, and type unifica-	582
	tion conflicts (e.g., using an event-typed variable	583
	where an entity is expected). Such errors are often	584
	directly indicated by the theorem prover’s feedback	585
	(parse errors, unknown constants, failed unification,	586
	or unsatisfied type-class constraints), making them	587
	amenable to automatic filtering.	588
	<b>Implication errors.</b> Implication errors occur	589
	when the generated formula has a valid syntax and	590
	type, but its <i>logical structure</i> encodes the wrong	591
	entailment relation compared to the intended tem-	592
	plate. This includes incorrect direction of implica-	593
	tion, incorrect connective choice (e.g., $\wedge$ vs. $\rightarrow$ ), or	594
	wrong nesting/precedence that changes the mean-	595
	ing. These errors are difficult for theorem provers	596
	to flag because the statement may still be perfectly	597
	consistent and provable under some background	598
	axioms, while representing a different claim.	599
	<b>Quantifier errors.</b> Quantifier errors arise when	600
	the autoformalisation introduces incorrect quanti-	601
	fiers ( $\forall$ vs. $\exists$ ), incorrect quantifier scope, or incor-	602
	rect ordering of quantifiers. These errors are es-	603
	pecially harmful because they often produce state-	604
	ments that are strictly stronger or weaker than in-	605

---

**Algorithm 1** VerifySubtree( $\mathcal{T}_{\text{sub}}, TP$ )

---

```
1:  $h \leftarrow \text{root}(\mathcal{T}_{\text{sub}})$  {current hypothesis node}
2:  $E \leftarrow \text{leaves}(\mathcal{T}_{\text{sub}})$  {current explanation nodes}
3: Construct  $A(E) = \bigcup_{e \in E} \{\Phi(a) \mid a \in D(e)\}$  {axioms}
4: for each atom  $a \in D(h)$  do
5:   Define the goal  $\tau \leftarrow \Phi(a)$ 
6:   if  $TP.\text{prove}(A(E) \vdash \tau) = \text{FAIL}$  then
7:      $\text{diag} \leftarrow TP.\text{diagnostics}()$ 
8:     Identify implicated explanation nodes  $\hat{E} \subseteq E$  from
        $\text{diag}$ 
9:     Refine  $\hat{E}$  using an LLM conditioned on  $(\hat{E}, \text{diag})$ 
10:    Update  $\mathcal{T}_{\text{sub}}$ ; recompute  $D(\cdot)$  and re-autoformalise
       $\Phi(\cdot)$ 
11:    return VerifySubtree( $\mathcal{T}_{\text{sub}}, TP$ ) {retry locally}
12:   end if
13: end for
14: Mark  $h$  as verified
15: return  $h$  {promote as support for parent level}
```

---

862 tended while remaining syntactically valid. More-  
863 over, theorem provers typically do not report quan-  
864 tifier mismatches as errors; the formula is well-  
865 formed, yet its truth conditions can differ dramati-  
866 cally.

867 **Variable errors.** Variable errors occur when log-  
868 ical variables are incorrectly assigned to predicate  
869 arguments, resulting in role confusion or broken  
870 discourse links. This category covers a wide range  
871 of semantic binding mistakes, including swapping  
872 semantic roles (e.g., agent vs. patient), incorrectly  
873 linking pronouns/coreference, using the wrong  
874 event/entity variable in a clause, or mistakenly  
875 reusing a variable that should be distinct. Vari-  
876 able errors may be subtle because they can remain  
877 type-correct (e.g., swapping two entity variables  
878 still type-checks) yet invert the intended meaning.

## 879 C Algorithm

### 880 D Isabelle/HOL and LLMs 881 Implementation Detail

882 We use the `isabelle-client` Python pack-  
883 age (Shminke, 2022) to run Isabelle/HOL as a real-  
884 time server and to retrieve prover responses. We  
885 access all LLM backbones via API calls with the  
886 following model versions: GPT-4o (`gpt-4o`), GPT-  
887 5 nano (`gpt-5-nano`), Grok-4 fast (`grok-4-fast`),  
888 Deepseek-V3.1 (DeepSeek-v3.1-terminus), and  
889 Qwen3-Max (`qwen3-max-preview`). For all non-  
890 thinking models, we set the temperature to 0. We  
891 also do not limit any thinking effort for Deepseek-  
892 V3.1 and Qwen3-Max.

893 We build the Isabelle/HOL theory following the  
894 approach in (Quan et al., 2025b). An example of

Dataset	LT		w/o Atom		w/o $\theta$	
	Init.	Fin.	Init.	Fin.	Init.	Fin.
<b>Backbone: GPT-4o</b>						
FOLIO	81.15	90.16	71.31	86.07	75.41	88.52
			$\downarrow 9.84$	$\downarrow 4.09$	$\downarrow 5.74$	$\downarrow 1.64$
			87.33	92.00	89.33	94.67
ProofWriter	91.33	95.33	$\downarrow 4.00$	$\downarrow 3.33$	$\downarrow 2.00$	$\downarrow 0.66$
			86.67	92.67	90.00	94.00
PrOntoQA	92.00	97.33	$\downarrow 5.33$	$\downarrow 4.66$	$\downarrow 2.00$	$\downarrow 3.33$
			17.33	53.33	19.33	61.33
EntailmentBank	21.33	70.00	$\downarrow 4.00$	$\downarrow 16.67$	$\downarrow 2.00$	$\downarrow 8.67$
<b>Backbone: Qwen3-max</b>						
FOLIO	85.25	95.08	72.13	86.89	79.51	90.98
			$\downarrow 13.12$	$\downarrow 8.19$	$\downarrow 5.74$	$\downarrow 4.10$
			88.00	93.33	88.00	96.00
ProofWriter	92.00	98.00	$\downarrow 4.00$	$\downarrow 4.67$	$\downarrow 4.00$	$\downarrow 2.00$
			86.67	94.00	92.00	96.00
PrOntoQA	98.00	100.00	$\downarrow 11.33$	$\downarrow 6.00$	$\downarrow 6.00$	$\downarrow 4.00$
			18.00	61.33	20.00	71.33
EntailmentBank	22.67	78.67	$\downarrow 4.67$	$\downarrow 17.34$	$\downarrow 2.67$	$\downarrow 7.34$

Table 4: Ablations of LLM-TP Tree on GPT-4o and Qwen3-max (Init/Fin verified rate, %). Red arrows indicate the absolute decrease compared to LT under the same backbone.

895 the constructed Isabelle/HOL theory is shown in  
896 Figure 6. We then apply atomic decomposition to  
897 each axiom and to the hypothesis. The decomposed  
898 hypothesis sentence and the decomposed axioms  
899 are then forming a new theory, which we use to  
900 prove theorems from each decomposed hypothesis  
901 atoms.

## 902 E Ablation Study on Explanation 903 Refinement Task

904 Table 4 reports ablations of LLM-TP Tree (LT) on  
905 GPT-4o and Qwen3-max, using *Init./Fin.* logically  
906 valid explanation rates. Removing either of the two  
907 proposed components, atomic decomposition or  
908 multi-step  $\theta$ -substitution autoformalisation, consis-  
909 tently degrades across all datasets. Across datasets,  
910 *Init.* decreases by 4.00 to 13.12 points and *Fin.*  
911 decreases by 3.33 to 17.34 points, with the most  
912 pronounced degradation on EntailmentBank. Re-  
913 moving  $\theta$ -substitution also leads to a consistent  
914 drop, although smaller than w/o Atom indicating  
915 that multi-step substitution contributes materially  
916 to producing prover-compatible formalisations.

## 917 F Full results: Model efficiency Across 918 LLM Backbones

919 Figure 7a, 7b and 7c shows average number of  
920 iterations required to refine an explanation across  
921 all LLMs. Figure 7e, 7f, 7g and 8 shows average  
922 inference time per iteration in both thinking and  
923 non-thinking models.

## 924 G Full Results: Autoformalisation 925 Faithfulness Across LLM Backbones

926 Figure 9 compares the average faithfulness of  
927 the autoformalisation process across different ap-

```

theory data_0
imports Main

begin

typedecl entity
typedecl event

consts
  Melting :: "event  $\Rightarrow$  bool"
  Change  :: "event  $\Rightarrow$  bool"
  Source  :: "event  $\Rightarrow$  entity  $\Rightarrow$  bool"
  Destination :: "event  $\Rightarrow$  entity  $\Rightarrow$  bool"
  Solid   :: "entity  $\Rightarrow$  bool"
  Liquid  :: "entity  $\Rightarrow$  bool"
  IncreaseHeatEnergy :: "entity  $\Rightarrow$  bool"
  By      :: "event  $\Rightarrow$  entity  $\Rightarrow$  bool"
  Chocolate :: "entity  $\Rightarrow$  bool"
  Melts   :: "event  $\Rightarrow$  bool"
  Agent   :: "event  $\Rightarrow$  entity  $\Rightarrow$  bool"
  In      :: "event  $\Rightarrow$  entity  $\Rightarrow$  bool"
  Sunlight :: "entity  $\Rightarrow$  bool"

(* Explanation 1: melting means changing from a solid to a liquid by increasing heat energy *)
axiomatization where
  explanation_1: " $\forall$  e x y z. Melting e  $\leftrightarrow$  (Change e  $\wedge$  Source e x  $\wedge$  Destination e y  $\wedge$  Solid x  $\wedge$ 
    Liquid y  $\wedge$  IncreaseHeatEnergy z  $\wedge$  By e z)"

(* Explanation 2: chocolate melts in the sunlight *)
axiomatization where
  explanation_2: " $\exists$  x e. Chocolate x  $\wedge$  Melts e  $\wedge$  Agent e x  $\wedge$  In e Sunlight"

theorem hypothesis:
  assumes asm: "Chocolate x  $\wedge$  Solid y  $\wedge$  Liquid z"
  (* Hypothesis: chocolate changes from a solid to a liquid in the sunlight *)
  shows False
  sledgehammer
  oops

end

```

Figure 6: An example of the constructed Isabelle/HOL Theory.

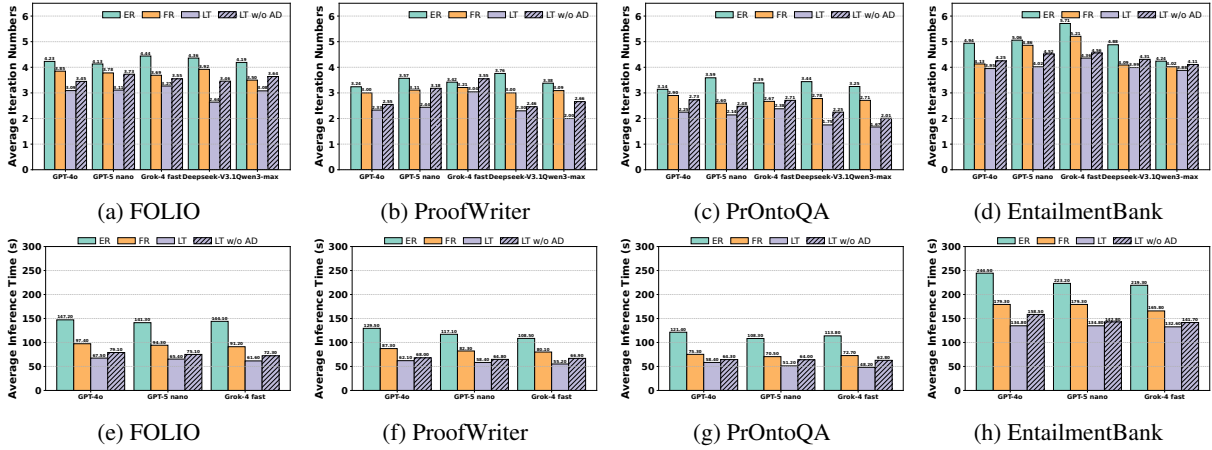


Figure 7: Top: Average number of iterations required to refine an explanation. Bottom: Average inference time per iteration. Comparison between ER (Explanation-Refiner), FR (Faithful-Refiner), LT (LLM-TP Tree) and LT w/o AD (LT without atomic decomposition).

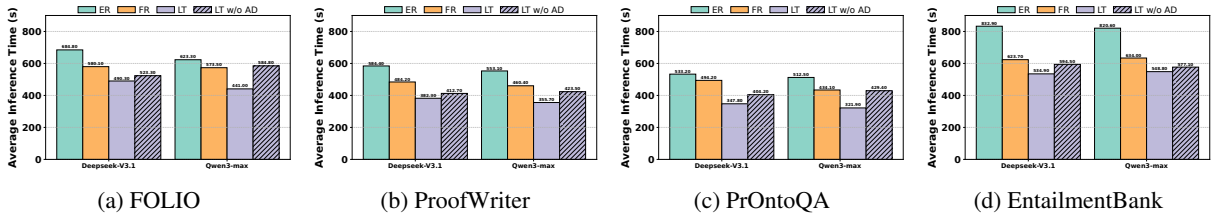


Figure 8: Average inference time per iteration in thinking mode LLMs. Comparison between ER (Explanation-Refiner), FR (Faithful-Refiner), LT (LLM-TP Tree) and LT w/o AD (LT without atomic decomposition).

proaches, averaged over all LLM backbones.

## H Solver Proof Depths

Table 5 and Table 6 shows Per-LLM depth alignment details for ProofWriter and EntailmentBank in refined and unrefined cases. Each cell reports the mean used proof depth averaged over instances within the corresponding gold depth bucket, for a fixed framework and LLM backbone.

## I Atomic Decomposition Prompt

We build the atomic decomposition prompt based on Srikanth and Rudinger (2025) as:

**SYSTEM:** You are an expert in symbolic/logical reasoning, linguistic and natural language inference. You are given a sentence and its logical information. Generate a list of atomic facts that are strictly logically entailed from the given sentence.

**Instructions:**

1. Keep each fact independent and self-contained.
2. Each fact should make sense when read on its own.
3. Only write facts that are directly described or supported by the sentence.
4. The atomic facts must be logically entailed from the given sentence.

**USER:** Here are some examples:

**Example:**

**Provided Sentence:**

Professional actors are in a summer performance.

**Answer:**

Atom 1: The people are professional.

Atom 2: The people are actors.

Atom 3: The performance is during summer.

**Task:**

**Provided Sentence:** {sentence}

**Answer:**

## J Entailment Tree Construction Prompt

**SYSTEM:** You are an expert in natural language inference and textual entailment. Given the following premise sentences and a final conclusion, generate some step-by-step intermediate conclusion sentences to finally infer the final conclusion.

**Instructions:**

1. The intermediate conclusion sentences must be strictly logical entailed from the premise sentences.
2. Since the reasoning is step-by-step, the intermediate conclusion sentences can be the new premise sentences for next step.
3. One intermediate conclusion sentence can be generated from multiple premise sentences.
4. If the final conclusion can be directly inferred from the premise sentences, you need to state the intermediate

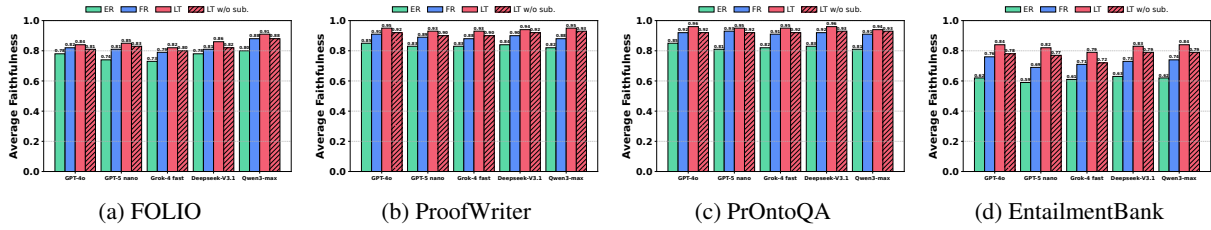


Figure 9: Comparison of the average faithfulness of the autoformalisation process across different approaches.

(a) ProofWriter: mean used proof depth by gold depth						(b) EntailmentBank: mean used proof depth by gold depth							
Framework	LLM	Gold depth $d$					Framework	LLM	Gold Depths $d$				
		1	2	3	4	5			1	2	3	4	5
Explanation-Refiner	GPT-4o	1.00	2.32	2.51	2.92	3.23	Explanation-Refiner	GPT-4o	1.23	2.14	2.65	3.15	3.21
	GPT-5 nano	1.14	2.00	2.67	2.61	2.94		GPT-5 nano	1.21	2.23	2.43	2.91	2.94
	Grok-4 fast	1.08	2.04	2.53	2.62	2.92		Grok-4 fast	1.13	2.21	2.52	2.93	2.92
	Deepseek-V3.1	1.28	2.41	2.23	2.87	3.04		Deepseek-V3.1	1.16	2.28	2.67	3.04	3.16
	Qwen3-max	1.10	2.43	2.24	2.90	3.11		Qwen3-max	1.28	2.32	2.63	3.18	3.20
Faithful-Refiner	GPT-4o	1.00	2.26	2.76	3.43	3.94	Faithful-Refiner	GPT-4o	1.12	2.10	2.87	3.33	3.67
	GPT-5 nano	1.02	2.04	3.13	3.34	3.63		GPT-5 nano	1.11	2.12	2.64	3.22	3.40
	Grok-4 fast	1.14	2.09	2.67	3.18	3.87		Grok-4 fast	1.17	2.13	2.72	3.58	3.42
	Deepseek-V3.1	1.00	2.17	2.78	3.41	3.73		Deepseek-V3.1	1.03	2.25	2.78	3.40	3.69
	Qwen3-max	1.00	2.22	2.77	3.58	3.79		Qwen3-max	1.17	2.20	2.93	3.47	3.73
LLM-TP Tree	GPT-4o	1.00	1.99	2.94	3.76	4.72	LLM-TP Tree	GPT-4o	1.00	2.10	2.85	3.90	4.81
	GPT-5 nano	1.08	2.12	2.98	3.68	4.53		GPT-5 nano	1.00	1.84	2.84	3.86	4.71
	Grok-4 fast	1.17	1.89	3.18	3.49	5.03		Grok-4 fast	1.00	2.04	2.92	3.88	4.73
	Deepseek-V3.1	1.00	2.03	3.08	3.84	4.94		Deepseek-V3.1	1.00	1.97	3.03	4.03	4.85
	Qwen3-max	1.00	2.00	3.00	3.89	4.88		Qwen3-max	1.00	1.96	2.94	3.98	4.92

Table 5: Per-LLM depth alignment details for ProofWriter and EntailmentBank in refined cases.

conclusion sentences are empty.

5. There might be redundant premise sentences.

**USER:** Here are some examples:

**Example:**

Initial Premises:

1. Monkeypox is an infectious disease caused by the monkeypox virus.
2. Monkeypox virus can occur in certain animals, including humans.
3. Humans are mammals.
4. Mammals are animals.
5. Symptoms of Monkeypox include fever, headache, muscle pains, feeling tired, and so on.
6. People feel tired when they get a flu.

Final Conclusion:

There is an animal.

Answer:

Monkeypox is an infectious disease caused by the monkeypox virus.

Monkeypox virus can occur in certain animals, including humans.

Conclusion: Humans can get monkeypox.

Humans are mammals.

Mammals are animals.

Conclusion: Humans are animals.

...

Humans are animals.

Therefore, there is an animal (humans) that can get monkeypox and feel tired. Conclusion: Therefore, there is an animal.

Initial Premises:

Final Conclusion:

Answer:

(a) ProofWriter: mean used proof depth by gold depth						(b) EntailmentBank: mean used proof depth by gold depth							
Framework	LLM	Gold depth $d$					Framework	LLM	Gold Depths $d$				
		1	2	3	4	5			1	2	3	4	5
Explanation-Refiner	GPT-4o	1.00	2.34	2.49	2.57	2.33	Explanation-Refiner	GPT-4o	1.53	2.84	2.48	2.44	2.21
	GPT-5 nano	1.12	1.99	2.57	2.46	2.24		GPT-5 nano	1.39	2.99	2.33	2.31	1.95
	Grok-4 fast	1.10	1.87	2.48	2.38	2.32		Grok-4 fast	1.43	2.83	2.24	2.21	1.92
	Deepseek-V3.1	1.18	2.37	1.98	2.84	2.54		Deepseek-V3.1	1.48	2.58	2.33	2.14	1.96
	Qwen3-max	1.00	2.33	2.11	2.47	2.62		Qwen3-max	1.88	3.12	2.32	2.19	2.12
Faithful-Refiner	GPT-4o	1.02	2.16	2.73	3.44	3.12	Faithful-Refiner	GPT-4o	1.34	2.42	2.48	3.23	3.55
	GPT-5 nano	1.02	1.95	3.23	3.24	3.02		GPT-5 nano	1.33	2.32	2.41	3.32	3.30
	Grok-4 fast	1.16	2.46	2.62	3.09	3.06		Grok-4 fast	1.49	2.28	2.43	3.54	3.42
	Deepseek-V3.1	1.00	1.98	2.78	2.98	3.27		Deepseek-V3.1	1.10	2.85	2.34	3.38	3.45
	Qwen3-max	1.00	2.22	2.67	3.15	3.43		Qwen3-max	1.67	2.74	2.53	3.43	3.42
LLM-TP Tree	GPT-4o	1.02	2.13	3.01	3.72	4.72	LLM-TP Tree	GPT-4o	1.03	2.13	2.83	3.93	4.62
	GPT-5 nano	1.09	2.14	2.95	3.54	4.65		GPT-5 nano	1.13	1.86	2.80	3.78	4.51
	Grok-4 fast	1.15	1.79	3.16	3.50	5.13		Grok-4 fast	1.12	2.04	2.95	3.91	4.43
	Deepseek-V3.1	1.02	2.13	3.04	3.86	4.87		Deepseek-V3.1	1.03	1.89	3.12	4.12	4.67
	Qwen3-max	1.00	1.98	3.00	3.88	4.91		Qwen3-max	1.00	2.01	2.99	4.08	4.74

Table 6: Per-LLM depth alignment details for ProofWriter and EntailmentBank in unrefined cases.

## K $\theta$ -substitution Autoformalisation Prompt

Answer:

**SYSTEM:** You are an expert in symbolic/logical reasoning and autoformalisation. You are given a sentence. Extract the logical relation from a given natural language sentence and construct a logical template representing its structure, following the format demonstrated in the examples.

### Instructions:

1. Read and understand the provided sentence.
2. Identify key entities, relationships, and logical structure (e.g., quantifiers, conjunctions, disjunctions, implications).
3. Abstract the sentence into a logical template using placeholders (e.g.,  $\forall x y z. P(x) \wedge Q(y) \wedge R(z) \rightarrow S$ ) that capture the logical form, not the specific content.

Follow these steps:

1. Parse the sentence and reason step-by-step about its logical components: What are its predicates, quantifiers, variables, and connectives?
2. Formulate the logical template, using generic predicate letters (P, Q, R, S, etc.) and variable placeholders (x, y, z, etc.), mirroring the schema shown in the examples.

...

Example 1:

Sentence: If someone wins the lottery, they will buy a new house or a car.

Example 1 Answer:

Logical Template:  $\forall x y z. P(x) \wedge Q(y) \rightarrow (R(z) \vee S(z))$

...

**Provided Sentence:** {sentence}