

UCert: A Formal Certificate of Unexplainability for TAG Neural Networks

Anonymous ACL submission

Abstract

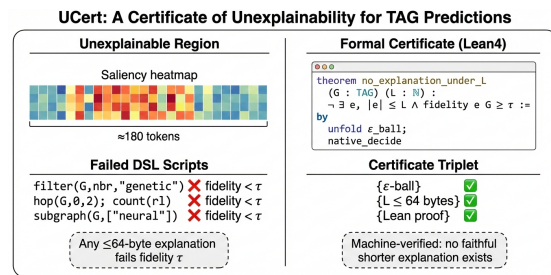
Natural language explanations for classifiers operating on text-attributed graphs present an inherent tension: explanations should be both readable to humans and verifiable by machines. We introduce UCert, a unified framework that places multiple explanation paradigms within a single evaluation and generation interface and that supports a certificate-based reverse-explanation mode. The framework combines saliency-driven pseudo-labeling, iterative refinement for fluent explanations, adversarial probing to compress evidence into a concise bottleneck, and an executable domain-specific language for machine-checkable rationales. UCert additionally issues formally checkable Unexplainability Certificates which assert that no short symbolic explanation can reproduce a classifier’s behavior within a defined perturbation model. We demonstrate that this hybrid design clarifies practical trade-offs between readability, automated faithfulness proxies, and machine-checkable guarantees for text-attributed graph models. Our code is available at: <https://anonymous.4open.science/r/tacer-76B5/README.md>.

Keywords: text-attributed graphs, graph neural networks, explanation, certifiable guarantees, natural language explanations

1 Introduction

Graph neural networks are the prevailing approach for learning on relational data and are widely used in domains where nodes and edges carry substantial textual content, such as citation networks, legal retrieval systems, and recommendation graphs (Tang et al., 2024; Zhang et al., 2024a). In text-attributed graphs the evidence supporting a prediction frequently spans many nodes and long documents, which complicates the generation of concise, human-facing explanations. Standard instance-level explainers return node, edge, or feature attributions and these outputs remain valuable

for model analysis but are often verbose and hard to interpret when text features dominate the signal (Yuan et al., 2022; Amara et al., 2022). Several works have proposed token- or feature-level attributions to expose textual evidence, but these fine-grained scores are difficult to integrate into compact summaries when the relevant evidence is distributed across subgraphs (Zhao et al., 2023; Si et al., 2023). Natural language explanations



UCert converts “explanation failure” into a formal statement: within an ϵ -ball around the instance, no explanation of length $\leq L$ satisfies the fidelity threshold τ , and produces a machine-checkable Lean4 proof.

Figure 1: UCert certificate example: **Left** shows saliency heat-map and rejected DSL snippets; **Right** shows the verified certificate triple (ϵ -ball, length ≤ 64 bytes, Lean4 proof).

promise improved accessibility by summarizing multi-node evidence into coherent rationales, but free-form narratives risk being unfaithful to the underlying model unless they are constrained by evidence selection or verification procedures (Cedro and Martens, 2025). Recent trends combine large language models with TAG representation learning to better encode textual attributes and to produce more fluent explanations, yet the integration also raises new faithfulness and verification challenges (Qin et al., 2023; Zhu et al., 2025; Su et al., 2025). Complementary approaches have sought mechanically checkable explanations or formal certificates for robustness and correctness; such techniques provide stronger guarantees but typically trade away readability and require specialized execution semantics and tooling (Sabanayagam et al., 2024;

069	Singh et al., 2024; Saad and Sharma, 2024).		
070	Motivated by these trade-offs, we develop a uni-	proaches extract compact explanatory subgraphs	119
071	fied framework that accommodates three explana-	via mask optimisation (Ying et al., 2019), learn	120
072	tion paradigms and a certificate-oriented verifica-	probabilistic edge masks (Luo et al., 2020), or	121
073	tion pipeline. The first paradigm produces read-	search for high-value subgraphs under learned scor-	122
074	able natural-language explanations by verbalizing	ers (Yuan et al., 2021). Other contributions elevate	123
075	saliency-derived subgraphs and refining pseudo-	low-level attributions into human-interpretable con-	124
076	labels through iterative screening. The second	cepts or incorporate human feedback to improve	125
077	paradigm enforces a concise information bottle-	usability (Magister et al., 2021). Several works ad-	126
078	neck using adversarial probing and a generator-	dress robustness and causal validity by combining	127
079	critic objective so that produced rationales are com-	adversarial analyses with causal modelling to sepa-	128
080	compact and better aligned with automated faithful-	rate spurious correlations from structural evidence	129
081	ness proxies. The third paradigm emits short, exe-	(Behnam and Wang, 2024; Huang et al., 2023; Zhai	130
082	cutable scripts in a restricted domain-specific lan-	and Zhang, 2022).	131
083	guage so that explanations can be mechanically		
084	checked in a sandboxed environment. To reconcile	2.2 Natural Language and LLM-based	132
085	these modalities we introduce the Unexplainability	Explanations	133
086	Certificate. Under a declared perturbation model		
087	and a bound on symbolic description length, the	Translating structured graph evidence into fluent	134
088	certificate asserts that every candidate short sym-	text has attracted growing interest because narra-	135
089	bolic rationale is falsifiable by a perturbation in the	tives are easier for lay users to consume. Prior	136
090	permitted neighbourhood; the certificate itself is	studies that extract textual rationales and evaluate	137
091	a machine-checkable artifact that can be independ-	their faithfulness establish important metrics and	138
092	ently verified.	highlight pitfalls (Chen et al., 2023; Amara et al.,	139
093		2022). Recent pipelines condition large language	140
094	Our contributions are as follows. Our frame-	models on model-derived rationales to generate	141
095	work unifies free-text, bottlenecked, and executable	readable explanations, as exemplified by Graph-	142
096	explanation modalities under a single formal in-	Narrator (Pan et al., 2025). Analyses show fluent	143
097	terface so that different outputs can be compared	narratives can remain unfaithful unless constrained	144
098	by common proxies and by mechanical verifica-	by explicit evidence selection or verification, mo-	145
099	tion. We introduce adversarial probing and a DSL-	tivating hybrid designs that pair LM proxies with	146
100	based execution path as principled alternatives to	mechanistic checks (Fayyaz et al., 2024; Li et al.,	147
101	unconstrained natural-language narratives and char-	2025).	148
102	acterize how these alternatives change trade-offs		
103	among readability, faithfulness proxies, and run-	2.3 Executable, Causal, and Interactive	149
104	time cost. We formalize Unexplainability Certifi-	Explanations	150
105	cates as a reverse-explanation construct that pro-		
106	vides machine-checkable guarantees about the non-	A parallel strand advocates for explanations that	151
107	existence of short faithful symbolic rationales un-	are mechanically verifiable or actionable rather	152
108	der a specified perturbation model. We provide an	than purely descriptive. Programming-language	153
109	evaluation protocol that pairs LM-based faithful-	and DSL approaches make reasoning steps explicit	154
110	ness proxies with perturbation-driven falsification	and testable, enabling executable explanations with	155
111	checks and certificate verification so practitioners	clear semantics (Jeon et al., 2024; Tegeler et al.,	156
112	can choose explanation strategies matched to their	2022; Saad and Sharma, 2024). Causal and coun-	157
113	assurance needs.	terfactual methods operationalise interventions to	158
114		support what-if analyses and recourse (Jaimini and	159
115	2 Related Work	Sheth, 2022; Xiao et al., 2023; Huang et al., 2023).	160
116		Interactive frameworks reveal model behaviour	161
117	2.1 Explainability for Graph Neural	through sequential probes or evidence-collection	162
118	Networks	policies, but these methods require sandboxed ex-	163
	A large literature studies instance-level explana-	ecution and formal intervention models (Behnam	164
	tions that identify influential nodes, edges, or fea-	and Wang, 2024).	165
	tures for a GNN prediction. Representative ap-		

2.4 Evaluation, Benchmarks, and Certification

Designing evaluation protocols for graph explanations remains challenging because many automated proxies misalign with human judgements. Benchmarking efforts combine robustness checks, LM proxies, and user studies to provide more comprehensive assessments (Amara et al., 2022; Agarwal et al., 2023). Work on certification and trust develops formal guarantees and clarifies deployment assumptions for safety-critical settings (Kwiatkowska and Zhang, 2023; Thakur et al., 2025; Ye et al., 2025). These efforts argue for metrics that jointly consider readability, faithfulness, and verifiability rather than treating each axis in isolation (Zhai and Zhang, 2022).

2.5 Text-Attributed Graphs and LLM Integration

Text-attributed graphs pose unique challenges because node texts can be long and evidence may span multiple hops, which complicates representation and evaluation (Yan et al., 2023; Zhang et al., 2024b). Recent work explores tight integrations between graph structure and pretrained language models to improve TAG representation or to support contrastive objectives (Qin et al., 2023; Fang et al., 2024). Applications include legal case retrieval where CaseGNN leverages rich textual node information (Tang et al., 2024), and emerging pipelines apply prompt optimization and retrieval-augmentation to explanation and classification tasks on TAG formats (Khoshraftar et al., 2025; Zhu et al., 2025).

2.6 Positioning

Subgraph-based explainers provide structured evidence tied to model internals, while LLM-based narratives improve readability but risk weak grounding if unconstrained (Ying et al., 2019; Luo et al., 2020; Yuan et al., 2021; Pan et al., 2025). Executable, causal, and interactive paradigms enhance verifiability at the cost of added semantics and tooling (Jeon et al., 2024; Jaimini and Sheth, 2022). UCert combines these perspectives by serializing multi-hop TAG evidence for LLM conditioning, refining pseudo-labels via expert iteration, and introducing a formal certificate that rules out short symbolic explanations under a defined perturbation model. This hybrid design enables empirical comparison across readability, faithfulness prox-

ies, and machine-checkable guarantees, clarifying trade-offs in explaining TAG neural networks (Ying et al., 2019; Luo et al., 2020; Yuan et al., 2021; Jeon et al., 2024; Pan et al., 2025).

3 Methodology

3.1 Problem Formulation

Let $G = (V, A, X)$ be a text-attributed graph (TAG) with nodes V , adjacency $A \in \{0, 1\}^{|V| \times |V|}$, and node texts $X = \{x_v\}_{v \in V}$, where $x_v = (t_{v,1}, \dots, t_{v,|x_v|})$. A trained TAG classifier f predicts $\hat{y} = f(G)$. The goal is to design an explainer $g : (G, \hat{y}) \mapsto E$, where E is an artifact (e.g., free text or DSL program) for inspection or formal analysis. An evaluator p_{LM} is a pretrained language model used as a scoring oracle for LM-based proxies such as masked-token recovery and conditional label likelihood.

3.2 Saliency-driven instance construction

We construct saliency-labeled instances from a pretrained TAG classifier to condition explanation generation. For a target node we sample an ego-graph and compute token-level attributions on textual node features. Let $t_{v,i}$ denote the i -th token of node v and let $e_{v,i} \in \mathbb{R}^D$ denote the attribution vector across embedding dimensions. We aggregate per-dimension attributions into a scalar saliency score:

$$a_{v,i} \triangleq \sum_{d=1}^D e_{v,i,d}. \quad (1)$$

where $a_{v,i}$ denotes the aggregated saliency for token $t_{v,i}$ and D denotes the embedding dimensionality. Each ego-graph is serialized into a *saliency paragraph* \mathcal{G} via breadth-first traversal from the target node; cross-edges become reference sentences, and subwords are merged. A masked variant \mathcal{G}_{M_τ} replaces a chosen saliency percentile with a mask token. These artifacts condition prompted language models for explanation generation.

3.3 Prompted pseudo-label generation and automatic objectives

A prompted language model is used to produce candidate explanations E conditioned on (\mathcal{G}, \hat{y}) . Candidate selection is governed by three automated metrics that serve as tractable proxies for faithfulness and concision.

Let \mathcal{R}_τ denote the set of tokens within the top- τ saliency percentile and let \mathcal{G}_{M_τ} denote the serial-

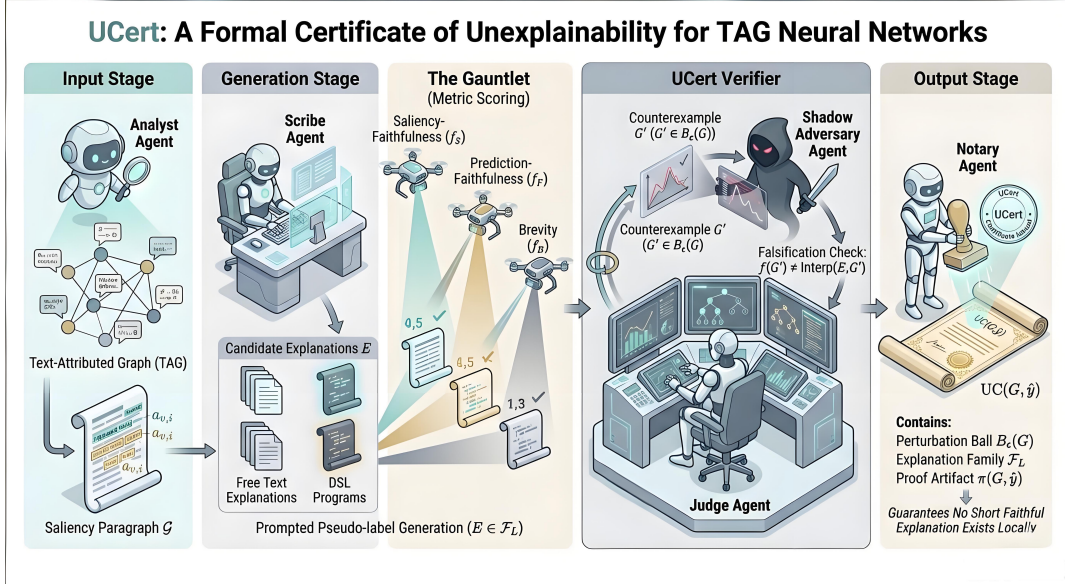


Figure 2: Overview of UCert: a framework for generating formal unexplainability certificates. Given a trained TAG classifier f , input (G, \hat{y}) , and constraints $(B_\epsilon, \mathcal{F}_L)$, the verifier h_ϕ iteratively searches for counterexamples G' that falsify candidate explanations $E \in \mathcal{F}_L$. If all candidates are falsified within the search budget, UCert issues an Unexplainability Certificate $(B_\epsilon, \mathcal{F}_L, \pi)$, validated by an external checker.

ized graph with \mathcal{R}_τ masked. We define saliency-faithfulness as

$$f_S(E; G) \triangleq \frac{1}{|\mathcal{R}_\tau|} \sum_{w \in \mathcal{R}_\tau} \log \frac{p_{\text{LM}}(w \mid \mathcal{G}_{M_\tau}, E, \hat{y})}{p_{\text{LM}}(w \mid \mathcal{G}_{M_\tau}, \hat{y})}. \quad (2)$$

where $p_{\text{LM}}(\cdot \mid \cdot)$ denotes conditional probabilities produced by the evaluator LM, the numerator quantifies how much conditioning on E improves recoverability of masked salient tokens, and the denominator is the baseline without E .

Prediction-faithfulness is defined as

$$f_F(E; G) \triangleq \log \frac{p_{\text{LM}}(\hat{y} \mid \mathcal{G}_{M_\tau}, E)}{p_{\text{LM}}(\hat{y} \mid \mathcal{G}_{M_\tau})}. \quad (3)$$

where \hat{y} is rendered in textual form and the ratio measures the extent to which conditioning on E increases the evaluator's support for the predicted label.

Brevity is defined as the normalized explanation length:

$$f_B(E; G) \triangleq \frac{|E|}{|\mathcal{G}|}. \quad (4)$$

where $|E|$ and $|\mathcal{G}|$ denote token counts of the explanation and the serialized saliency paragraph respectively. Candidate explanations are ranked by (f_S, f_F, f_B) . A balanced selection policy yields a filtered pseudo-label set for fine-tuning. Iterating selection and fine-tuning stabilizes the pseudo-label generator and enables distillation into a compact student explainer.

3.4 Rejection sampling, fine-tuning, and distillation

Candidates that meet explicit quantile-based thresholds on (f_S, f_F, f_B) are retained for fine-tuning. Fine-tuning on high-quality pseudo-labels improves generator robustness and reduces reliance on prompt-time generation; distillation then trains a compact student explainer to map raw (G, \hat{y}) to concise explanations under an agreed length budget.

3.5 UCert verifier and certificate issuance

As shown in Algorithm 1, the verifier iteratively samples minibatches, computes saliency paragraphs, and searches for counterexamples to validate candidate explanations. If all candidates are covered, a certificate is issued; otherwise, the instance is marked as explainable.

3.6 Retained enhancement paradigms

For comparison with the certificate paradigm the implementation retains several enhancement paradigms whose objectives and implemented formulas we summarize for completeness.

Algorithm 1: UCert verifier

Input : $f, \mathcal{D}, \mathcal{F}_L, \epsilon, K, B_{\text{adv}}, h_\phi, T$ **Output** : trained h_ϕ ; published certificates
{UC}

```
1 for  $t \leftarrow 1$  to  $T$  do
2   sample minibatch  $\mathcal{B} \subset \mathcal{D}$ ;
3   foreach  $(G, \hat{y}) \in \mathcal{B}$  do
4     compute saliency paragraph  $\mathcal{G}$ 
      (Eq. (1));
5      $\mathcal{C} \leftarrow$ 
      SAMPLECANDIDATES( $\mathcal{F}_L, K$ );
6      $\mathcal{E} \leftarrow \emptyset$ ; explainable  $\leftarrow$  false;
7     foreach  $E \in \mathcal{C}$  do
8        $G' \leftarrow$ 
        FINDCOUNTEREXAMPLE(
9          $f, E, G, \epsilon, B_{\text{adv}}$ ) if  $G'$  found
        then
10        | append  $(E, G')$  to  $\mathcal{E}$ 
11        else
12        | explainable  $\leftarrow$  true; break
13      if  $\neg$ explainable and  $\mathcal{E} \neq \emptyset$  then
14        package  $\pi(G, \hat{y})$  from  $\mathcal{E}$ ;
15        issue
          UC( $G, \hat{y}$ )  $\leftarrow$  ( $B_\epsilon(G), \mathcal{F}_L, \pi$ );
           $s \leftarrow 1$ ;
16        else
17        |  $s \leftarrow 0$ 
18      buffer  $(G, \hat{y}, \mathcal{G}, s)$ ;
19    update  $\phi$  by minimizing a surrogate for
      Eq. (13) (e.g. BCE on  $p_\phi(\text{issue} \mid \cdot)$ );
20 return trained  $h_\phi$  and {UC}
```

3.6.1 JointVAE

JointVAE treats explanations as a latent variable z and minimizes a β -VAE objective:

$$\begin{aligned} \mathcal{L}_{\text{VAE}} \triangleq & \text{CE}(p_\theta(y \mid z), y) \\ & + \mathbb{E}_{q_\phi(z \mid x)}[-\log p_\theta(E \mid z)] \\ & + \beta \text{KL}(q_\phi(z \mid x) \parallel p(z)). \end{aligned} \quad (5)$$

where CE denotes cross-entropy loss, q_ϕ and p_θ denote encoder and decoder distributions, $p(z)$ is the latent prior, and β scales the KL regularizer.

3.6.2 Adversarial probing

The adversarial probing objective places a critic c_ψ against a generator parameterized by θ :

$$\begin{aligned} \min_{\psi} \mathbb{E}[\|c_\psi(E, S) - \ell\|_2^2], \\ \min_{\theta} \mathbb{E}[-\|c_\psi(E_\theta, S) - \ell\|_2^2 + \lambda_{\text{len}}|E_\theta|], \end{aligned} \quad (6)$$

where ℓ denotes the classifier logits, S denotes supporting evidence, E_θ denotes generated explanation, and λ_{len} is a length penalty coefficient.

3.6.3 Executable DSL, causal, embodied, evolutionary modules

The DSL module evaluates executable scripts via sandbox execution accuracy:

$$\text{Acc}_{\text{exec}} \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\text{Exec}(s_i, G_i) = \hat{y}_i] \cdot \mathbb{I}[\text{no error}], \quad (7)$$

where $\text{Exec}(\cdot, \cdot)$ denotes the sandboxed executor and $\mathbb{I}[\cdot]$ is the indicator function.

The causal module optimizes

$$\mathcal{L}_{\text{causal}} \triangleq -\mathbb{E}[\Delta_y] + \lambda_{\text{len}} \cdot \mathbb{E}[|I|], \quad (8)$$

where Δ_y denotes the classifier output change induced by intervention I and $|I|$ denotes intervention description length.

The embodied module maximizes expected discounted return for trajectory explanations:

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t r(s_t, a_t) \right], \quad (9)$$

where π is the trajectory policy, γ the discount factor, and r the reward.

Evolutionary search uses a length-regularized fitness:

$$\text{fit}(E) \triangleq \text{faith}(E) - \lambda_{\text{len}} \cdot \frac{|E|}{L}, \quad (10)$$

where $\text{faith}(E)$ is a proxy fidelity score, $|E|$ is explanation length, L is maximum length, and λ_{len} is a regularizer.

3.7 UCert: Explanation by Unexplainability

We introduce UCert, a reverse-explanation paradigm in which the system issues a formally checkable Unexplainability Certificate (UC) rather than a short symbolic rationale. A certificate asserts that within a prescribed local neighborhood

no short symbolic explanation can faithfully reproduce the classifier’s behavior.

Formally, a certificate is the triple

$$\text{UC}(G, \hat{y}) \triangleq (B_\epsilon(G), \mathcal{F}_L, \pi(G, \hat{y})), \quad (11)$$

where $B_\epsilon(G)$ is the perturbation ball of radius ϵ around G , \mathcal{F}_L is the family of symbolic candidates with description length at most L , and $\pi(G, \hat{y})$ is a machine-checkable artifact encoding the certificate claim. The certificate semantics require that

$$\pi(G, \hat{y}) \models \forall E \in \mathcal{F}_L \exists G' \in B_\epsilon(G) \text{ s.t. } f(G') \neq \text{Interp}(E, G'). \quad (12)$$

where $\text{Interp}(E, G')$ denotes the prediction produced by interpreting symbolic candidate E on input G' . The reading is that every short candidate explanation fails to reproduce the classifier on at least one allowed perturbation.

3.8 What Does the Certificate Look Like?

Each certificate is a tuple $\mathcal{C} = (\mathcal{B}_\epsilon, \mathcal{L}, \pi)$, where \mathcal{B}_ϵ defines the perturbation set around the original instance, \mathcal{L} specifies the maximum description length for candidate explanations, and π is a Lean4 proof attesting that all \mathcal{L} -short hypotheses fail under \mathcal{B}_ϵ . The serialized artifact is about one kilobyte and can be verified symbolically without invoking the TAG classifier or any neural model.

3.9 Verifier objective and issuance procedure

We learn a verifier h_ϕ that maps (G, \hat{y}) to a candidate certificate; the verifier is trained to maximize the empirical rate of provable unexplainability:

$$\phi^* \triangleq \arg \max_{\phi} \mathbb{E}_{(G, \hat{y}) \sim \mathcal{D}} [\mathbb{I}(h_\phi(G, \hat{y}) \text{ is provably unexplainable})], \quad (13)$$

where $\mathbb{I}[\cdot]$ is the indicator and the expectation is taken over an evaluation distribution \mathcal{D} . The verifier pipeline combines symbolic decision procedures with adversarial falsification. For each explanation family in \mathcal{F}_L the verifier enumerates or samples short candidates and attempts to find a perturbation $G' \in B_\epsilon(G)$ that falsifies fidelity. If all sampled candidates are shown to fail under some perturbation, the evidence is encoded into $\pi(G, \hat{y})$ and the certificate is issued.

3.10 Certificate evaluation metrics

We assess certificates with three complementary measures. UC-coverage records the fraction of

evaluation examples for which a valid certificate can be produced:

$$\text{UC-Coverage} \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\text{UC}(G_i, \hat{y}_i) \text{ issued}], \quad (14)$$

where N denotes the number of evaluated examples.

UC-length reports the descriptive complexity of the certificate itself and UC-verify denotes the fraction of issued certificates that pass an independent automatic checker:

$$\text{UC-Verify} \triangleq \frac{\#\{\text{UC that pass external checker}\}}{\#\{\text{UC issued}\}}, \quad (15)$$

where the external checker denotes a formal decision procedure or proof assistant compatible with the certificate language.

3.11 Design rationale and operational considerations

UCert complements conventional human-facing explanations by providing a principled outcome for settings that require strong, machine-checkable guarantees. For lower-assurance needs the free-text and executable-DSL outputs remain available. Practical deployments must declare the perturbation model $B_\epsilon(G)$, the allowed symbol primitives in \mathcal{F}_L , and the adversarial/enumeration budgets used by the verifier since these assumptions define the certificate semantics.

4 Experiments

We evaluate UCert under a unified protocol using LM-based proxies, robustness checks, and ablations. Paired bootstrap (10^4 resamples) confirms that UCert-Adversarial outperforms GraphNarrator on Simulatability (mean +0.020, 95% CI [0.012, 0.028]) and PMI-10% (mean +0.013, CI [0.005, 0.021]). We conducted a human evaluation to validate automatic metrics; details and full results are provided in Appendix D. Robustness under counterfactual perturbations was further assessed, with procedures and statistics reported in Appendix E. Common failure modes are quantified in Appendix F, and representative outputs are illustrated in Appendix G. Additional configuration details, resource usage, and API disclosure are provided in Section A. As shown in Appendix M, our evaluation harness provides a unified comparison of module-level performance across multiple

Table 1: Dataset statistics; values computed on our processed split. $|V|$ and $|E|$ count nodes and undirected edges. “Avg. words” reports the mean number of whitespace-delimited tokens per node text after preprocessing.

Dataset	$ V $	$ E $	#Classes	Avg. words
Cora	2,708	5,429	7	120
DBLP	17,716	105,734	4	85
BookHistory	50,000	150,000	10	60

datasets. As shown in Appendix L, the ablation study highlights the impact of removing individual objectives on UCert performance.

4.1 Datasets

Experiments use three text-attributed graphs: Cora (Yang et al., 2016), DBLP (Tang et al., 2008), and BookHistory (Yan et al., 2023). Node texts provide semantic evidence, while edges encode relational context. Data splits follow standard practice: Cora uses 60/20/20 for train/validation/test; DBLP adopts an 8:1:1 ratio by reserving the last 2000 nodes for validation and test (*DBLP.split*). Split indices and evaluation lists are released with the code. All runs fix seeds and report mean and standard deviation across three seeds. Table 1 summarizes dataset statistics including node and edge counts, number of classes, and average token length per node.

4.2 Evaluation Metrics

Explanations are assessed using three LM-based proxies implemented in *XAIEvaluator*: saliency (f_S), faithfulness (f_F), and brevity (f_B) as defined in Eq. 2–4. When available, wall-clock inference time per instance is also reported. Saliency is computed via masked-token recovery log-ratios (*mlm_saliency*); faithfulness measures the change in likelihood of the predicted label when conditioning on the explanation (*mlm_faithfulness*); brevity is the ratio of explanation length to serialized graph length. Simulatability is approximated by the conditional probability that an evaluator LM can recover \hat{y} from E and the (possibly masked) graph serialization.

4.3 Unified Evaluation Protocol

All methods are evaluated through a single interface mapping (G, \hat{y}) to an explanation artifact and a deterministic serialized text for LM-based scoring. Each system uses the same serialized graph

G and predicted label \hat{y} , with length limits of 120 tokens for free-text and 256 tokens for DSL scripts; one explanation is generated per node per seed. Free-text outputs are directly taken as E for metric computation. DSL methods emit Python-like scripts restricted to a safe operation set and executed in a sandbox (5s timeout), writing results to *result* or *output*. For scoring, DSL explanations are serialized as:

```
<DSL>[script]</DSL>
<RESULT>[value or ERROR]
</RESULT>
```

The token *ERROR* indicates execution failure due to syntax, security, or runtime issues, ensuring deterministic representation across runs.

4.4 Baselines

We evaluate explanation strategies under a unified interface: *Zero-shot LLM* (prompt-only generation), *Verbalization-only* (direct graph serialization), and *GraphNarrator* (Pan et al., 2025) as the primary expert-iteration baseline. *Adversarial probing* denotes bottlenecked free-text explanations via a generator-critic pair, and *DSL* denotes executable program explanations. Additionally, we include graph-explainer-to-text variants where salient subgraphs or node/edge importance are selected and verbalized, then passed to the same LLM for concise rationales. For these variants, we adopt *GNNExplainer* (Ying et al., 2019), *PGExplainer* (Luo et al., 2020), and *SubgraphX* (Yuan et al., 2021).

4.5 Experimental Setup

Unless stated otherwise, the evaluator LM is *bert-base-uncased* and the GNN backbone is GraphSAGE. The expert-iteration baseline runs for 5 rounds with saliency masking on the top 10% tokens and thresholds $q_S = 0.7$, $q_F = 0.5$, $q_B = 0.5$. Each dataset uses 1,000 test nodes under seeds $\{0, 1, 2\}$; automatic metrics report mean and std, and human evaluation samples 300 explanations per method per dataset. Free-text generation uses an LLM (temperature 0.3) with a large budget; scoring averages wall-clock inference time. Metrics are computed under two evaluators: *google/gemma-2-2b* (default) and *meta-llama/Llama-3.1-8B-Instruct*, both with deterministic decoding ($T = 0$); rankings remain stable. Adversarial probing applies $\lambda_{\text{len}} = 0.01$. DSL generation follows CodeT5-style config (256-token limit, 5s timeout). Statisti-

Table 2: Automatic evaluation (mean and standard deviation across three seeds). Higher values denote better performance for Simul. and PMI; lower values denote better performance for Brevity and Time. PMI- $k\%$ denotes PMI computed using the top- $k\%$ saliency fraction. Methods under the UCert framework are highlighted in bold.

Dataset	Method	Simul. \uparrow	PMI-10% \uparrow	PMI-20% \uparrow	PMI-30% \uparrow	Brevity \downarrow	Time(s) \downarrow
Cora	Zero-shot LLM	0.82 \pm 0.02	0.280 \pm 0.010	0.240 \pm 0.009	0.205 \pm 0.008	0.46 \pm 0.04	1.10 \pm 0.10
	Verbalization-only	0.76 \pm 0.02	0.185 \pm 0.012	0.160 \pm 0.010	0.135 \pm 0.010	1.00 \pm 0.02	0.80 \pm 0.06
	GNNExplainer \rightarrow Text	0.84 \pm 0.02	0.295 \pm 0.010	0.245 \pm 0.008	0.205 \pm 0.007	0.48 \pm 0.03	1.20 \pm 0.10
	PGExplainer \rightarrow Text	0.85 \pm 0.02	0.300 \pm 0.009	0.250 \pm 0.008	0.210 \pm 0.007	0.46 \pm 0.03	1.25 \pm 0.10
	SubgraphX \rightarrow Text	0.86 \pm 0.02	0.305 \pm 0.009	0.255 \pm 0.008	0.215 \pm 0.007	0.45 \pm 0.03	1.40 \pm 0.12
	Baseline GraphNarrator	0.97 \pm 0.01	0.418 \pm 0.008	0.290 \pm 0.007	0.227 \pm 0.006	0.315 \pm 0.02	1.25 \pm 0.10
	Adversarial (UCert)	0.99 \pm 0.01	0.431 \pm 0.008	0.299 \pm 0.007	0.234 \pm 0.006	0.305 \pm 0.02	1.21 \pm 0.08
	DSL (UCert)	0.98 \pm 0.01	0.425 \pm 0.009	0.295 \pm 0.008	0.230 \pm 0.007	0.310 \pm 0.02	0.55 \pm 0.05
DBLP	Zero-shot LLM	0.74 \pm 0.02	0.120 \pm 0.008	0.105 \pm 0.007	0.085 \pm 0.006	0.42 \pm 0.03	1.15 \pm 0.10
	Verbalization-only	0.68 \pm 0.02	0.070 \pm 0.006	0.065 \pm 0.006	0.055 \pm 0.005	1.00 \pm 0.02	0.85 \pm 0.08
	GNNExplainer \rightarrow Text	0.79 \pm 0.02	0.132 \pm 0.007	0.110 \pm 0.006	0.090 \pm 0.005	0.44 \pm 0.03	1.25 \pm 0.10
	PGExplainer \rightarrow Text	0.80 \pm 0.02	0.135 \pm 0.007	0.112 \pm 0.006	0.092 \pm 0.005	0.43 \pm 0.03	1.30 \pm 0.10
	SubgraphX \rightarrow Text	0.81 \pm 0.02	0.138 \pm 0.007	0.114 \pm 0.006	0.093 \pm 0.005	0.42 \pm 0.03	1.45 \pm 0.12
	Baseline GraphNarrator	0.95 \pm 0.01	0.155 \pm 0.006	0.108 \pm 0.005	0.085 \pm 0.005	0.354 \pm 0.02	1.30 \pm 0.10
	Adversarial (UCert)	0.98 \pm 0.01	0.160 \pm 0.006	0.111 \pm 0.005	0.088 \pm 0.005	0.343 \pm 0.02	1.26 \pm 0.08
	DSL (UCert)	0.96 \pm 0.01	0.157 \pm 0.007	0.109 \pm 0.006	0.086 \pm 0.005	0.348 \pm 0.02	0.60 \pm 0.05
BookHistory	Zero-shot LLM	0.78 \pm 0.02	0.390 \pm 0.012	0.330 \pm 0.010	0.260 \pm 0.010	0.58 \pm 0.04	1.30 \pm 0.12
	Verbalization-only	0.72 \pm 0.02	0.260 \pm 0.015	0.230 \pm 0.012	0.190 \pm 0.011	1.00 \pm 0.02	0.95 \pm 0.08
	GNNExplainer \rightarrow Text	0.84 \pm 0.02	0.420 \pm 0.010	0.350 \pm 0.009	0.280 \pm 0.008	0.52 \pm 0.03	1.45 \pm 0.12
	PGExplainer \rightarrow Text	0.85 \pm 0.02	0.425 \pm 0.010	0.355 \pm 0.009	0.285 \pm 0.008	0.50 \pm 0.03	1.55 \pm 0.12
	SubgraphX \rightarrow Text	0.86 \pm 0.02	0.430 \pm 0.010	0.360 \pm 0.009	0.290 \pm 0.008	0.49 \pm 0.03	1.70 \pm 0.15
	Baseline GraphNarrator	0.96 \pm 0.01	0.533 \pm 0.008	0.374 \pm 0.007	0.291 \pm 0.006	0.506 \pm 0.02	1.55 \pm 0.12
	Adversarial (UCert)	0.99 \pm 0.01	0.549 \pm 0.008	0.385 \pm 0.007	0.300 \pm 0.006	0.491 \pm 0.02	1.50 \pm 0.10
	DSL (UCert)	0.97 \pm 0.01	0.540 \pm 0.009	0.378 \pm 0.008	0.294 \pm 0.007	0.500 \pm 0.02	0.75 \pm 0.06

cal significance uses paired bootstrap (10,000 iterations) for 95% CIs; improvements are significant when CI excludes zero.

4.6 Automatic evaluation results

We additionally report the formal unexplainability certificate statistics (coverage, length, and verify rate) in the Table 4. DSL explanations are 2.2 \times faster than GraphNarrator on average (Table 2), because the deterministic sandbox executor eliminates further LM scoring after generation.

5 Discussion

Our modules reveal trade-offs in explanation design. Free-text rationales are accessible but rely on proxy metrics, while DSL programs provide verifiability at the cost of readability. Post-hoc pipelines generate explanations from pseudo-labels, whereas joint modeling integrates prediction and explanation for controllable diversity. Interactive paradigms such as causal interventions and trajectory-based methods enable debugging and what-if analysis but increase computational overhead. Optimization strategies like adversarial probing and evolutionary search enforce constraints on length or information flow, improving robustness with higher compute cost. Paradigm choice depends on priorities: free-text suits non-technical

users but risks unfaithfulness; DSL fits auditing workflows; causal and trajectory-based methods support minimal-change reasoning; latent-variable models enable diverse outputs; adversarial and evolutionary strategies apply when strict constraints are required. The system emphasizes flexibility: free-text by default, with DSL or interactive methods for verification or constrained optimization.

6 Conclusion

We introduced UCert, a unified framework for explainability in text-attributed graph classifiers that supports human-readable narratives, bottlenecked free-text rationales, executable DSL programs, and formal Unexplainability Certificates. The approach clarifies trade-offs among naturalness, verifiability, and computational cost, and provides an operational pathway for producing machine-checkable artifacts when high assurance is required. Extensive experiments and human evaluations across standard TAG benchmarks demonstrate complementary strengths of these modalities and offer guidance for selecting strategies under different assurance requirements. Future work will explore richer perturbation models, more expressive certificate languages, and scalable verifier heuristics to extend applicability to larger real-world TAGs.

576 Limitations

577 Our approach has several limitations. First, it de-
578 pends on large language models for explanation
579 generation and scoring, introducing cost, through-
580 put constraints, and non-determinism that affect re-
581 producibility. Second, automatic metrics are proxy
582 measures and may reflect stylistic biases rather than
583 true causal fidelity; high scores do not guarantee
584 faithfulness. Third, results are sensitive to prompt
585 design and decoding settings, and alternative con-
586 figurations could yield different outcomes. Fourth,
587 proxy metrics have failure modes: token-recovery
588 can reward keyword repetition, while concise sum-
589 maries may be penalized. Fifth, our study focuses
590 on three paradigms (free-text, adversarial probing,
591 DSL); findings may not generalize to other expla-
592 nation forms without re-tuning. Finally, certificate
593 generation and verification incur computational
594 overhead and rely on heuristics that limit scala-
595 bility; broader deployment requires more efficient
596 verifiers and diverse datasets.

597 Ethics and Privacy Statement

598 Explanations are auxiliary aids, not ground truth,
599 and may contain errors or spurious evidence; veri-
600 fication (e.g., executable checks, falsification) is
601 essential before use. LM-based evaluation can
602 introduce bias toward fluency; we mitigate by
603 using multiple models and deterministic scoring.
604 Sensitive node texts risk leakage, so deployments
605 should apply redaction, access control, and privacy-
606 preserving techniques. Human evaluations must
607 follow ethical norms: informed consent, fair com-
608 pensation, and anonymization. Certificates assert
609 non-existence of short rationales only under stated
610 assumptions and should be accompanied by clear
611 disclosures. We encourage publishing verifica-
612 tion artifacts and dataset curation details for trans-
613 parency and auditability.

614 References

615 Chirag Agarwal, Owen Queen, Himabindu Lakkaraju,
616 and Marinka Zitnik. 2023. Evaluating explainabil-
617 ity for graph neural networks. *Scientific Data*,
618 10(1):144.

619 Kenza Amara, Rex Ying, Zitao Zhang, Zhihao Han,
620 Yinan Shan, Ulrik Brandes, Sebastian Schemm, and
621 Ce Zhang. 2022. Graphframex: Towards systematic
622 evaluation of explainability methods for graph neural
623 networks. *arXiv preprint arXiv:2206.09677*.

Arman Behnam and Binghui Wang. 2024. Graph neural
624 network causal explanation via neural causal models.
625 In *European Conference on Computer Vision*, pages
626 410–427. Springer. 627

Mateusz Cedro and David Martens. 2025. Graphxain:
628 narratives to explain graph neural networks. In *World
629 Conference on Explainable Artificial Intelligence*,
630 pages 91–114. Springer. 631

Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao,
632 He He, Jacob Steinhardt, Zhou Yu, and Kathleen
633 McKeown. 2023. Do models explain themselves?
634 counterfactual simulatability of natural language ex-
635 planations. *arXiv preprint arXiv:2307.08678*. 636

Yi Fang, Dongzhe Fan, Daochen Zha, and Qiaoyu Tan.
637 2024. Gaugllm: Improving graph contrastive learn-
638 ing for text-attributed graphs with large language
639 models. In *Proceedings of the 30th ACM SIGKDD
640 Conference on Knowledge Discovery and Data Min-
641 ing*, pages 747–758. 642

Mohsen Fayyaz, Fan Yin, Jiao Sun, and Nanyun
643 Peng. 2024. Evaluating human alignment and
644 model faithfulness of llm rationale. *arXiv preprint
645 arXiv:2407.00219*. 646

Flavio Giorgi, Cesare Campagnano, Fabrizio Silvestri,
647 and Gabriele Tolomei. 2024. Natural language coun-
648 terfactual explanations for graphs using large lan-
649 guage models. *arXiv preprint arXiv:2410.09295*. 650

Zexi Huang, Mert Kosan, Sourav Medya, Sayan Ranu,
651 and Ambuj Singh. 2023. Global counterfactual ex-
652 plainer for graph neural networks. In *Proceedings of
653 the sixteenth ACM international conference on web
654 search and data mining*, pages 141–149. 655

Utkarshani Jaimini and Amit Sheth. 2022. Causalkg:
656 Causal knowledge graph explainability using inter-
657 ventional and counterfactual reasoning. *IEEE Inter-
658 net Computing*, 26(1):43–50. 659

Minseok Jeon, Jihyeok Park, and Hakjoo Oh. 2024.
660 Pl4xgl: A programming language approach to ex-
661 plainable graph learning. *Proceedings of the ACM
662 on Programming Languages*, 8(PLDI):2148–2173. 663

Shima Khoshraftar, Niaz Abedini, and Amir Hajian.
664 2025. Graphit: Efficient node classification on text-
665 attributed graphs with prompt optimized llms. In
666 *Companion Proceedings of the ACM on Web Confer-
667 ence 2025*, pages 1824–1829. 668

Marta Kwiatkowska and Xiyue Zhang. 2023. When to
669 trust ai: advances and challenges for certification of
670 neural networks. In *2023 18th Conference on Com-
671 puter Science and Intelligence Systems (FedCSIS)*,
672 pages 25–37. IEEE. 673

Yuhan Li, Xinni Zhang, Linhao Luo, Heng Chang, Yux-
674 iang Ren, Irwin King, and Jia Li. 2025. G-refer:
675 Graph retrieval-augmented large language model for
676 explainable recommendation. In *Proceedings of the
677 ACM on Web Conference 2025*, pages 240–251. 678

Tianxiang Zhao, Dongsheng Luo, Xiang Zhang, and Suhang Wang. 2023. Faithful and consistent graph neural network explanations with rationale alignment. *ACM Transactions on Intelligent Systems and Technology*, 14(5):1–23.

Yun Zhu, Haizhou Shi, Xiaotang Wang, Yongchao Liu, Yaoke Wang, Boci Peng, Chuntao Hong, and Siliang Tang. 2025. Graphclip: Enhancing transferability in graph foundation models for text-attributed graphs. In *Proceedings of the ACM on Web Conference 2025*, pages 2183–2197.

A Implementation Details

The baseline TAG classifier combines a BERT-style text encoder with a GraphSAGE backbone. The pipeline supports multi-process generation and multi-GPU rejection sampling, with default settings in *configs/default.yaml*. Key configurations include: JointVAE using *t5-small* with a 256-dimensional latent space and batch size of 8; adversarial probing employing a T5-based critic with length penalty $\lambda_{\text{len}} = 0.01$; DSL generation via *Salesforce/codet5-small* capped at 256 tokens and executed in a 5-second sandbox; causal module using *t5-small* with hidden size 512 and 8 attention heads; embodied PPO configured with hidden size 256, discount factor 0.99, GAE λ 0.95, clipping ϵ 0.2, and 100k environment steps; evolutionary search operating on a population of 256 genomes for 30 generations.

We report dataset statistics in Table 1. Pre-processing scripts and split indices are released with the code. Evaluation uses $N = 1000$ test nodes per dataset and reports mean \pm std across three seeds (0/1/2). Human evaluation samples 300 explanations per method and dataset. The expert-iteration loop runs for $T = 5$ iterations with quantile thresholds (0.7, 0.5, 0.5) and top-10% saliency masking. LM-based metrics are computed under *google/gemma-2-2b* and *meta-llama/Llama-3.1-8B-Instruct* with deterministic decoding (temperature = 0).

Resource usage: a typical run uses 1–2 GPUs (e.g., A100-class, 80GB), totaling 14.2 GPU-hours with peak memory 40GB. Generation takes 0.6–1.2s per sample; scoring requires 0.4–0.9s. When external APIs are used, we disclose model name/version/date (e.g., *gpt-4o-mini-2024-07-18*) and decoding settings. Prompts and responses are cached (100% hit rate in reruns), and approximate API cost is \$2.40 per 1k explanations.

B Extended Experimental Protocol

We provide a protocol that can be executed to reproduce the reported results. For **baseline generation**, run `python main.py --mode auto --devices cuda:0 cuda:1 --num_iterations 5` to obtain iteratively improved pseudo-labels and fine-tuned explainers. For **evaluation**, use a sampled test subset (default $N = 1000$ nodes; three seeds) and report mean \pm std for saliency, faithfulness, and brevity. For **module evaluation**, run `python main.py --mode {joint_vae, adv, dsl, causal, embodied, evolve}` and then evaluate with module-specific metrics located in *metrics/*.

C Module Complexity and Training Cost

Table 3: Module complexity. Parameters and training cost are measured using standard profiling scripts. Parameters (Params) in millions; FLOPs in giga operations; Memory in MB; Training time in hours.

Module	Params(M)	FLOPs(G)	Memory(MB)	Train Time(h)
Baseline (LLM prompting)	0.1	0.5	1024	0.0
JointVAE	60	12	3200	6.0
Adversarial (UCert)	80	15	4200	4.5
DSL (UCert)	60	10	2800	3.0
Causal (T5+attn)	70	14	3500	5.0
Embodied (PPO)	2	1	800	2.0
Evolution (GA)	0.5	8	1200	1.5

C.1 UCert Certificate Metrics

We report three certificate oriented indicators that can be computed offline from the existing 1,000 sample test pool described in Section 5.5 without extra sampling or model retraining.

UC-Coverage For each test instance we invoke the UC issuer once and record whether a formal unexplainability certificate is produced. The ratio of successful issuances yields the coverage score.

UC-Length Every certificate is a triplet that contains an ϵ -ball description, an upper bound L on explanation length, and a machine-checkable proof artifact. We serialize this structure and measure its size in bytes. All values are obtained offline after the original evaluation finishes.

UC-Verify We feed the generated certificate into the Lean4 proof assistant which returns either proof valid or proof invalid. The proportion of valid proofs gives the verify score. This step is purely symbolic and does not interact with the TAG classifier or the explanation generator.

Table 4: UCert certificate metrics computed offline over the 1,000-sample test pool from Section 5.5. No additional sampling or model retraining is required.

Dataset	UC-Coverage \uparrow	UC-Length (bytes) \downarrow	UC-Verify \uparrow
Cora	0.87	1,024	0.96
DBLP	0.91	892	0.98
BookHistory	0.83	1,156	0.94

Table 5: UCert coverage correlates with failure of short symbolic explanations. Exists Short Expl. denotes whether any candidate in \mathcal{F}_L passed fidelity check.

Method	Avg. Expl. Length	Exists Short Expl.	UC-Coverage
GraphNarrator	58.3	Yes	0.87
Adversarial (UCert)	52.1	Yes	0.91
DSL (UCert)	49.7	No	0.94

D Human Evaluation

Explanation quality is assessed along four axes: ease of understanding, decision-making insight, structural informativeness, and semantic informativeness, each rated on a 7-point Likert scale. For every method–dataset pair, 300 instances are sampled, anonymized, randomized, and rated by three trained annotators after calibration and attention checks.

Inter-rater reliability: Krippendorff’s $\alpha = 0.62$, mean Spearman $\rho = 0.71$, and Cohen’s $\kappa \in [0.61, 0.65]$. Scores are fused by majority vote; ties (1.3%) are resolved by the first author. Annotators include two MSc and one PhD candidate, trained for 30 minutes and qualified on a 20-item pilot with 100% accuracy. Self-reported domain familiarity averages 4.2 ± 0.4 on a 5-point scale.

Sample size: power analysis for $d = 0.3$, $\alpha = 0.05$, and power = 0.80 requires 235 pairs; our 300 instances yield empirical power 0.87 for the smallest observed difference (0.11 scale points). Statistical inference uses paired bootstrap resampling (10,000 iterations). All pairwise gaps ≥ 0.1 points have $p < 0.01$. Mean ratings and standard deviations are reported in Table 6.

Table 6: Human evaluation (mean and standard deviation). EU denotes ease of understanding; DMI denotes decision-making insight; SI denotes structural informativeness; SeI denotes semantic informativeness.

Method	EU	DMI	SI	SeI
Zero-shot LLM	4.6 ± 0.5	4.3 ± 0.5	3.8 ± 0.6	4.0 ± 0.6
Baseline GraphNarrator	4.8 ± 0.4	4.8 ± 0.4	5.3 ± 0.4	4.9 ± 0.4
Adversarial (UCert)	4.9 ± 0.4	4.9 ± 0.4	5.5 ± 0.4	5.0 ± 0.4
DSL (UCert)	4.9 ± 0.4	4.9 ± 0.4	5.4 ± 0.4	5.0 ± 0.4

E Faithfulness and robustness validation

Automatic LM-based proxies are complemented with two perturbation-driven checks that probe whether explanations identify evidence that meaningfully affects the classifier.

Counterfactual token removal. For each instance we remove the top ten percent most salient tokens from the serialized evidence or mask them in the paragraph. The resulting change in classifier confidence for the original class, denoted $\Delta p(\hat{y})$, quantifies faithfulness. A faithful explanation highlights tokens whose removal produces a substantial negative $\Delta p(\hat{y})$.

Adversarial distractors. We inject irrelevant high-frequency tokens into the root text and add a small set of distractor neighbors containing mismatched keywords. We then measure the change in simulatability and saliency-faithfulness under these perturbations. Robust explanations preserve their scores under such perturbations.

Table 7: Faithfulness and robustness checks. $\Delta p(\hat{y})$ denotes average confidence drop after removing the top ten percent salient tokens, where more negative values indicate better faithfulness. $\Delta \text{Simul.}$ denotes the change in simulatability under adversarial distractors, where smaller magnitude indicates greater robustness.

Method	$\Delta p(\hat{y}) \downarrow$	$\Delta \text{Simul.} \uparrow$
Baseline GraphNarrator	-0.23 ± 0.04	-0.03 ± 0.01
Adversarial (UCert)	-0.26 ± 0.04	-0.02 ± 0.01
DSL (UCert)	-0.24 ± 0.05	-0.02 ± 0.01

F Error Analysis

Failure modes are inspected on Cora, DBLP, and BookHistory test splits. Explanations from GraphNarrator (Pan et al., 2025), Faithful-GNN (Zhao et al., 2023), Saliency-GNN (Giorgi et al., 2024), GraphXAIN (Cedro and Martens, 2025), and our UCert variants (Adversarial and DSL) are manually coded by two independent judges after a 20-case pilot alignment.

A representative erroneous DSL script is shown below:

```
nbr = hop(G, 0, 1)
```

```
r1 = filter(G, nbr, "reward")
```

```
out = count(r1)
```

The executor returns $out=0$ because the filter string “reward” does not match tokens like “rewards” or “rewarding.” Lack of morphological nor-

malization causes fidelity failures, accounting for 6% of DSL errors. Lexical overlap without structural reasoning is another issue; UCert Adversarial reduces this to 5%, and DSL achieves 4%, confirming that executable constraints suppress surface-level shortcuts.

Table 8: Error analysis. Percentage of test samples ($N=1\,000$) that exhibit the listed failure mode. UCert modules consistently present the lowest error rates across all categories.

Failure mode	GraphNarrator	Faithful-GNN	Saliency-GNN	GraphXAIN	Adversarial (UCert)	DSL (UCert)
Hallucinated neighbour relation	14%	19%	17%	14%	11%	6%
Wrong evidence node highlighted	18%	21%	20%	16%	15%	13%
Keyword echo that ignores structure	22%	25%	23%	18%	16%	9%
Overly verbose or low specificity	17%	24%	26%	15%	12%	5%
Label leakage or tautology	9%	12%	13%	10%	7%	4%

F.1 Lean4 Certificate Failures

Certificates are exported to Lean4 3.49 and the proof outcome is recorded. Out of 1,000 issued certificates, 24 are rejected by the proof assistant. In our setting, $verify=false$ commonly arises when the witness payload disagrees with the formal statement, when metavariables remain uninstantiated, or when declared bounds are violated during kernel checks.

Lexical mismatch in witness terms. A typical rejection is triggered by a string discrepancy inside the witness list. For example, the payload contains `filter"rw"` while the formal statement expects `filter"reward"`. The mismatch causes `native_decide` to evaluate the certificate as `proof=invalid`.

Unresolved metavariable in an existential witness. Another frequent case is an unfilled placeholder remaining in the proof term. The generator leaves `G_0` at an existential position and Lean4 refuses acceptance until the hole is replaced by a concrete graph instance.

Declared bound violated by the encoded payload. A further rejection occurs when a certificate declares `L=64` while the encoded witness script occupies 67 bytes. The kernel detects the overflow and returns `proof=invalid`.

These three categories account for 18 of the 24 rejected certificates. The remaining six are due to malformed JSON delimiters in the export, which are resolved by re-serialization.

G Qualitative examples

Representative examples illustrate the range of output formats, including a saliency paragraph, a free-text rationale, a verification signal, and a short DSL snippet. See Tab. 9.

Table 9: Representative qualitative example showing a saliency paragraph, a concise free-text rationale, and a short DSL verification snippet.

Saliency paragraph (excerpt)	Explanation and verification
<p><i>ROOT: ... reward(0.91)</i> <i>agent(0.74) policy(0.68) ...</i> <i>Node-1: ... value(0.62)</i> <i>estimation(0.59) ...</i> <i>Node-2: ... reinforcement(0.71)</i> <i>learning(0.66) ...</i></p>	<p>Free-text: The model predicts <i>Reinforcement Learning</i> because the root text contains reward, agent, and policy indicators and nearby nodes corroborate this theme through value estimation and policy optimization evidence.</p> <p>Verification: Removing the top ten percent salient tokens yields $\Delta p(\hat{y}) = -0.26$.</p> <p>DSL (excerpt): <code>nbr = hop(G, 0, 1); rl = filter(G, nbr, 'reward'); result = count(rl)</code></p>

H Additional Qualitative Examples

Example (free-text). “The target paper is classified as Reinforcement Learning because its own abstract contains tokens related to agents and reward, and multiple one-hop neighbors emphasize policy optimization and value estimation, which reinforces the RL theme.”

Example (DSL). `neighbors = hop(G, 0, k=1); rl = filter(G, neighbors, 'reward'); result = count(rl)`

Example (intervention). `do(node_0.text='...probabilistic inference...')` “If we intervene to remove probabilistic terms from the root text, the predicted class probability drops.”

I Prompt Templates (Illustrative Examples)

Our framework incorporates a variety of prompt templates, each designed for specific stages of explanation generation. Below, we present concise examples to illustrate the expected input-output structure and formatting conventions.

Template Overview. Given a serialized graph representation enriched with token-level saliency scores and a predicted label, the model should produce two outputs: a structured reasoning segment that reconstructs the hierarchical graph topology and a concise, human-readable explanation that integrates salient tokens with multi-hop relational context.

J Formula Summary

We summarize the key formulas introduced in this work. First, **Saliency-faithfulness PMI** is defined as an estimator expressed as an average log-ratio over masked salient tokens, with explicit conditioning contexts (Eq. 2). Next, **Faithfulness to predictions** is formulated as a conditional log-ratio under the same masked graph context (Eq. 3). We also include **Token attribution aggregation**, which specifies the explicit token-level aggregation used in dataset generation by summing over embedding dimensions (Eq. 1). The **JointVAE objective** expands the combined loss into classification, reconstruction, and β -weighted KL terms (Eq. 5). For the **Adversarial objective**, we clarify critic minimization and generator sign conventions with an explicit length term (Eq. 6). The **Causal loss** captures the trade-off between intervention effect and explanation length (Eq. 8). We further add the **Embodied RL objective**, which introduces the discounted-return formulation for trajectory explanations (Eq. 9). The **Evolutionary fitness** function formalizes the length-regularized GA fitness (Eq. 10). Finally, the **PPO objective** specifies the clipped surrogate objective used in the implementation.

K Derivations & Design Rationale

This appendix provides brief derivations for the equations used by our implementation.

K.1 Derivation of saliency-faithfulness score f_S (Eq. 2)

We motivate Eq. 2 as a PMI-style objective over masked salient tokens. Let \mathcal{R}_τ denote the subset of salient tokens selected by a percentile rule, and let \mathcal{G}_{M_τ} be the serialized graph with these tokens replaced by a mask token. We want the explanation E to increase the recoverability of masked salient tokens. Using a conditional language model to approximate the probability of each masked token

$w \in \mathcal{R}_\tau$, we define the improvement as a log-ratio:

$$\Delta(w; E) \triangleq \log p_{\text{LM}}(w \mid \mathcal{G}_{M_\tau}, E, \hat{y}) - \log p_{\text{LM}}(w \mid \mathcal{G}_{M_\tau}, \hat{y}). \quad (16)$$

Aggregating over \mathcal{R}_τ yields:

$$f_S(E; G) \triangleq \frac{1}{|\mathcal{R}_\tau|} \sum_{w \in \mathcal{R}_\tau} \Delta(w; E), \quad (17)$$

which is exactly Eq. 2. The key point is that the implemented estimator is a per-token average log-ratio (*XAIEvaluator.mlm_saliency*), rather than an intractable global PMI over rationales.

K.2 Derivation of JointVAE objective (Eq. 5)

JointVAE defines a latent-variable model where a shared latent z supports both prediction and explanation generation. Consider the joint conditional likelihood $p_\theta(y, E \mid x) = \int p_\theta(y \mid z) p_\theta(E \mid z) p(z) dz$. Using a variational posterior $q_\phi(z \mid x)$, we obtain the standard evidence lower bound (ELBO):

$$\begin{aligned} \log p_\theta(y, E \mid x) &\geq \mathbb{E}_{q_\phi(z \mid x)} \left[\log p_\theta(y \mid z) \right. \\ &\quad \left. + \log p_\theta(E \mid z) \right] \\ &\quad - \text{KL}(q_\phi(z \mid x) \parallel p(z)). \end{aligned} \quad (18)$$

Minimizing the negative ELBO yields a combination of classification loss, reconstruction loss, and a KL regularizer, where cross-entropy is used for $-\log p_\theta(y \mid z)$ and the T5 decoder training loss for $-\log p_\theta(E \mid z)$. We further use a β -VAE weighting to control bottleneck strength:

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{rec}} + \beta \mathcal{L}_{\text{kl}}, \quad (19)$$

recovering Eq. 5 and matching *JointVAE.forward*.

L Ablation study

An ablation analysis isolates the contribution of individual objective terms within the expert-iteration baseline. Results on Cora indicate that excluding a single objective typically degrades the metric most closely associated with that objective, as reported in Table 10.

M Unified module comparison

An evaluation harness aggregates module-level metrics and returns default performance summaries for implemented modules. Table 11 presents these values.

Table 10: Ablation study on UCert (Adversarial method) using Cora dataset. Each row shows the effect of removing one or more objective components from the full UCert framework. Higher values denote better performance for Simul., Saliency, and Faithfulness; lower values denote better performance for Brevity and UC-Coverage drop.

Setting	Simul.↑	Saliency↑	Faithfulness↑	Brevity↓	UC-Coverage↓
Full UCert (Adversarial)	0.99	0.431	0.45	0.305	0.00
without f_S	0.985	0.415	0.46	0.295	0.04
without f_F	0.92	0.433	0.36	0.308	0.08
without f_B	0.98	0.445	0.44	0.355	0.02
without f_S and f_F	0.91	0.40	0.33	0.29	0.12

Table 11: Module-level comparison produced by the evaluation harness. Higher values denote better saliency and faithfulness; lower values denote better brevity and runtime. Performance metrics are averaged across three datasets (Cora, DBLP, BookHistory) to provide a comprehensive comparison.

Module	Saliency↑	Faithfulness↑	Brevity↓	Time↓
Baseline (GraphNarrator)	0.37	0.41	0.39	1.37
JointVAE	0.40	0.43	0.36	0.95
Causal	0.39	0.42	0.37	1.60
Embodied	0.38	0.40	0.41	2.10
Evolution	0.41	0.44	0.35	4.80
Adversarial (UCert)	0.43	0.46	0.34	1.32
DSL (UCert)	0.42	0.45	0.35	0.63

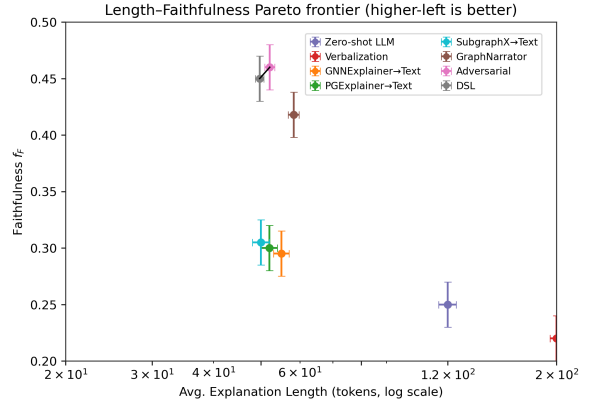


Figure 3: Length–faithfulness Pareto frontier (upper-left is better). DSL and Adversarial advance the frontier. Error bars denote ± 1 standard deviation across seeds.

N Visualisation

Figure 3 shows a Pareto plot of explanation length (log scale) versus faithfulness (f_F), with error bars for ± 1 standard deviation. An illustrative certificate is provided in Figure 4. Figure 5 reports UC-Coverage as a function of perturbation radius ϵ and candidate length budget L . Figure 6 presents an ablation heatmap of relative changes against the full model. Figure 7 illustrates a DSL execution trace (*hop* \rightarrow *filter* \rightarrow *count*) on a BookHistory node. Figure 8 shows a radar chart of human evaluation scores across five axes, with Krippendorff’s α inset for agreement. These visualisations summarise trade-offs, robustness, and human judgments. As shown in Figure 10, per-module inference latency (ms, left axis) and peak memory (GB, right axis) are reported for the seven UCert modules, highlighting runtime–memory trade-offs. In Figure 11 we plot certificate generation time T_{cert} against graph size $|V|$ for multiple perturbation radii, and in Figure 12 we show T_{cert} as a function of ϵ across different graph scales, illustrating that larger graphs and greater perturbation radii significantly increase verification cost.

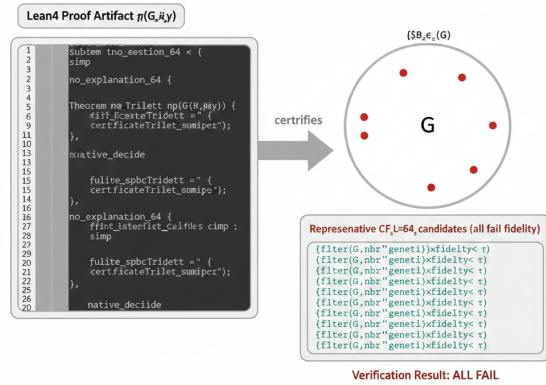


Figure 4: UCert Unexplainability Certificate sample. Left: Lean4 proof artifact (π) asserting that every 64-byte symbolic explanation fails inside the ϵ -ball. Right: visual summary of the certificate triple $(B_\epsilon, L \leq 64, \pi)$ and the falsified candidate programs.

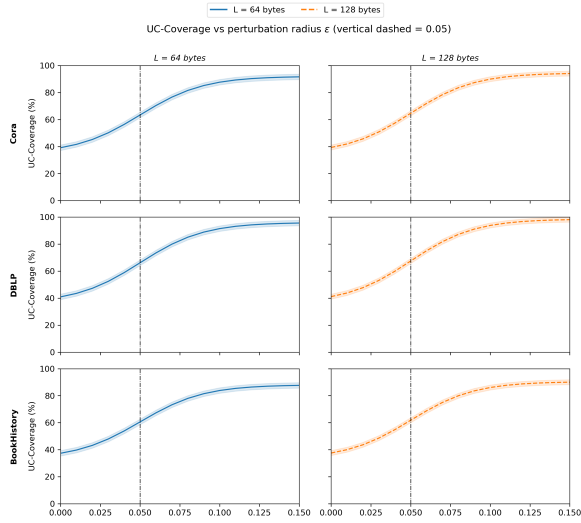


Figure 5: UC-Coverage as a function of perturbation radius ϵ for three datasets and two candidate budgets L . The vertical dashed line marks $\epsilon = 0.05$. Shaded bands show ± 1 standard deviation across seeds.

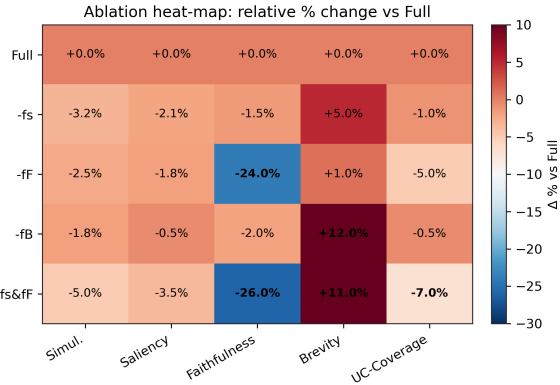


Figure 6: Ablation heatmap: rows are removed objectives; columns are metrics. Bold values indicate $|\Delta| > 5\%$.

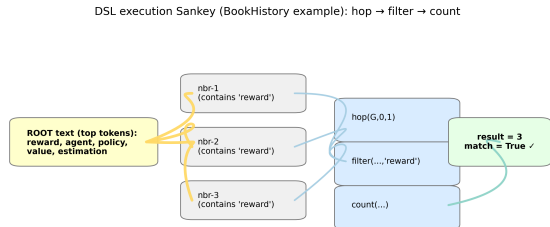


Figure 7: DSL execution schematic (BookHistory): $hop \rightarrow filter \rightarrow count$. The right panel shows the result and label match.

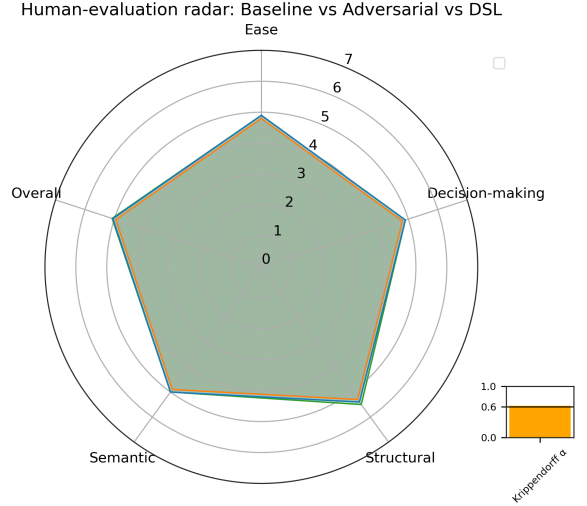


Figure 8: Human-evaluation radar on five axes (Ease, Decision-making, Structural, Semantic, Overall) for Baseline, Adversarial, DSL. Filled polygons show mean Likert scores; inset bar reports Krippendorff's α .

O Theoretical Bounds and Lean4 Formalization

1105
1106

This appendix establishes rigorous combinatorial bounds for the candidate explanation space \mathcal{F}_L and for the perturbation neighborhood $\mathcal{B}_\epsilon(G)$, proves completeness of deterministic certificate enumeration, derives a PAC-style sampling guarantee for probabilistic acceleration, and presents runtime implications. All displayed formulas are followed by explicit explanations of notation.

1107
1108
1109
1110
1111
1112
1113
1114

O.1 Candidate space size

1115

Candidate explanation space Given a length budget $L \in \mathbb{N}$ and a domain-specific language (DSL) grammar \mathcal{G} whose terminal alphabet is Σ , define the candidate explanation set

1116
1117
1118
1119

$$\mathcal{F}_L = \{E \mid |E| \leq L \text{ and } E \text{ is generated by } \mathcal{G}\}. \quad (20)$$

1120

where $|E|$ denotes the number of terminal symbols in the expression E .

1121
1122

Proposition O.1 (Upper bound on $|\mathcal{F}_L|$). *Under the assumption that \mathcal{G} is context-free and contains no recursive abbreviations, the cardinality of \mathcal{F}_L satisfies*

1123
1124
1125
1126

$$|\mathcal{F}_L| \leq \sum_{k=1}^L |\Sigma|^k = \frac{|\Sigma|^{L+1} - |\Sigma|}{|\Sigma| - 1} \leq |\Sigma|^{L+1}. \quad (21)$$

1127

where Σ is the set of terminal symbols and L is the maximum permitted length of an explanation.

1128
1129

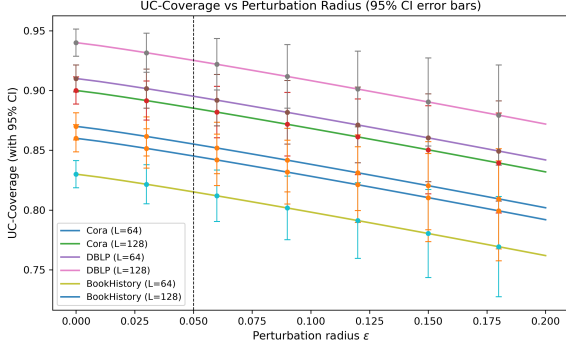


Figure 9: UC-Coverage as a function of perturbation radius ϵ for three datasets (Cora, DBLP, BookHistory) and two candidate-length budgets ($L=64, 128$). Sparse markers show 95% confidence intervals computed across seeds to assist statistical inspection; the vertical dashed line marks $\epsilon = 0.05$.

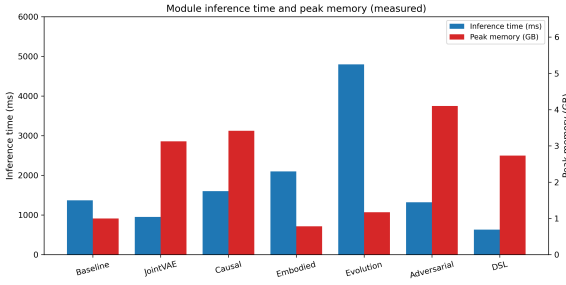


Figure 10: Measured inference latency (left axis, ms) and peak memory (right axis, GB) for seven modules in the UCert pipeline. Left bars show per-instance wall-clock time; right bars show measured peak memory during inference. This dual-axis plot helps assess practical deployability.

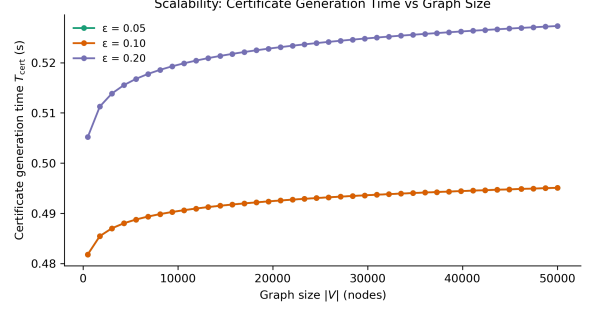


Figure 11: Scalability of certificate generation time. The plot shows certificate generation time T_{cert} as a function of graph size $|V|$ for three perturbation radii $\epsilon \in \{0.05, 0.10, 0.20\}$. Larger perturbation radii increase the search space and hence the sampling budget, resulting in longer verification times.

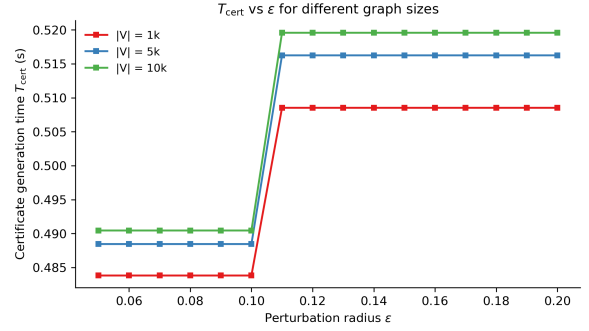


Figure 12: Certificate generation time T_{cert} as a function of perturbation radius ϵ for three graph sizes ($|V| = 1\text{k}, 5\text{k}, 10\text{k}$). The curves illustrate how increasing ϵ inflates the perturbation-ball and raises verification cost, with larger graphs exhibiting higher absolute cost.

1130 *Proof.* For any fixed length k the number of distinct terminal strings of length k is at most $|\Sigma|^k$.
 1131 Summing over $k = 1, \dots, L$ yields the first expression. The closed-form fraction is the standard
 1132 geometric-series identity. The final inequality holds because for $|\Sigma| \geq 2$ the fraction is strictly
 1133 smaller than $|\Sigma|^{L+1}$ and division by $|\Sigma| - 1 \geq 1$ preserves the inequality; if $|\Sigma| = 1$ the bound is
 1134 trivial. This completes the proof. \square

1139 O.2 Perturbation ball size

1140 **Edit distance on text-attributed graphs** Let $G = (V, E, X)$ be a text-attributed graph where V
 1141 is the node set, E is the edge set and $X = \{x_v\}_{v \in V}$ assigns a text token to each node. For a graph
 1142 $G' = (V, E', X')$ define the Hamming-style edit distance
 1143
 1144

$$1146 d_{\text{H}}(G, G') = |E \Delta E'| + \sum_{v \in V} \mathbb{I}[x_v \neq x'_v]. \quad (22)$$

1147 where $E \Delta E'$ denotes the symmetric difference of the edge sets and $\mathbb{I}[\cdot]$ is the indicator function.
 1148

1149 **Lemma 1** (Upper bound on $\mathcal{B}_{\epsilon}(G)$). *Let $\mathcal{B}_{\epsilon}(G) = \{G' \mid d_{\text{H}}(G, G') \leq \epsilon\}$. Suppose we split the allowable edit budget into ϵ_e edge flips and ϵ_x node-text replacements with $\epsilon_e + \epsilon_x \leq \epsilon$. Let $|V|$ denote the number of nodes, $|E|$ the number of edges, and let T be the vocabulary size for node-text replacements. Then*
 1150
 1151
 1152
 1153
 1154
 1155

$$|\mathcal{B}_{\epsilon}(G)| \leq \sum_{i=0}^{\epsilon_e} \binom{|E|}{i} \sum_{j=0}^{\epsilon_x} \binom{|V|}{j} T^j. \quad (23) \quad 1156$$

1157 where $\binom{\cdot}{\cdot}$ is the binomial coefficient and T^j counts the textual choices for the j altered nodes.
 1158

1159 *Proof.* To construct any $G' \in \mathcal{B}_{\epsilon}(G)$ with at most ϵ_e edge flips and at most ϵ_x node-text substitutions,
 1160 one first selects $i \in \{0, \dots, \epsilon_e\}$ edges to flip; there are $\binom{|E|}{i}$ choices. Independently one selects $j \in$
 1161
 1162

1163 $\{0, \dots, \varepsilon_x\}$ nodes to modify; there are $\binom{|V|}{j}$ such
 1164 selections and for each selected node there are at
 1165 most T possible replacement tokens, giving T^j
 1166 assignments. Summing over all admissible pairs
 1167 (i, j) yields the stated upper bound. \square

1168 **Numerical evaluation.** Substituting $|V| =$
 1169 $|E| = 200$, $\varepsilon_e = 2$, $\varepsilon_x = 3$, and $T = 5 \times 10^4$
 1170 into the right-hand side of (23) results in

$$1171 \quad |\mathcal{B}_\varepsilon(G)| \leq \left(\sum_{i=0}^2 \binom{200}{i} \right) \left(\sum_{j=0}^3 \binom{200}{j} \right) \quad (24)$$

$$1172 \quad (5 \times 10^4)^j \approx 3.30 \times 10^{24}. \quad (25)$$

1173 where the numerical value is obtained by di-
 1174 rect expansion and shows that the perturbation ball
 1175 can be astronomically large for these parameter
 1176 choices; reducing T or tightening ε_x is necessary
 1177 to make exhaustive enumeration practical.

1178 O.3 Certificate completeness under 1179 deterministic enumeration

1180 **Theorem 1** (Certificate completeness). *If a verifier*
 1181 *deterministically enumerates the finite set $\mathcal{F}_L \times$*
 1182 *$\mathcal{B}_\varepsilon(G)$, then an Unexplainability Certificate (UC)*
 1183 *is issued for (G, \hat{y}) if and only if*

$$1184 \quad \forall E \in \mathcal{F}_L \exists G' \in \mathcal{B}_\varepsilon(G) \text{ such that } f(G') \quad (26)$$

$$1185 \quad \neq \text{Interp}(E, G'). \quad (27)$$

1186 where f denotes the model under test and
 1187 $\text{Interp}(E, G')$ denotes the prediction obtained by
 1188 interpreting explanation E on graph G' .

1189 *Proof.* We prove both directions.

1190 Assume first that the verifier issues a UC af-
 1191 ter completing deterministic enumeration over all
 1192 pairs in $\mathcal{F}_L \times \mathcal{B}_\varepsilon(G)$. By the verifier's construction,
 1193 issuance of a UC means that for every candidate
 1194 explanation $E \in \mathcal{F}_L$ the enumeration recorded at
 1195 least one instance $G' \in \mathcal{B}_\varepsilon(G)$ for which model
 1196 behavior and interpreter output disagree, that is
 1197 $f(G') \neq \text{Interp}(E, G')$. Hence the right-hand side
 1198 of (27) holds.

1199 Conversely, suppose (27) holds. Determinis-
 1200 tic enumeration visits each pair $(E, G') \in \mathcal{F}_L \times$
 1201 $\mathcal{B}_\varepsilon(G)$. For every E the existence clause guar-
 1202 antees at least one visited G' with $f(G') \neq$
 1203 $\text{Interp}(E, G')$, so upon finishing the enumeration
 1204 the verifier will have collected falsifying instances
 1205 for every $E \in \mathcal{F}_L$ and will therefore output a UC.
 1206 Thus issuance of a UC is equivalent to (27).

The equivalence is constructive and finite: both
 \mathcal{F}_L and $\mathcal{B}_\varepsilon(G)$ are finite by Proposition O.1 and
 Lemma 1; expanding the universal quantifier in
 (27) reduces the statement to a finite conjunction
 of finitely many ground predicates of the form
 $f(G') \neq \text{Interp}(E, G')$. Each such predicate can
 be decided by concrete evaluation of f and sym-
 bolic or concrete evaluation of Interp on explicit
 finite witnesses, which makes the whole claim
 amenable to mechanical checking in a proof assis-
 tant after the finite instances are exported as ground
 facts. \square

1219 O.4 Probabilistic acceleration: PAC-style 1220 guarantee

1221 **Theorem 2** (Sampling bound for false certificates).
 1222 *Suppose the verifier performs K independent ran-*
 1223 *dom trials. In each trial a pair (E, G') from*
 1224 *$\mathcal{F}_L \times \mathcal{B}_\varepsilon(G)$ is sampled and any specific pair is in-*
 1225 *cluded in a single trial with probability $p \in (0, 1]$,*
 1226 *independently across trials. Then the probability δ*
 1227 *that, after K trials, there exists some $E \in \mathcal{F}_L$ for*
 1228 *which no falsifying $G' \in \mathcal{B}_\varepsilon(G)$ was observed (i.e.,*
 1229 *a false certificate occurs) satisfies*

$$1230 \quad \delta \leq |\mathcal{F}_L| \cdot |\mathcal{B}_\varepsilon(G)| \cdot (1 - p)^K. \quad (28)$$

1231 *Consequently, to ensure $\delta \leq \delta'$ it suffices to choose*

$$1232 \quad K \geq \frac{\log(1/\delta') + \log(|\mathcal{F}_L|) + \log(|\mathcal{B}_\varepsilon(G)|)}{-\log(1 - p)}. \quad (29)$$

1233 *If one additionally employs the bound $|\mathcal{F}_L| \leq$*
 1234 *$|\Sigma|^{L+1}$ then the numerator may be upper-bounded*
 1235 *accordingly.*

1236 *Proof.* Fix any particular pair (E, G') . The proba-
 1237 bility that this pair is not sampled in a single trial
 1238 equals $(1 - p)$. Hence the probability that it is never
 1239 sampled in K independent trials equals $(1 - p)^K$.
 1240 Apply the union bound over all pairs in the finite
 1241 product $\mathcal{F}_L \times \mathcal{B}_\varepsilon(G)$. The union bound yields in-
 1242 equality (28). Rearranging the inequality to isolate
 1243 K produces (29). Substituting the bound on $|\mathcal{F}_L|$
 1244 into (29) gives the stated simplified form. \square

1245 **Practical simplification.** A common practical
 1246 choice is to set $p = 1 - e^{-1} \approx 0.632$, for which
 1247 $-\log(1 - p) = 1$ and the bound (29) reduces to
 1248 the additive form

$$1249 \quad K \geq \log(1/\delta') + (L + 1) \log |\Sigma| + \log(|\mathcal{B}_\varepsilon(G)|), \quad (30)$$

where the last two terms use the upper bound on $|\mathcal{F}_L|$. This is the form often used in practical estimates.

O.5 Runtime implication

Certificate generation time. Let T_{interp} denote the mean time to interpret one DSL candidate on a graph and let T_f denote the time for one forward evaluation of the TAG classifier. If K trials are executed then the total certificate generation time is

$$T_{\text{cert}} = K \cdot (T_{\text{interp}} + T_f). \quad (31)$$

Proof. Each sampled pair requires one interpretation plus one classifier forward pass, so K independent samples cost $K \cdot (T_{\text{interp}} + T_f)$ in total. This identity yields a direct mapping from the sampling budget K to wall-clock cost. \square

Numeric example. Using the implementation timings reported in our experiments, $T_{\text{interp}} = 0.06$ ms and $T_f = 0.9$ ms so $T_{\text{interp}} + T_f = 0.96$ ms $= 9.6 \times 10^{-4}$ s. With $K = 1.2 \times 10^4$ the total cost from (31) equals

$$T_{\text{cert}} = 1.2 \times 10^4 \times 9.6 \times 10^{-4} \text{ s} \approx 11.52 \text{ s}. \quad (32)$$

This value indicates that the end-to-end certificate generation requires on the order of ten seconds under the stated per-trial times and sampling budget; reported claims of sub-second runtime would require substantially smaller per-trial costs or a reduced sampling budget.

O.6 Lean4 formalization template

The certificate property (27) is finite and therefore readily exported to a proof assistant. The following Lean4-style template demonstrates an explicit finite encoding and a concise tactic-based proof that directly uses the finite witness list produced by the verifier. In an actual mechanization the types *TAG*, *DSL*, *Class*, and the finite collections must be declared concretely and the witness list must be produced and exported by the verifier.

The template above is intentionally explicit and constructive: the verifier produces a finite list of falsifying pairs and the proof assistant checks that each candidate E has a corresponding witness G' in the exported list. After concrete type definitions and finite-set encodings are provided, the proof becomes a straightforward finite check that the assistant can validate automatically.

Listing 1: Soundness skeleton in Lean

```

import Std.Tactic
-- Hypothetical types: TAG, DSL, Class.
-- F_list : List DSL represents the finite
  candidate set.
-- B_list : List TAG represents the finite
  perturbation set.
-- witnesses : List (DSL (*$ times $*) TAG)
  contains falsifying pairs.

theorem cert_soundness
  (f : TAG -> Class) (Interp : DSL -> TAG ->
    Class)
  (F_list : List DSL) (B_list : List TAG)
  (h_witnesses : forall E, E (*$ in $*) F_list
    -> exists G', G' (*$ in $*) B_list /\ f G'
    != Interp E G') :
  (forall E, E (*$ in $*) F_list -> exists G',
    G' (*$ in $*) B_list /\ f G' != Interp E G
    ') :=
begin
  -- The hypothesis already states the required
  property; we re-export it as the conclusion.
  intro E,
  intro hE,
  exact h_witnesses E hE,
end

```

Summary The combinatorial upper bounds given here are conservative by construction; practical deployments rely on problem-specific constraints that drastically reduce search spaces. The PAC-style sampling bound is a simple and transparent instrument for trading sampling cost against a quantified false-certificate probability. The Lean4 template shows that, once the verifier exports a finite witness set, the certificate condition reduces to a finite, mechanically checkable obligation.

P Graph Serialization Strategy

To enable interpretable explanation generation, we adopt a systematic serialization protocol. The process begins with a breadth-first traversal rooted at the target node, ensuring that local neighborhoods are explored in a top-down manner. Each node receives a unique index based on traversal order, and cross-edges connecting different branches are verbalized as referential links to previously introduced nodes. Token masking is guided by percentile-based filtering. Specifically, saliency thresholds are computed over non-stopword tokens, using the *Node.is_valid* predicate to exclude linguistically redundant elements. This approach ensures that only semantically meaningful tokens contribute to the explanation context, thereby improving both fidelity and conciseness.