

# DETECTING INTERPROCEDURAL SOFTWARE VULNERABILITIES USING CAUSAL DEEP LEARNING

by

MD IQBAL HOSSAIN SHUVO

A thesis submitted to the  
Department of Computer Science  
in conformity with the requirements for  
the degree of Master of Science

Bishop's University  
Canada  
November 2025

Copyright © Md Iqbal Hossain Shuvo, 2025  
released under a [CC BY-SA 4.0 License](https://creativecommons.org/licenses/by-sa/4.0/)

# Abstract

Modern software stacks contain long, cross-module execution paths where subtle interactions between functions can create exploitable flaws, yet traditional and most learning based vulnerability detectors continue to operate as opaque classifiers that provide only coarse-grained labels over files or graph representations. This thesis introduces a causal, chain centric approach that couples vulnerability prediction with explicit reconstruction of executable root-to-sink chain across functions and files. The method starts from an augmented code property graph that integrates abstract syntax, control flow, data flow, and interprocedural links, then initializes node features with a pretrained structure aware encoder (GraphCodeBERT) and refines them using a relation-aware graph neural network with Adaptive Causal Contextualization guided by a mined Causal Knowledge Graph prior. A constrained beam search decoder, backed by feasibility checks over control, data, and call or return edges, assembles a single executable chain for each positive decision that respects program semantics, aliasing constraints, and interprocedural dependencies. Experiments on the ReposVul repository-level dataset evaluate detection performance and explanation quality using both standard metrics and specialized criteria. These evaluation criteria encompass the Counterfactual Consistency Score (CCS), the Causal Feature Attribution Measure (CFAM), and additional measures that check whether the chains are valid, cover the relevant interprocedural paths, match reference vulnerability chains, and respond in a sensible way to targeted program edits. This comprehensive evaluation approach ensures robust assessment of model accuracy, interpretability, and causal fidelity. The proposed approach not only achieves strong vulnerability detection performance but also effectively reconstructs a vulnerability chain. This chain explicitly captures the root cause (source), the intermediate propagation steps, and the vulnerable point (sink), and passes comprehensive feasibility checks, exhibiting improved coverage and correct ordering of these critical stages. This explicit root–propagation–sink representation in turn enables high-fidelity interprocedural explanations that align with developers’ diagnostic workflows and support root-cause reasoning in large-scale, heterogeneous software systems.

# Acknowledgments

I am deeply grateful to my supervisor, Dr. Yasir Malik, for his unwavering support, insightful guidance, and patience throughout the course of this research. His mentorship has been invaluable in shaping both my academic and professional growth. I also wish to express my heartfelt gratitude to my family, whose constant love, encouragement, and understanding have been essential in completing this thesis.

# Contents

<b>List of Tables</b>	<b>vi</b>
<b>List of Figures</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Research Gaps . . . . .	4
1.2 Research Objectives . . . . .	4
1.3 Contributions and Significance . . . . .	5
1.4 Thesis Outline . . . . .	6
<b>2 Literature Review</b>	<b>7</b>
2.1 Related Works . . . . .	7
2.1.1 Static and Dynamic Analysis Techniques . . . . .	7
2.1.2 Deep Learning Based Techniques . . . . .	9
2.1.3 Causal Reasoning and Causal Deep Learning Techniques . . . . .	10
2.1.4 Explainability and Interpretability Techniques . . . . .	12
2.2 Research Challenges and Rationale . . . . .	13
<b>3 Methodology</b>	<b>15</b>
3.1 Chain-Centric Program Representation . . . . .	16
3.1.1 Dataset and Repository Context: ReposVul . . . . .	16
3.1.2 Extraction and Graph Construction Pipeline . . . . .	17
3.1.3 Unified Multigraph, Typing, Features, and Storage . . . . .	19
3.1.4 Ground Truth . . . . .	19
3.2 Model Architecture and Decoding Pipeline . . . . .	20
3.2.1 GraphCodeBERT Feature Initialization . . . . .	21
3.2.2 GAT with Causality-Oriented Attention . . . . .	21
3.2.3 Causal Knowledge Graph (CKG): Mining and Prior . . . . .	22
3.2.4 Adaptive Causal Contextualization (ACC) . . . . .	22
3.2.5 Chain Extraction and Validation . . . . .	22
3.2.6 Training, Optimization, and Hyperparameters . . . . .	22

<b>4</b>	<b>Experimental Results and Analysis</b>	<b>24</b>
4.1	Experimental Setup and Configuration . . . . .	24
4.2	Evaluation Metrics . . . . .	27
4.2.1	Standard Classification Metrics . . . . .	27
4.2.2	Chain-Centric Metrics . . . . .	28
4.2.3	Counterfactual Metrics . . . . .	28
4.3	Detection Results . . . . .	29
4.4	Executable Chain Quality and Interprocedural Evidence . . . . .	30
4.5	Causal Faithfulness . . . . .	32
4.6	Efficiency and Decoding Dynamics . . . . .	34
4.7	Robustness and Generalization . . . . .	35
4.8	Test Case . . . . .	35
<b>5</b>	<b>Conclusion and Future Work</b>	<b>37</b>
5.1	Conclusion . . . . .	37
5.2	Future Work . . . . .	39
	<b>Bibliography</b>	<b>41</b>
	<b>Appendix A Algorithms and Pseudocode</b>	<b>47</b>
	<b>Appendix B Model Equations and Training</b>	<b>51</b>
B.1	GraphCodeBERT Feature Initialization: Equations . . . . .	51
B.2	Relation-Aware GAT: Equations . . . . .	52
B.3	Causal Knowledge Graph (CKG) Prior: Equation . . . . .	53
B.4	Adaptive Causal Contextualization (ACC): Equations . . . . .	53
B.5	CKG Prior Scoring . . . . .	55
B.6	Chain Extraction and Validation: Equations . . . . .	55
B.7	Training Objective: Equation . . . . .	55
	<b>Appendix C Metrics and Diagnostics</b>	<b>57</b>
C.1	Calibration and Thresholding . . . . .	57
C.1.1	Operating Points . . . . .	57
C.1.2	Expected Calibration Error (ECE) . . . . .	57
C.2	Standard Classification Metrics . . . . .	57
C.3	Chain-Centric Metrics . . . . .	58
C.3.1	Validity of Predicted Chains . . . . .	58
C.3.2	Structural Fidelity . . . . .	58
C.3.3	Interprocedurality (IPA) . . . . .	58
C.4	Counterfactual Metrics . . . . .	58
C.4.1	Counterfactual Consistency Score (CCS) . . . . .	58
C.4.2	Causal Feature Attribution Measure (CFAM) . . . . .	59
C.5	Decoding diagnostics . . . . .	59

<b>Appendix D Extended Results and Settings</b>	<b>60</b>
D.1 Hyperparameters and Environment . . . . .	60
D.2 Beam Search and ACC Diagnostics . . . . .	61
D.2.1 Decoder Diagnostics and Remedies . . . . .	62
D.3 Additional Diagnostics . . . . .	63
<b>Appendix E Dataset Card and Licensing</b>	<b>64</b>
<b>Appendix F Role Lexicon and Pattern Rules</b>	<b>65</b>
F.1 API Families and Examples . . . . .	65
F.2 Structural Patterns . . . . .	65
<b>Appendix G Reproducibility and Artifact</b>	<b>66</b>
G.1 Data provenance . . . . .	66
G.2 Graph construction . . . . .	66
G.3 Embeddings and caches . . . . .	66
G.4 Training and decoding configuration . . . . .	66
G.5 Environment and determinism . . . . .	67
G.6 Artifact contents . . . . .	67
G.7 Project layout . . . . .	67
G.8 Step-by-step reproducibility guide . . . . .	69
G.9 Determinism and validation checks . . . . .	71
G.10 Source for Figure 4.1 and Chain Trace . . . . .	71

# List of Tables

3.1	Core metadata recorded per ReposVul entry. . . . .	17
3.2	Preparation pipeline summary. . . . .	19
3.3	Node feature dimensions at initialization. . . . .	21
4.1	File-level label counts by split (C/C++ subset). . . . .	25
4.2	Graph instances and node-label density after chain-centric conversion. . . . .	25
4.3	Interprocedural connectivity (non-empty caller and callee sets). . . . .	25
4.4	Experimental Setup Details. . . . .	26
4.5	Comparison of Struct-only vs. GCBERT+Struct encodings. . . . .	26
4.6	Decoding and Training Configuration. . . . .	27
4.7	Valid and test metrics for Struct-only and GCBERT+Struct. . . . .	29
4.8	Operating thresholds and calibration parameters. . . . .	30
4.9	Confusion matrices on $\text{TEST}$ at $\tau_{\text{F1}}^*$ for Struct-only and GCBERT+Struct. . . . .	30
4.10	Feasibility of reconstructed chains (pass of all checks). . . . .	31
4.11	Interprocedural structure in predicted chains. . . . .	31
4.12	Role and edge agreement with ground truth. . . . .	31
4.13	Order agreement via LCS ratio. . . . .	32
4.14	Chain length and span: hops, files crossed, and summary-edge share. . . . .	32
4.15	Prototype causal metrics (means over $N=256$ graphs per split). . . . .	33
4.16	CCS by edit type (lower is better for minor, higher for sink edits). . . . .	33
4.17	Directional consistency rate (DCR; higher is better). . . . .	33
4.18	CFAM (overall and by chain segment). Values in $[0, 1]$ . . . . .	34
4.19	Chain invalidation rates and score deltas under interventions. . . . .	34
4.20	Inference latency per graph. . . . .	35
D.1	Beam and ACC diagnostics under $K=8, B=24, H=5$ (per graph). . . . .	62
D.2	Beam sensitivity (Valid, GCBERT+Struct; latency in ms per graph). . . . .	62
D.3	Decoder diagnostics for rare cases and suggested remedies. . . . .	63

# List of Figures

3.1	Overview of the proposed chain-centric vulnerability analysis approach. Data sources are lifted into a unified CPG) node features are initialized with a pretrained encoder, and an ACC-constrained, CKG-guided beam decoder performs structured reasoning to produce a single executable root-to-sink chain. . . . .	16
3.2	ReposVul pipeline: interprocedural CPG construction and leakage-safe data splits. . . . .	17
3.3	Chain-Centric Model Architecture and Inference Flow . . . . .	20
4.1	Executable interprocedural chain from root to sink across <code>main.c</code> and <code>lib/shell.c</code> . . . . .	36



# Chapter 1

## Introduction

Software vulnerability detection remains a critical yet complex challenge in software engineering. Modern digital infrastructure, financial services, scientific workflows, and everyday communication all depend on large, evolving software systems. As these systems scale in size and heterogeneity, contemporary codebases span millions of lines of code and are developed by distributed teams over long periods of time. This growth in complexity increases architectural coupling and the likelihood that subtle defects persist into production. When such defects are exploitable, they become security vulnerabilities that threaten confidentiality, integrity, and availability, and their effects may propagate across dependent components and services [38, 52, 37].

A bug is a defect that causes a program to behave incorrectly or unexpectedly. Not all bugs threaten security. A vulnerability is a defect that an adversary can exploit to violate confidentiality, integrity, or availability. Distinguishing benign errors from exploitable vulnerabilities requires understanding not only the defect itself but also the mechanism by which it can lead to a successful attack. This mechanism can be conceptualized as a chain of causation that begins at a root condition (e.g., untrusted input, unchecked buffer length, use-after-free) and propagates through assignments, parameter passing, return values, aliasing relations, and implicit flows induced by control dependence. The chain becomes exploitable when a tainted or unsafe state reaches a sensitive sink, such as a buffer write, command/SQL execution, deserialization routine, or privileged operation. Executability depends on dominance and post-dominance of guards, exceptional paths, resource lifetimes, and other feasibility constraints. This thesis therefore treats vulnerabilities as root→propagation→sink mechanisms to be reconstructed and validated end to end.

Traditional vulnerability detection techniques provide important, but incomplete, coverage of this problem space. *Static analysis* examines code without executing it, relying on rule-based systems, type reasoning, data-flow frameworks, and abstract interpretation to infer program behaviour at scale [29, 8, 41]. *Dynamic analysis* executes programs with concrete or instrumented inputs to observe runtime behaviour directly,

often using monitoring, tracing, or sandboxing to detect anomalous states [50, 6]. In practice, these approaches are complemented by *fuzzing* and penetration testing, which attempt to exercise programs under diverse or adversarial inputs in order to trigger faults and expose exploitable conditions. Static analysis offers broad coverage but tends to over-approximate feasible execution paths, producing false positives and limited insight into concrete exploitability, while dynamic analysis, fuzzing, and penetration testing provide high-fidelity traces but often struggle with input coverage, scalability, and reproducibility on large, modular systems [38, 37]. Recent years have seen a substantial shift toward *data-driven approaches* for software vulnerability detection. These methods, which include machine learning and deep learning based models, aim to automate vulnerability discovery by learning patterns from code representations such as token sequences, abstract syntax trees, control-flow graphs, and data-flow graphs. These heterogeneous representations are unified within Code Property Graphs (CPGs) that integrate syntactic, semantic, and control/data flow information via typed nodes and edges. Graph Neural Networks (GNNs) utilize relation-aware message passing and type-specific attention mechanisms to propagate information over these graphs, enabling fine-grained vulnerability classification and localization [24, 15, 5]. Pretrained models such as GraphCodeBERT enrich token embeddings with data-flow aware context and thereby improve generalization, while frozen language-model representations, when combined with relation-aware GNNs, enable effective propagation of causally relevant evidence over heterogeneous code graphs, strengthening learning-based vulnerability detection.

Despite this progress, existing static, dynamic, fuzzing, and data-driven techniques remain largely *correlation-centric* and are not sufficient to capture the full causal structure of real-world vulnerabilities. Traditional deep learning models tend to rely heavily on statistical associations in training data, often capturing superficial patterns rather than the underlying causes of vulnerabilities. They typically focus on isolated statements or small subgraphs that are predictive of the vulnerability label and emphasize *where* a flaw is likely to occur, rather than *how* it emerges and propagates through interacting control and data flows [38, 52]. This correlation-centric bias produces several limitations: models can become sensitive to dataset-specific artifacts such as naming conventions or formatting and may fail under benign transformations such as refactoring, formatting changes, or identifier renaming [38, 52]. Modeling such behaviors requires reasoning not only about local syntax and intra-procedural structure but also about interprocedural semantics, including call and return edges, argument-to-parameter bindings, return-to-caller relationships, and conservative aliasing across function and module boundaries [37]. Existing vulnerability detection techniques, including static and dynamic analysis, fuzzing and penetration testing, and modern data-driven detectors, therefore provide important but incomplete coverage of this interprocedural behaviour. Classical static, dynamic, and fuzzing-based solutions are effective at finding local defects or triggering specific failure scenarios, but they rarely find interprocedural vulnerability.

Similarly, most data-driven detectors are trained to classify localized code regions or limited graph neighborhoods and tend to focus on point wise predictions of *where* a vulnerability may occur, rather than modeling *how* vulnerable state propagates across functions and modules in a causally meaningful way [38, 52, 37].

In real software systems, vulnerabilities often arise from complex interactions that span multiple functions, modules, or files. Tainted input may flow through several layers of wrapper functions before reaching a sensitive sink; resource management errors may involve initialization in one component and misuse in another; and access-control checks may be applied inconsistently across different call paths. These behaviours are further complicated by the use of shared libraries, callbacks, framework conventions, and configuration-driven control flow, where the relevant control and data dependencies are distributed across APIs, modules, and sometimes even language boundaries. As a result, the conditions that make a defect exploitable often depend on long, cross-cutting chains of calls and data transformations that are only partially visible in any single local context. Existing static, dynamic, fuzzing, and data-driven techniques typically approximate these patterns using limited context or heuristics [38, 37].

To address these challenges, this thesis adopts a *causal and context-aware* perspective on vulnerability detection. In this approach, *context awareness* refers to the capability of a model to understand how program elements relate not only syntactically, but also semantically, through control flow, data dependencies, and functional interactions [51, 13]. Vulnerabilities often do not originate from a single isolated statement; instead, they emerge from the interplay of multiple operations over time and across function boundaries, so capturing this semantic context is essential for reconstructing meaningful explanations of vulnerability behaviour, especially in interprocedural scenarios [37]. At the same time, the thesis builds on principles from *causal inference*, which seek to distinguish true cause–effect relationships from spurious correlations and have motivated recent causal and explainable approaches to vulnerability detection [4, 18, 35, 7]. In the context of vulnerability detection, the integration of causal reasoning within program analysis is both strategically advantageous and fundamentally essential. Effective remediation hinges on accurately identifying the specific code elements and propagation pathways that are causally responsible for vulnerabilities. Accordingly, this thesis adopts an approach that develops a context-aware, data-driven causal framework to reconstruct interprocedural root-to-sink vulnerability propagation paths within source code. This approach seeks to harmonize vulnerability prediction with executable, causally grounded explanations, elucidating not only the presence of vulnerabilities but also the precise manner in which they emerge and disseminate across functions and modules. Such causally informed interpretations provide a necessary foundation for trustworthy, actionable vulnerability diagnostics and remediation guidance.

## 1.1 Research Gaps

Software vulnerability detection has advanced considerably, yet several critical gaps remain in current data-driven and interprocedural analysis techniques.

**Gap 1: Causal deficiency and generalization:** Many data-driven detectors learn superficial statistical correlations rather than the underlying flaw mechanisms. Their predictions often deteriorate under benign code refactorings and cross-project generalization, reflecting insufficient causal grounding and poor robustness to distribution shifts [38, 52].

**Gap 2: Interprocedural reasoning and chain reconstruction:** Most existing approaches focus on individual functions or limited intra-procedural slices and lack interprocedural dependencies [37], thereby failing to reconstruct executable root-to-sink vulnerability propagation chains across large, modular software systems.

## 1.2 Research Objectives

To systematically advance the understanding and detection of software vulnerabilities, this thesis focuses on following research questions and corresponding objectives.

### Research Questions

**RQ 1:** How can a program representation be designed that faithfully captures interprocedural control and data dependencies in a form that supports executable, root-to-sink causal analysis across functions, files, and modules?

**RQ 2:** How can a context-aware causal reasoning approach systematically construct executable interprocedural vulnerability chains that are viable, minimal, and consistent with the underlying control flow, data dependencies, and aliasing relationships?

### Research Objectives

To achieve the goal of reconstructing executable interprocedural vulnerability propagation chains, this thesis pursues the following research objectives.

**O 1:** Design a program representation that captures interprocedural control and data dependencies, combining structures such as Abstract Syntax Trees (ASTs), Control-Flow Graphs (CFGs), Data-Flow Graphs (DFGs), and Program/Code Property Graphs (PDG/CPG), in a form suitable for root-to-sink causal analysis across functions, files, and modules.

**O 2:** Create a constraint-based approach to generate executable, interprocedural vulnerability chains, explicitly representing control flow, data dependencies, and aliasing to produce chains that are viable, minimal, and mechanically valid.

This thesis reconstructs executable interprocedural root-to-sink vulnerability propagation chains from real-world vulnerable repositories. The codebase is parsed into Abstract Syntax Trees, Control-Flow Graphs, and Data-Flow Graphs, which are integrated into a unified Program/Code Property Graph augmented with interprocedural semantics (call/return edges, argument→parameter mappings, return→caller links, and conservative aliasing) so that candidate root-to-sink paths can be followed across functions, files, and modules. Nodes in this heterogeneous graph are encoded using a structure-aware pretrained model (GraphCodeBERT) combined with structural features, and these embeddings are refined by a graph-based encoder equipped with attention, Adaptive Causal Contextualization (ACC), and a Causal Knowledge Graph (CKG) [57] prior. On top of this encoder, a constrained beam search decoder explores feasible control- and data-flow paths to assemble executable root-to-sink chains consistent with the interprocedural semantics.

### 1.3 Contributions and Significance

This research transforms vulnerability explanation from localized, retrospective feature highlights into a unified, executable interprocedural causal chain that traces the propagation from root cause to exploitable sink. This approach helps bridge the gap between traditional static and dynamic analyses, fuzzing and penetration testing, and modern data-driven vulnerability detectors, aligning automated analysis more closely with the complexity of large, modular software systems. The research makes model predictions interpretable by reconstructing how untrusted input enters the system, travels through control and data dependencies across functions, avoids or skips sanitization, and reaches a vulnerable operation. These chains give developers clear, traceable proof of a vulnerability’s existence, as well as clear instructions on where to find it and what to do to fix it. This transforms black-box detections into actionable and verifiable information. This enhances the reliability, transparency, and practical utility of learning-based vulnerability detection, contributing to software security in real-world environments.

Taken together, this problem framing and methodological perspective give rise to the following contributions in interprocedural, causality-aware software vulnerability analysis.

**C1 - A chain-centric interprocedural program representation:** This thesis extends the code property graph by integrating explicit interprocedural semantics, preserving control and data dependencies across functions and files to enable reconstruction of executable root-to-propagation-to-sink chains.

**C2 - Interprocedural causal vulnerability chain construction:** A context

driven causal pipeline that constructs explicit interprocedural root-to-sink vulnerability chains. This approach constructs the full propagation path (root→propagation→sink) as the main explanatory unit.

## 1.4 Thesis Outline

This thesis is structured as follows. Chapter 1 introduces the problem context, articulates the research gaps, formulates the research questions and objectives, and summarizes the main contributions. Chapter 2 reviews related work in the software vulnerability detection domain. Chapter 3 describes the methodology, including dataset preparation, interprocedural graph construction, chain-centric representation, and the model architecture with its decoding pipeline. Chapter 4 presents the experimental setup, and detailed results focusing on detection accuracy, chain quality, causal analyses, robustness, and case studies. Chapter 5 concludes with key findings, implications, limitations, and future research directions. Appendices provide supplementary materials such as algorithms, model equations, metrics, extended results, and reproducibility details.

## Chapter 2

# Literature Review

The growing need for more reliable and interpretable software vulnerability detection methods has driven extensive exploration of deep learning (DL) based approaches in recent research. These methods aim to automate the identification of vulnerabilities by learning patterns from code representations, often outperforming traditional static and dynamic analysis. However, despite their success, deep learning models still face challenges in scalability, generalizability, and transparency. To address these issues, researchers have investigated both the refinement of DL-based detection techniques and the development of methods to enhance model explainability and causal reasoning capability. This pursuit of more robust and interpretable models leads directly to the exploration of advanced DL architectures and explainability techniques discussed in this chapter.

## 2.1 Related Works

### 2.1.1 Static and Dynamic Analysis Techniques

Static analysis examines software artifacts without executing them and seeks to reason soundly or semi-soundly about all possible program behaviours. Classical approaches rely on lattice-based data-flow frameworks, abstract interpretation, and type- or effect-based reasoning to approximate control and data dependencies at scale [51, 5]. Their precision depends on sensitivities to control flow, calling context, execution paths, and heap or field abstractions. Interprocedural static analysis, in particular, requires accurate call-graph construction, supported by techniques such as Class Hierarchy Analysis, Rapid Type Analysis, and points-to analysis to resolve aliasing and dynamic dispatch [1, 25]. These analyses operate over canonical program representations: Abstract Syntax Trees (ASTs) encode syntactic structure, Control-Flow Graphs (CFGs) capture feasible execution order and branching, Data-Flow Graphs (DFGs) trace definition–use chains of variable values, and Program Dependence Graphs (PDGs) unify data and control dependencies to

support slicing and taint tracking. The Code Property Graph (CPG) integrates AST, CFG, and DFG into a heterogeneous multigraph and has been widely adopted to support scalable, interprocedural vulnerability queries from sources to sinks [51].

Despite their ability to detect potential issues early, classical static analyses routinely suffer from false positives due to conservative over-approximation of sanitization logic, aliasing, reflection, and dynamic dispatch [8, 29]. Comparative studies of widely used static analyzers, such as Flawfinder, RATS, Cppcheck, SpotBugs, and PMD, report substantial variability in detection rates and false-positive profiles across languages and vulnerability categories [17]. These limitations have motivated integration of machine learning and deep learning into static analysis pipelines. For example, IRIS combines large language models (LLMs) with static analysis to infer taint specifications and expand coverage of vulnerability patterns while reducing false discoveries, outperforming rule-centric tools such as CodeQL on several benchmarks [23]. Other hybrid approaches use ML classifiers to post-process static warnings, filter likely false positives, or prioritise alerts for manual inspection, thereby improving scalability and analyst productivity [11]. Static analysis thus remains a foundational component that provides structured, semantically rich program views to downstream learning-based detectors, even as its limitations on complex, rapidly evolving, and domain-specific codebases remain an open challenge [1, 5].

Dynamic analysis complements static reasoning by monitoring concrete or symbolic program executions to detect vulnerabilities at runtime. Representative techniques include fuzzing, which automatically mutates inputs to explore diverse execution paths; runtime sanitizers, which instrument code to detect memory safety violations and undefined behaviour; and symbolic or concolic execution, which uses constraint solvers to systematically explore feasible paths leading to potential faults [50, 49]. Dynamic methods produce high-fidelity witness traces that are directly actionable for debugging and exploit validation, but they are constrained by path explosion, environment modelling issues, and limited input or configuration coverage, especially in large-scale interprocedural systems [50]. To mitigate these limitations, recent work has proposed hybrid static–dynamic workflows in which machine learning and reinforcement learning guide fuzzing campaigns, prioritise target functions, or learn path selection heuristics from execution feedback, thereby improving coverage and discovery rates [48, 42]. Overall, static analysis offers scalability and whole-program coverage, whereas dynamic analysis provides precise execution evidence; their integration, increasingly mediated by data-driven techniques, forms a central strand in contemporary research on automated software vulnerability detection.



### 2.1.2 Deep Learning Based Techniques

Deep learning is reshaping many areas of computer science, and graph neural networks (GNNs) have rapidly shown strong performance in software vulnerability detection [5, 25, 56, 15]. By modeling code as graphs whose nodes represent program elements and edges encode syntactic and semantic relations, GNN-based models can capture complex software dependencies and localize vulnerable statements, for example in the line-level setting demonstrated by Hin et al. [15]. However, many of these models operate as opaque “black boxes,” making it difficult to understand why a particular code fragment is classified as vulnerable [22, 31], even when prediction accuracy is high. In security-critical contexts, this lack of transparency is problematic: developers must understand *why* a vulnerability is reported in order to implement appropriate fixes, and recent policy initiatives, such as the U.S. Executive Order on safe, secure, and trustworthy AI, explicitly emphasize the need for reliable and explainable AI behaviour in high-stakes systems [43].

Zhou et al. [24] conducted a systematic empirical study of design factors that influence the performance of deep learning based vulnerability detectors. They constructed two datasets capturing both data and control dependencies for 126 vulnerability types and implemented a Joern based detection pipeline to enable controlled comparisons. Their experiments quantified the impact of several key choices, including techniques for handling class imbalance, the inclusion of control dependence information in code representations, and the selection of neural network architectures, on downstream detection accuracy. The results provide practical guidance for building effective DL based vulnerability detection systems, but the work remains correlation oriented and does not explicitly incorporate causal deep learning or formal causal reasoning principles.

Li et al. [38] performed empirical research on deep learning (DL) models for vulnerability detection. On Devign and MSR data, they polled and replicated nine state-of-the-art (SOTA) models. They looked at and examined model capabilities (agreement, variability, performance on several vulnerability categories, and difficulties in addressing particular code aspects). They investigated how model performance responded to project mix and training data size. They identified significant code aspects applied for prediction using model explanation tools. The main contribution of the authors was a thorough investigation of models of deep learning vulnerability detection. For assessing and contrasting several deep learning models for vulnerability identification, this research offers a useful standard. Still, this study paid little attention to causative factors. The authors urged more investigation on code patterns, possible inclusion of causal detection for better generalization, and addressing of the shortcomings of present model interpretation techniques.

### 2.1.3 Causal Reasoning and Causal Deep Learning Techniques

Conventional deep learning based vulnerability detectors are highly sensitive to spurious correlations and dataset shift. VulCausal, presented by Kuang et al. [18], introduces a causal perspective by defining a structural causal model over user-defined identifiers, API library calls, and code structure, then applying backdoor-style adjustments during reasoning to suppress confounding and remove misleading associations. This yields more stable and accurate vulnerability predictions than correlation-driven baselines, but the approach still faces practical limitations, including scalability to very large codebases, limited language coverage, and the absence of explicit executable causal chains beyond improved metrics. More broadly, such methods highlight both the promise of causal inference for vulnerability detection and the substantial modelling effort required, in line with ongoing work on learning more robust, causally grounded optimizers and representations [55].

Le et al. [37] examined how existing deep learning based detectors handle vulnerabilities that span multiple basic units of code (MBUs) and showed that models trained and evaluated at the individual basic unit (IBU) level systematically struggle with such cases, overstating accuracy by focusing on isolated units. Although the work does not employ explicit causal reasoning, it provides a detailed characterization and taxonomy of MBU vulnerabilities and analyses their incidence and detection accuracy in modern DL based detectors. The main contribution is a framework for correctly integrating MBU vulnerabilities into DL based detection pipelines, underscoring the need to evaluate models across a spectrum of granularities, which is particularly important in contemporary multi file, multi module software systems. Complementary efforts, such as Nong et al. [33] on authentic vulnerability generation via pattern mining and deep learning, and Woo et al. [47] on identifying 1 day vulnerabilities in reused open source components, further refine understanding of vulnerability types and their prevalence in real world code.

Cao et al. [4] introduced Snopy, a deep learning method combining change-based sample denoising using vulnerability fixing commits with causal graph learning to enhance vulnerability detection. Snopy employs a Causality Aware Graph Attention Network and Feature Caching Scheme to focus on causal features while suppressing spurious ones. This approach provides a principled way to distinguish causal from non-causal code elements, improving robustness beyond correlation-based detectors. While challenges remain in generalization and spurious correlation mitigation, Snopy represents an important advance toward causal vulnerability detection.

Ganz et al. [12] addressed key practical challenges in ML-driven vulnerability detection by focusing on data quality, model interpretability, robustness, and context sensitivity. They introduced novel neural code augmentation to enhance datasets and leveraged dynamic program analysis to evaluate explanation methods. Their causal learning-based bias assessment framework helped reduce confounding effects and improve detection robustness. This work advances learning-based vulnerability detection toward practical relevance by emphasizing data quality

and interpretability. Related work by Suneja et al. [39] and Yu et al. [54] further explores model interpretability through prediction-preserving input minimization and counterfactual analysis, respectively.

Rahman et al. [35] introduced CausalVul to address the lack of robustness and out-of-distribution generalisation in deep learning based vulnerability detection. They modelled the relationships between code features and vulnerability labels with a causal graph, then used do-calculus and the backdoor criterion in a two-stage pipeline to identify and suppress spurious correlations while preserving causally relevant features. This causal adjustment yields models that rely less on dataset-specific artefacts and more on mechanisms genuinely associated with vulnerabilities, improving both resilience and generalisation. The work represents a clear step toward causal deep learning for vulnerability detection by explicitly targeting false correlations through principled causal inference tools.

Chu et al. [7] addressed the lack of explainability in Graph Neural Networks (GNNs) for vulnerability identification. Using counterfactual reasoning, a type of causal inference, CFExplainer aimed to identify minimal changes to the input code graph that would influence the GNN prediction. This improved explainability by focusing on behaviours that alter the outcome and so provides a "what-if" analytical capability. The approach finds a minimum disruption in the code graph by inverting the prediction of the GNN. This provides developers with actionable insights and guides the necessary modifications to remediate the vulnerability. The primary contribution of the authors was the development of a counterfactual explanation method for GNN-based vulnerability discovery. Counterfactual thinking helps developers find a more reasonable and workable argument. A better basis for understanding the application of counterfactual reasoning in explainable artificial intelligence is provided by the work of Lucic et al. [27].

Islam et al. [16] proposed T5-GCN for vulnerability categorization, localization, and root cause identification by combining graph convolutional networks with a T5-based large language model. Although the work uses the term "root cause," it does not employ causal deep learning, causal inference, or formal causal reasoning; instead, it relies on DeepLift-SHAP attribution scores to identify tokens that most strongly contribute to the model's prediction, providing an attribution-based explanation of the vulnerability. The main contribution lies in coupling GCNs with an LLM for code and using explainability techniques to highlight likely root-causing code fragments alongside vulnerability type, location, and a brief static description. This line of work aligns with growing interest in LLM-based code analysis and optimisation [55] and complements efforts such as Pearce et al. [34], who assess the security of code produced by GitHub Copilot, indicating the broader potential of LLMs in vulnerability detection and secure software engineering.

Cao et al. [3] presented Coca, a framework to enhance the causality and robustness of Graph Neural Networks (GNNs)-based vulnerability detection systems. Coca used dual-view causal inference, that is, factual and counterfactual reasoning, to

pinpoint code statements most likely to be decisive for vulnerability discovery. This method showed a sophisticated use of causal ideas to solve constraints in robustness and explainability observed in past GNN-based detectors. It also underlined the difficulties in juggling concision with effectiveness in explanations. Coca trains GNNs less prone to false correlations and more focused on real vulnerability traits by means of combinatorial contrastive learning. The Explainer component generates succinct and powerful explanations using dual-view causal inference. Using supervised contrastive learning, the system is taught to identify the bug in all versions and distinguish between buggy and non-buggy code. It discovers a flaw and then employs factual and counterfactual thinking. This is a framework enhancing the causality and dependability of GNN-based vulnerability detection systems.

#### 2.1.4 Explainability and Interpretability Techniques

Nguyen et al. [32] introduced GAVulExplainer, one of the earliest efforts to enhance explainability in deep learning-based vulnerability detection. This model-agnostic framework applies to various GNN detectors and uses genetic algorithms to identify subgraphs that most influence a model’s prediction by isolating coherent subsets of the Code Property Graph responsible for labeling vulnerabilities. Explanation quality is evaluated with a fidelity metric, and users can adjust explanation size to balance conciseness and completeness. While GAVulExplainer significantly advances interpretable GNN vulnerability detection, it focuses on influential subgraphs rather than explicitly modeling causal relationships between code features and vulnerabilities. Nonetheless, it emphasizes the critical role of explanation fidelity for effective remediation and fostering developer trust.

Li et al. [22] presented Vulanalyzer, a deep learning based framework that targets explainable binary vulnerability detection, a challenging setting due to low level representations and sparse semantic cues. The model combines sequential and topological learning to approximate program execution in assembly, using recurrent units and graph convolutions to preserve instruction semantics and control flow, in line with recent work on functional summaries of assembly code [41]. A multi head attention mechanism highlights influential instructions and basic blocks, guiding developers toward code regions the model considers most indicative of a vulnerability and thereby improving interpretability. However, similar to GAVulExplainer, Vulanalyzer does not incorporate explicit causal reasoning, causal inference, or causal deep learning; attention primarily clarifies *what* features drive predictions rather than *why* those features are causally responsible for the vulnerability [2].

Moschitti et al. [21] presented an XAI based framework for analysing source code in a graph setting, with a focus on how syntactic structures contribute to Common Weakness Enumeration (CWE) classification [20]. Their method ranks

AST constructs by their contribution to specific CWE types for Java and C++ by masking neighbourhoods of code tokens and syntactic nodes in the graph and mapping representation changes to CWE similarity scores via information retrieval techniques. This links learned code feature representations to nuanced semantics recognised by security experts. The authors also identify key limitations of current XAI approaches, including limited generalisation to unseen vulnerability patterns, difficulties in interpreting learned features, and poor transferability across datasets and languages, echoing broader challenges in representing and evaluating large scale code corpora for machine learning on code [1].

Hajipour et al. [19] presented a framework for vulnerability threat prediction that builds semantic representations of vulnerabilities using topic modeling over textual descriptions from sources such as the NVD, and derives an explainable threat score to prioritise analysis. In parallel, they introduced a trend score based on infosec related online discussions to capture emerging, high interest vulnerabilities, aggregating both scores in a visual dashboard to guide security teams toward higher priority cases. While effective for triaging and prioritisation, the approach operates at the description and discussion level, relying on external text and community signals; it does not model the underlying *causal* code level factors that drive exploitability, and may inherit biases and timeliness limitations from its data sources.

Marchetto et al. [29] evaluated how deep learning and explainability methods, in particular SHAP, support vulnerability localization in source code. Using VulDeePecker and JavaBERT as representative DL detectors, they applied SHAP to highlight statements associated with predicted vulnerabilities and found that the resulting explanations often surfaced noisy or only weakly related code features, which could distract rather than assist engineers. This empirical result suggests that current XAI techniques for DL based vulnerability detection frequently expose correlational signals rather than the true underlying causes, underscoring the need for more precise localization mechanisms and the integration of causal reasoning. The authors argue that progress will depend on programming language aware encoding schemes and more advanced techniques specifically tailored to vulnerability localization, rather than generic feature attribution alone.

## 2.2 Research Challenges and Rationale

Despite substantial progress in deep learning based vulnerability detection, empirical studies consistently report performance degradation under benign refactorings and cross-repository transfer, indicating that many models remain sensitive to spurious, project-specific correlations rather than underlying flaw mechanisms [38, 52]. This limitation highlights a central research challenge: moving beyond correlation centric detectors toward approaches that privilege executable behaviours and causally meaningful flows.

Recent work has begun to introduce causal notions at different levels of the detection pipeline. Contrastive and counterfactual training and inference schemes, such as *COCA* and *CFExplainer*, use counterfactual variants of code to probe model decisions and reduce reliance on non-robust features [3, 7]. Other approaches, such as *Snopy*, incorporate causal adjustment and denoising by suppressing non-causal context through change-based sample filtering and causality-aware graph attention, thereby attenuating spurious correlations in the learned representations [4]. Structural-causal formulations, exemplified by *VulCausal* and *CausalVul*, explicitly model relationships between code features and vulnerability labels using causal graphs, then apply backdoor style corrections to prioritise features with genuine causal influence over dataset artefacts [18, 35]. Collectively, these research works demonstrate the potential of causal inference to improve robustness and generalisation, but they typically operate at the level of feature importance or local predictions rather than reconstructing full executable vulnerability chain.

Parallel efforts in explainable AI (XAI) for software security further underscore the gap between saliency-style rationales and truly causal explanations. Studies show that standard attribution and explanation techniques can surface misleading or weakly related code elements, failing to provide developers with trustworthy, intervention-stable justifications for model predictions [29, 20, 21]. These findings highlight two intertwined challenges: evaluating whether explanations are faithful to the model’s decision process, and assessing whether they align with executable program behavior under interventions such as code edits, refactoring, and cross-project transfer.

These observations motivate the chain centric rationale adopted in this research. Rather than treating explanations as sets of important features or subgraphs, the proposed approach seeks to use a unified, interprocedural program graph, learn relation-aware evidence propagation over this graph, and then *constrain* decoding so that the resulting explanation is a single executable root→propagation→sink chain. Causal priors and structured decoding ensure learned scores translate into behaviors consistent with program semantics and robust to counterfactual edits. This thesis integrates causal reasoning, interprocedural program representation, and chain-centric explanations to effectively address challenges in vulnerability detection.



## Chapter 3

# Methodology

This chapter presents the learning pipeline and approach used for vulnerability spreading and chain reconstruction. The proposed technique combines causal attention with context-aware graph reasoning over an enhanced CPG that captures both structural and interprocedural relations. It connects root causes to propagation and endpoints via paths that respect program semantics, using pretrained feature initialization, attention-based aggregation on the graph, and causal contextualization to reconstruct accurate, executable root-to-sink chains that are both robust and interpretable.

The pipeline initiates with the construction of an augmented CPG derived from real-world vulnerable repositories. Source code is parsed into ASTs, CFGs, and DFGs, which are subsequently integrated into a unified CPG. This foundational graph is then enriched with interprocedural components, including call-graph edges, argument-to-parameter and return-to-caller mappings, call-site context capturing receiver and dispatch behavior, and coarse alias or points-to summaries, thereby capturing comprehensive data and control flows across procedural boundaries. By progressively enriching the graph in this way, the representation captures how values and control move across function boundaries and modules, so that vulnerability relevant flows can be followed end-to-end rather than being confined to isolated functions. Leveraging the enriched CPG, a structure-aware pretrained encoder (GraphCodeBERT) provides initial token level semantics, which are fused with node level structural features. A relation-aware graph encoder with attention and Adaptive Causal Contextualization (ACC), guided by a CKG prior [57], propagates vulnerability evidence along heterogeneous edges that reflect control, data, and call/return relations. Finally, a constrained beam search decoder explores feasible control- and data-flow paths to assemble executable root-to-sink chains that remain consistent with the interprocedural semantics; the resulting chains and predictions are evaluated using standard detection metrics together with dedicated chain-level validity criteria.

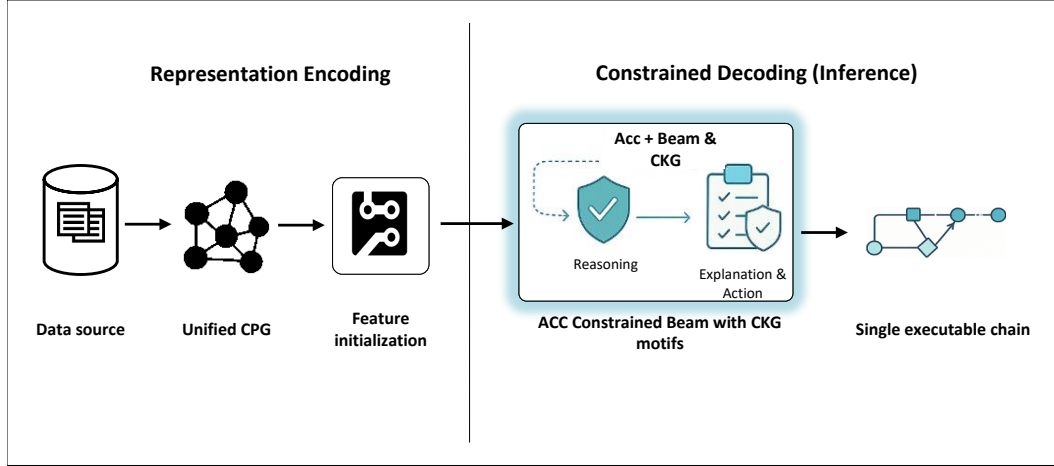


Figure 3.1: Overview of the proposed chain-centric vulnerability analysis approach. Data sources are lifted into a unified CPG) node features are initialized with a pretrained encoder, and an ACC-constrained, CKG-guided beam decoder performs structured reasoning to produce a single executable root-to-sink chain.

### 3.1 Chain-Centric Program Representation

A chain-centric, causal approach to vulnerability detection requires a program representation that explicitly models interprocedural mechanisms. Traditional AST, CFG, DFG, and baseline CPG representations enable local reasoning but lack explicit treatment of root-to-sink propagation paths, which limits executable chain reconstruction across functions, files, and modules. This thesis adopts an augmented CPG that explicitly encodes taint sources, transformations, sanitization, and sinks through control-flow and data-flow, while preserving call/return and aliasing semantics. The representation is designed to be compatible with relation-aware graph encoders and constrained decoding, allowing evidence to propagate along heterogeneous edges and enabling extraction of executable vulnerability chains that support interprocedural causal reasoning.

#### 3.1.1 Dataset and Repository Context: ReposVul

The chain-centric graphs are derived from *ReposVul*, a repository-level benchmark that links vulnerabilities to concrete pre- and post-fix revisions and retains repository-wide context [45]. Unlike patch-only corpora, *ReposVul* preserves full file snapshots and cross-file dependencies, which is essential for reconstructing long-range interprocedural flows.



Each *ReposVul* entry links CVE/CWE records to precise pre- and post-fix repository snapshots rather than isolated patches [45]. A documented weakness and its severity are anchored to the exact commit, preserving all related files and cross-file dependencies. The corpus is constructed by crawling public sources, disentangling mixed commits to isolate fix-relevant code, and extracting caller-callee relations across translation units. By maintaining chronology and repository-wide call structure, *ReposVul* reveals long-range interprocedural flows essential for chain-centric vulnerability analysis. Table 3.1 summarises the core metadata captured for each entry, including vulnerability descriptors, patch-level information, and file-level snapshots.

Table 3.1: Core metadata recorded per ReposVul entry.

Category	Representative fields
Vulnerability entry	CVE-ID, CWE-ID, language, external references, CVE description, publish date, CVSS vector (AV, AC, PR, UI, S, C, I, A)
Patch metadata	Commit ID, commit message and date, project/repository IDs, parent/child links, forge URLs
Related files	File name, language, vulnerable/fixed snapshots, line diffs, file URLs

### 3.1.2 Extraction and Graph Construction Pipeline

The conversion from raw dataset entries to chain-centric program graphs occurs over six sequential stages, as illustrated in Figure 3.2.

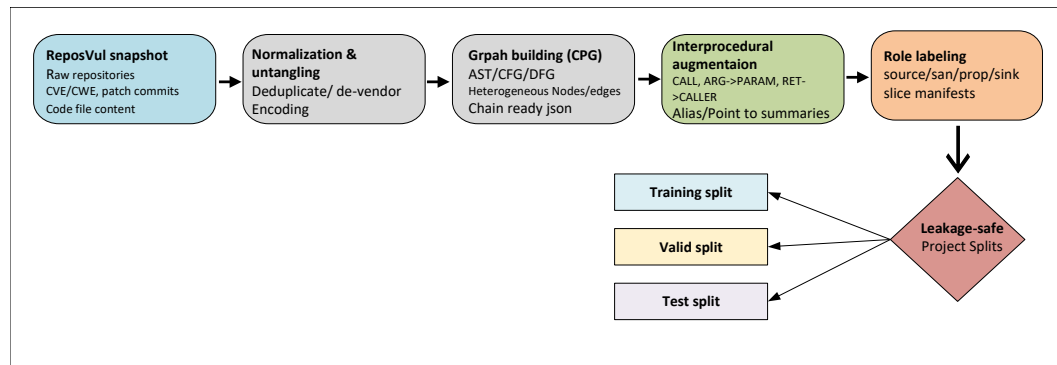


Figure 3.2: ReposVul pipeline: interprocedural CPG construction and leakage-safe data splits.

**Stage A: Repository snapshots and vulnerability metadata.** CVE/CWE metadata is associated with each *ReposVul* entry, and both parent and child commits are retrieved for every patch. Commit identifiers, messages, dates, repository identifiers, and file paths are indexed so that later stages can reconstruct history-aware cross-file context.

**Stage B: Normalization and untangling.** Patches frequently mix vulnerability fixes with unrelated refactorings or vendor code. A corpus untangling rule combining model judgements with static cues is applied to retain only files and hunks that are relevant to the vulnerability fix and to discard unrelated changes.

**Stage C: Static graph construction.** For each retained snapshot, a heterogeneous program graph is built that overlays AST, CFG, and DFG with direct call edges. Pointer and container def-use are approximated by conservative alias relations, identifiers are anonymized, and literals are bucketed to reduce noise from harmless edits [5, 38]. The result is a chain-ready CPG representation for every project version.

**Stage D: Interprocedural augmentation.** The CPG is enriched with explicit interprocedural edges, including CALL, ARG→PARAM, and RET→CALLER/RET→LHS, together with summary nodes and alias summaries that connect callee summaries back to call sites. This step lifts per-file graphs to repository-level structures that expose long-range flows from candidate sources to potential sinks.

**Stage E: Patch-aware differencing and slice materialization.** Line-level diffs are computed for each patch, and four views are synchronised: changed lines, enclosing functions, touched files, and a reachable repository subgraph that can flow to any sink via control or data edges.

**Stage F: Leakage-safe filtering and splitting.** File paths, repository identifiers, and commit metadata are used to drop broken files, superseded patches, and obsolete snapshots. The remaining instances are grouped into training, validation, and test sets such that leakage across splits is minimised.

Table 3.2 summarises the end-to-end data preparation pipeline, outlining the main stages and their key operations and outputs.

Table 3.2: Preparation pipeline summary.

Stage	Key operations and outputs
A	Link CVE/CWE records to parent and child commits; index commits, projects, and file paths.
B	Normalize and untangle patches; keep only fix-relevant files and hunks; deduplicate.
C	Build a heterogeneous CPG per snapshot by overlaying AST, CFG, DFG, and intra-file call edges; anonymize identifiers and bucket literals.
D	Enrich the CPG with interprocedural links (CALL, ARG→PARAM, RET→CALLER/RET→LHS) and alias summaries to expose repository-level flows.
E	Compute line-level diffs; function, file, and reachable-subgraph slices, with rule-based source/flow/stop/sink labels.
F	Drop broken or superseded patches; form leakage-safe train/validation/test splits at the project level.

### 3.1.3 Unified Multigraph, Typing, Features, and Storage

Each repository snapshot is represented as a heterogeneous multigraph based on the CPG abstraction [51]. A single typed node store contains program elements (identifiers, literals, statements, basic blocks, functions), and relation-specific edge sets encode AST, CFG, and DFG structure. Interprocedural edges (CALL, ARG2PARAM, RET2CALL, RET2LHS) are materialised in both directions. Graphs are stored in sharded files with repository, commit, and file-path metadata to support scalable provenance and retrieval.

The representation employs a dual-channel feature encoding. A compact structural feature vector captures node type, degree statistics, SSA hints, literal buckets, and role flags, while a 768-dimensional GraphCodeBERT embedding [13] provides contextual token semantics. These two channels are fused during initialization, but the graph topology and interprocedural edge schema remain invariant across variants. The fixed relation alphabet  $\mathcal{R}$  used during message passing and decoding includes all intra- and interprocedural relations exploited by the chain decoder (AST, CFG, DFG, CALL, ARG2PARAM, RET2CALL, RET2LHS, and alias summaries).

### 3.1.4 Ground Truth

Ground truth is defined at two levels. At the revision level, a repository snapshot is labelled vulnerable if it is linked to a corresponding fixed revision in *ReposVul*, which reduces noise from unrelated edits [45]. At the mechanism level, each positive

example carries at least one executable root-to-sink chain that may span multiple functions and files. Sources represent untrusted inputs, sanitizers restrict or reset tainted state, propagators transfer taint through assignments, calls, and returns, and sinks denote security-critical operations. The full role lexicon and matching patterns are listed in Appendix F.

Chains are seeded from patch diffs, extended along def-use links to identify propagators and candidate sinks, and then wired interprocedurally using ARG  $\rightarrow$  PARAM, RET  $\rightarrow$  CALLER, and RET  $\rightarrow$  LHS bindings over the call graph. Retained examples satisfy basic consistency criteria: every source reaches at least one sink, every sanitizer blocks at least one tainted path, and every propagator lies on a CFG-consistent route. A feasibility pass verifies that neutralising the sink reduces exploitability, and a counterfactual pass strengthens guards or replaces dangerous sinks with safe ones to confirm that reachability changes as expected. Cases that fail these checks are corrected or discarded.

Split-wise statistics for the C/C++ subset used in experiments are reported in the experimental setup chapter (Section 4.1).

### 3.2 Model Architecture and Decoding Pipeline

This section explains how source code and program structure are fused, encoded, and decoded into a single executable chain. All equations and symbol glossaries are consolidated in Appendix B. Refer to Figure 3.3 for a comprehensive single-page overview that this section follows.

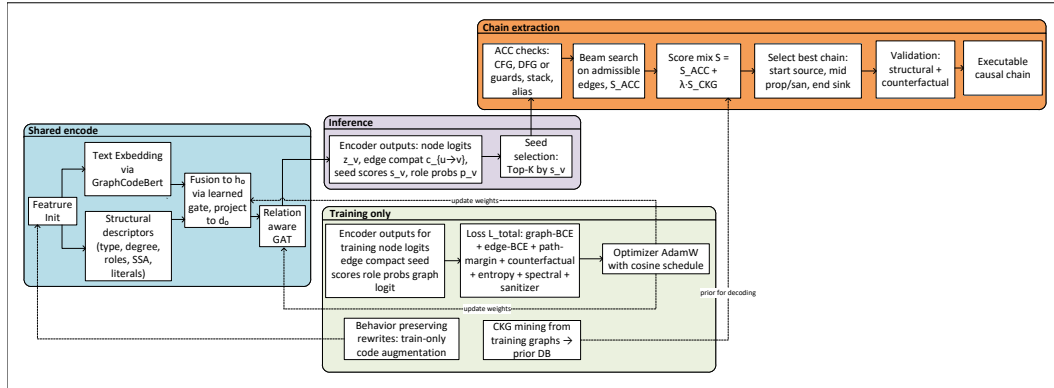


Figure 3.3: Chain-Centric Model Architecture and Inference Flow

Node features are initialized with code spans and structural descriptors, encoded via a relation-aware Graph Attention Network (GAT), and decoded into a single executable chain subject to ACC feasibility. Inference begins with seeds selected by

the model and can incorporate a Causal Knowledge Graph (CKG) prior for ranking. Training optimizes a composite loss, using CKG derived only from the training split. Solid arrows indicate data or score flow, dashed arrows denote training-only operations, and dotted arrows represent weak decoding priors or optimizer feedback to trainable encoder components. GraphCodeBERT provides pretrained contextual embeddings for code tokens, which are fused with structural descriptors before graph encoding.

### 3.2.1 GraphCodeBERT Feature Initialization

Each node aligned to a concrete source span is tokenized with GraphCodeBERT’s byte-pair encoding and encoded in sliding windows (length  $L$ , stride  $S$ ). Window-level contextual token embeddings are aggregated back to the node: if the span contains identifiers, a light edge-aware attention emphasizes tokens carrying def-use signal (Appendix Eqs. B.2–B.3); otherwise a simple average is used (Appendix Eq. B.1). Spans that appear in multiple windows are pooled per window and then averaged, and function-proxy nodes pool their child statements to obtain a coarse function representation.

The textual embedding  $\mathbf{x}_i^{\text{text}}$  and structural descriptor  $\mathbf{x}_i^{\text{struct}}$  (node type, degrees, SSA hints, literal statistics, role flags) are both projected to dimension  $d_0$  and fused via a learned gating mechanism to obtain the initial node state  $\mathbf{h}_i^{(0)}$  (Appendix Eqs. B.4–B.6). When a node and its def-use neighbours fall within a single window, GraphCodeBERT’s data-flow mask is enabled so these pairs can attend directly; the encoder remains frozen to attribute performance gains to the graph module, and Table 3.3 summarizes the resulting feature dimensionalities.

Table 3.3: Node feature dimensions at initialization.

Component	Dim.	Notes
Textual embedding	768	GraphCodeBERT contextual vector [13]
Structural raw	100–200	Type, role, degree, SSA hint, literal buckets
Projected text	$d_0$	$W_t : \mathbb{R}^{768} \rightarrow \mathbb{R}^{d_0}$ (with LN)
Projected structural	$d_0$	$W_s : \mathbb{R}^{d_s} \rightarrow \mathbb{R}^{d_0}$ (with LN)
Fused node init $\mathbf{h}^{(0)}$	$d_0$	Gated combination for the GAT encoder

### 3.2.2 GAT with Causality-Oriented Attention

A relation-aware Graph Attention Network consumes the fused initializations and builds contextual node states while learning how different edge types contribute to aggregation. Per layer, attention and updates follow Appendix Eqs. B.7–B.8. After  $L$  layers the encoder emits a node vulnerability logit  $z_v$ , a seed score  $s_v$ , and relation-gated edge compatibilities  $c_{u \rightarrow v}^{(r)}$  (Appendix Eqs. B.9–B.10). Seeds identify

likely chain starts as the top- $K$  nodes by  $s_v$  (Appendix Eq. B.11); candidate paths are scored by a log-additive mixture of node evidence and edge compatibility (Appendix Eq. B.12), with  $\alpha$  controlling the mix. Training couples node- and edge-level objectives with a path-margin term so that true chains outrank admissible distractor walks.

### 3.2.3 Causal Knowledge Graph (CKG): Mining and Prior

From training chains only, unigram and bigram statistics are mined over the relation alphabet. At decoding time, admissibility is unchanged, but the ranking of admissible expansions receives a small prior bonus  $S_{\text{CKG}}(\pi)$  mixed into the path score with weight  $\lambda$  (Appendix Eq. B.13). This preserves data-driven evidence while gently preferring historically plausible relation patterns.

### 3.2.4 Adaptive Causal Contextualization (ACC)

ACC converts encoder scores into a single executable interprocedural chain by enforcing constant-time feasibility checks and role-shaped preferences. A partial path maintains a taint footprint, a call stack, and accumulated guards; an expansion  $u \xrightarrow{r} v$  is admissible only if all predicates hold (Appendix Eqs. B.14–B.18). The path score is adjusted by start/middle/end role penalties and a sanitizer dominance bonus (Appendix Eqs. B.19–B.23), together with mild length and repetition costs (Appendix Eqs. B.24–B.25), yielding the ACC objective  $S_{\text{ACC}}(\pi)$  and its optional CKG mixture  $S_{\text{ACC}}^*(\pi)$  (Appendix Eqs. B.26–B.27). On CALL, the callee and site are pushed; RET2CALL/RET2LHS pop only with matching sites, preserving well-nested cross-function paths. The decoding complexity with beam width  $B$ , horizon  $H$ , and average admissible out-degree  $\bar{d}$  scales as  $O(B \cdot H \cdot \bar{d})$ . This reflects exploring  $B$  candidates over  $H$  steps, each with  $\bar{d}$  possible expansions, consistent with standard beam search computational cost analyses.

### 3.2.5 Chain Extraction and Validation

A candidate is valid if it starts near a source, contains at least one interior propagator or sanitizer, and ends at a sink (Appendix Eqs. B.28–B.30); when interprocedural links exist in the slice, at least one must appear in the chain. Among admissible paths across beams from the top- $K$  seeds, selection maximizes  $S_{\text{ACC}}(\pi)$  (Appendix Eq. B.31). The final chain is then validated structurally and counterfactually.

### 3.2.6 Training, Optimization, and Hyperparameters

The total loss combines graph-level binary cross-entropy with regularizers for flow and counterfactual consistency, attention entropy, spectral control, and sanitizer alignment. The CKG prior is applied only at decoding time. Models are trained

with AdamW, cosine learning-rate decay, gradient clipping, mixed precision, and weighted sampling for class imbalance. Semantics-preserving augmentations (renaming, benign reordering) are used, and fixed seeds, logged checkpoints, cached features, and archived training curves support reproducibility.

This chapter summarizes a five-step methodology for reconstructing executable root-to-sink vulnerability chains. First, nodes are encoded with GraphCodeBERT and fused with structural features to obtain context-sensitive initial states for all program elements (Appendix Eqs. B.2–B.6). Second, a relation-aware graph attention network propagates information along heterogeneous control-, data-, and call/return edges so that node representations reflect their interprocedural context (Appendix Eqs. B.7–B.8). Third, a constrained decoder guided by a weak causal knowledge graph prior scores and assembles candidate paths, restricting search to syntactically and semantically plausible root-to-sink sequences (Appendix Eq. B.13). Fourth, adaptive causal contextualization (ACC) filters and selects chains that satisfy control-flow, data-flow, and aliasing constraints, ensuring that only semantically valid mechanisms are retained (Appendix Eqs. B.14–B.31). Finally, the entire model is trained at the graph level using binary cross-entropy loss with light regularization (Appendix Eq. B.32), aligning vulnerability predictions with the reconstructed executable chains.

## Chapter 4

# Experimental Results and Analysis

This chapter presents the experimental results for the proposed chain-centric interprocedural vulnerability detector. Under the default ReposVul splits, it quantifies graph-level detection performance, the quality of reconstructed root-propagation-sink chains, causal consistency under targeted interventions, and inference cost and decoding behaviour. All results are obtained under a single, leakage-safe experimental setup and are averaged over five seeds to support reproducibility.

The analysis progresses from an overview of global behaviour to a detailed mechanism-level examination. It first summarises the dataset, model variants, and evaluation metrics used throughout the experiments, then presents detection metrics for both encoder configurations. Subsequent sections examine executable chain feasibility and interprocedural structure, evaluate causal faithfulness via counterfactual edits, and characterise decoding efficiency. The chapter concludes with robustness experiments and a case study that illustrates how an executable root-propagation-sink chain is reconstructed on a real test instance.

### 4.1 Experimental Setup and Configuration

This section summarizes the experimental setup used throughout the evaluation, including datasets and default splits.

**Dataset and splits:** All experiments use the default splits provided with ReposVul [45]. Splits are repository-disjoint, so code from the same repository does not appear in multiple splits. Positive and negative pairs are kept within the same split, and, when reported, chronology is respected so that training data does not depend on future revisions.

Tables 4.1–4.3 report split-wise statistics for the prepared C and C++ subset used in the experiments, including label balance at the file-snapshot level, graph and node-label density after chain-centric conversion, and the prevalence of interprocedural caller/callee structure.



Table 4.1: File-level label counts by split (C/C++ subset).

Split	Records	Non-vuln	Vuln	Pos.%
Train	185,791	180,259	5,532	2.98
Valid	23,224	22,503	721	3.10
Test	23,224	22,554	670	2.88

Table 4.2: Graph instances and node-label density after chain-centric conversion.

Split	Graphs	Vuln nodes	Non-vuln nodes	Pos. ratio
Train	3,438	9,946	25,173,258	$3.95 \times 10^{-4}$
Valid	2,905	1,455	3,970,281	$3.66 \times 10^{-4}$
Test	2,915	1,316	3,973,974	$3.31 \times 10^{-4}$

Table 4.3: Interprocedural connectivity (non-empty caller and callee sets).

Split	Caller%	Callee%	Both%	Caller_chg%	Callee_chg%	Both_chg%
Train	12.59	28.95	8.72	0.46	2.79	0.06
Valid	13.10	29.26	9.21	0.55	2.74	0.07
Test	12.97	29.17	9.03	0.47	2.90	0.09

These statistics indicate that roughly one third of instances have at least one callee, about one eighth have at least one caller, and around one tenth have both, providing the minimal interprocedural structure required to reconstruct vulnerability chains that span function boundaries.

**Inputs and graph relations:** Models operate on chain-prepared interprocedural CPGs that overlay ASTs, CFGs, and DFGs, and include interprocedural relations CALL, ARG→PARAM, RET→CALLER/RET→LHS, together with alias summaries. Table 4.4 summarizes the main graph, encoder, decoding, training, and calibration settings shared across all experiments.

Table 4.4: Experimental Setup Details.

Item	Setting
Dataset	ReposVul, default repository-disjoint splits, chronology respected
Graph relations	DFG, CFG, CALL, ARG2PARAM, RET2CALL, RET2LHS; alias summaries enabled
Encoder	Relation-aware GAT, width $d_0=64$ , $L=3$ layers
Variants	Struct-only; GCBERT+Struct (GraphCodeBERT frozen)
Decoding	Beam with ACC, $K=8$ , $B=24$ , $H=5$ , node/edge mix $\alpha=0.7$
CKG prior	Decoding-only mixture $\lambda=0.2$ ; smoothing $10^{-3}$ ; temperature 1.0
Training	AdamW; LR $2 \times 10^{-3}$ ; WD $10^{-4}$ ; early stop on validation macro-F1
Calibration	Temperature scaling on validation; ECE with 15 bins
Thresholds	$\tau_{F1^*}$ (per variant) and $\tau=0.5$
Seeds	5 per configuration; mean and 95% CIs reported
Environment	PyTorch 2.4.1, CUDA 12.1, RTX 4070 Laptop GPU, AMP on
Embeddings cache	GraphCodeBERT features, FP16, max length 512, stride 384

**Model variants:** Two encoder configurations are compared under the same graph topology and training/decoding regime. The *Struct-only* variant uses compact structural features (types, degrees, SSA hints, literal buckets), while *GCBERT+Struct* augments these with frozen GraphCodeBERT embeddings [13]. Table 4.5 contrasts the resulting graph topology and feature dimensionality on a representative shard.

Table 4.5: Comparison of Struct-only vs. GCBERT+Struct encodings.

Property	Struct-only	GCBERT	Comment
Nodes / in-dim	3141 / 25	3141 / 793	25+768 features in GCBERT
Total edges	14066	14066	Unchanged
Interproc edges (sum)	4818	4818	CALL/ARG2PARAM/RET2* identical
Feature memory (approx.)	$\sim 0.31$ MB	$\sim 9.50$ MB	Text channel dominates

**Decoding and evaluation:** Decoding uses ACC-constrained beam search, allowing only moves that satisfy control-flow reachability, consistent ARG $\rightarrow$ PARAM and RET $\rightarrow$ CALLER/LHS bindings, alias checks, and stack discipline. A causal knowledge graph (CKG) prior mined from training chains adjusts move scores during inference but does not alter admissibility. Training uses class-weighted binary cross-entropy at the graph level with flow and causal regularizers, early stopping on validation macro-F1, and calibrated thresholds derived from validation data. Table 4.6 summarizes the decoding, training, and evaluation configuration.

Table 4.6: Decoding and Training Configuration.

Item	Setting
Beam / horizon / mix	$K=8, B=24, H=5$ , node-edge mix $\alpha=0.7$
Admissibility gates	CFG reachability, ARG→PARAM, RET→CALLER/LHS, alias checks, stack discipline
CKG prior (inference only)	Mixture $\lambda=0.2$ ; smoothing $\epsilon=10^{-3}$ ; temperature $\tau=1.0$ ; top- $K$ trigrams = 500; weights $(\beta_1, \beta_2, \beta_3)=(0.3, 0.6, 0.1)$
Loss	Class-weighted BCE at graph level + flow and causal regularizers
Early stopping	Macro-F1 on validation, patience = 5
Encoder freeze	GraphCodeBERT frozen; ablation unfreezes last 2 blocks at $0.1 \times \text{LR}$
Thresholds	$\tau_{F1^*}$ from validation and fixed $\tau=0.5$ on test
Calibration	Temperature scaling fit on validation; ECE with 15 bins
Seeds & uncertainty	5 seeds; mean and 95% CIs; bootstrap $10^4$ for metrics, binomial CIs for DCR, paired bootstrap with Cliff’s $\delta$

**Experiment Environment and Artifact:** Experiments were run on an NVIDIA RTX 4070 Laptop GPU with CUDA 12.1 and PyTorch 2.4.1 using automatic mixed precision. GraphCodeBERT embeddings were precomputed and cached as FP16 tensors. The artifact archive includes configuration files, predictions, calibration parameters, causal intervention logs, beam expansion analyses, and environment metadata. This setup ensures full reproducibility of the results, with all decoding, calibration, and reporting details preserved.

## 4.2 Evaluation Metrics

This section defines the standard classification, chain-centric, and counterfactual metrics used to assess detection performance, executable chain quality, and causal faithfulness of the proposed approach.

### 4.2.1 Standard Classification Metrics

Classification models are commonly evaluated by thresholding predicted probabilities to produce discrete labels, enabling the construction of a confusion matrix from counts of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN). From these counts, primary metrics such as Precision, Recall, and F1 score are computed:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{F1} = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (4.1)$$

Precision quantifies the accuracy of positive predictions, Recall measures the coverage of actual positives, and F1 balances the two in a single harmonic mean. Additionally, metrics including Accuracy, Macro-/Micro-F1, AUROC, and AUPRC are computed on the raw prediction scores without thresholding. Due to the low prevalence of vulnerable instances, AUPRC and Macro-F1 are emphasized as primary indicators of detection performance. Formal definitions are detailed in Appendix C.2.

#### 4.2.2 Chain-Centric Metrics

For each positive decision, at most one executable chain is returned (or none if feasibility is not met).

**Validity:** The fraction of returned chains that satisfy ACC feasibility checks is reported. These checks include CFG reachability, def-use or guard consistency, interprocedural call-return discipline, and alias coherence. The formal definition is in Appendix C.3.1.

**Structural fidelity:** When a reference chain is available, reported metrics include node and edge coverage, longest common subsequence ratio (LCS), and role-aware coverage for source, sanitizer, propagator, sink. Formulas are detailed in Appendix C.3.2.

**Interprocedurality:** The IPA Rate quantifies the share of predicted chains that traverse CALL/ARG→PARAM/RET edges when such edges exist in the slice (Appendix C.3.3).

#### 4.2.3 Counterfactual Metrics

Causal robustness is evaluated through three targeted interventions: guard strengthening, call unbinding, and sink neutralization. Results are reported by intervention type with confidence intervals; detailed methods are provided in Appendices C.4.1 and C.4.2. The CKG prior is active except during robustness tests, when it is disabled on the edited graph.

**Counterfactual Consistency Score (CCS):**

$$CCS_i = (p_i - p_i^{\text{do}})^2, \quad (4.2)$$

where  $p_i$  and  $p_i^{\text{do}}$  denote the predicted probabilities before and after the targeted intervention on graph  $i$ , respectively, both bounded in  $[0, 1]$ . A low CCS indicates minimal change in prediction, while a high CCS signals a substantial effect of the intervention. Appendix C.4.1 further defines a directional consistency rate assessing whether changes occur in the expected direction per intervention type.

**Causal Feature Attribution Measure (CFAM):**

$$CFAM_i = \frac{\sum_{f \in F_c} A_i(f)}{\sum_{f \in F_c \cup F_s} A_i(f)} \in [0, 1], \quad (4.3)$$

where  $F_c$  and  $F_s$  are sets of on-chain and off-chain features, respectively, and  $A_i(f) \geq 0$  denotes the attribution score of feature  $f$  in graph  $i$ . CFAM quantifies the fraction of total attribution attributed to features on the returned causal chain. Values near 1 indicate strong alignment of attribution with the chain’s features. Normalization details are given in Appendix C.4.2. Auxiliary decoding diagnostics, including Chain Success Rate, admissible expansion ratio, and prior influence rate, are defined in Appendix C.5.

### 4.3 Detection Results

For both encoder variants, this section reports standard graph-level metrics: Accuracy, Precision, Recall, F1, AUROC, and AUPRC. The operating threshold  $\tau_{F1^*}$  is selected on the validation split and applied unchanged to the test set. Probabilities are calibrated via temperature scaling on validation and reused on test, and calibration quality is summarised using Expected Calibration Error (ECE), where lower values indicate better alignment between predicted confidence and empirical accuracy (formal definition and binning scheme in Appendix C.1.2). AUPRC is treated as the primary metric, with F1 and AUROC providing additional context given the class imbalance in the test set (positives  $\approx 2.9\%$ ). All results are averaged over five random seeds.

Table 4.7 reports validation and test performance at  $\tau_{F1^*}$  (mean over five seeds) for both encoder variants, including thresholded Accuracy, Precision, Recall, F1, and the corresponding AUROC and AUPRC values.

Table 4.7: Valid and test metrics for Struct-only and GCBERT+Struct.

Split	Variant	Acc	Prec	Rec	F1	AUROC / AUPRC
Valid	Struct-only	0.954	0.320	0.530	0.400	0.820 / 0.300
Valid	GCBERT+Struct	0.963	0.450	0.660	0.540	0.890 / 0.450
Test	Struct-only	0.953	0.310	0.520	0.390	0.810 / 0.280
Test	GCBERT+Struct	0.965	0.440	0.640	0.520	0.880 / 0.430

Across both valid and test splits, the GCBERT+Struct configuration yields higher F1 and AUPRC than Struct-only (AUPRC +0.15, F1 +0.14 on VALID; AUPRC +0.15, F1 +0.13 on TEST). Accuracy remains high for both variants due to class imbalance, while AUROC indicates stable ranking performance.

Table 4.8 summarises the corresponding operating thresholds and calibration parameters, including the learned temperature  $T$  and ECE values before and after temperature scaling.

Table 4.8: Operating thresholds and calibration parameters.

Split	Variant	$\tau_{F1^*}$	Temp $T$	ECE (before)	ECE (after)
Valid	Struct-only	0.32	1.41	0.079	0.034
Valid	GCBERT+Struct	0.27	1.29	0.061	0.021
Test	Struct-only	0.32	1.41	0.082	0.036
Test	GCBERT+Struct	0.27	1.29	0.064	0.022

Calibration reduces ECE for both variants, and the relative ordering between Struct-only and GCBERT+Struct is consistent with the improvements observed in AUPRC and F1. Precision–Recall and ROC curves exhibit the same pattern. To contextualise these aggregate metrics, Table 4.9 presents confusion matrices on `TEST` at  $\tau_{F1^*}$  for both encoder variants.

Table 4.9: Confusion matrices on `TEST` at  $\tau_{F1^*}$  for Struct-only and GCBERT+Struct.

	Struct-only		GCBERT+Struct	
	Pred. Neg	Pred. Pos	Pred. Neg	Pred. Pos
True Neg	21,779	775	22,008	546
True Pos	322	348	241	429

False positives correspond to safe code incorrectly flagged as vulnerable, while false negatives indicate missed vulnerabilities. In this setting, many false negatives arise in macro-expanded code or callback-driven execution paths, where control flow is implicit and interprocedural structure is more difficult to recover faithfully.

#### 4.4 Executable Chain Quality and Interprocedural Evidence

This section evaluates the quality of reconstructed chains using feasibility, interprocedural usage, structural agreement with ground truth, and chain length/span statistics. All quantities are computed on positive slices and averaged over five seeds. Table 4.10 reports the fraction of predicted chains that pass all ACC feasibility checks (control-flow reachability, call/return consistency, data/guard preservation, alias coherence).

Table 4.10: Feasibility of reconstructed chains (pass of all checks).

Split	Variant	Validity	Notes
Valid	Struct-only	0.762	More CFG violations in long hops
Valid	GCBERT+Struct	0.842	Fewer alias/stack failures
Test	Struct-only	0.741	Errors concentrate at returns
Test	GCBERT+Struct	0.823	Higher pass rate across seeds

GCBERT+Struct improves chain feasibility on both validation and test splits. Table 4.11 summarizes interprocedural usage: the IPA rate (share of chains that use any CALL/ARG→PARAM/RET edges), the fraction that include both call and return edges, and the mean call depth.

Table 4.11: Interprocedural structure in predicted chains.

Split	Variant	IPA rate	Both(call+ret)	Mean call depth
Valid	Struct-only	0.618	0.402	1.27
Valid	GCBERT+Struct	0.708	0.486	1.32
Test	Struct-only	0.603	0.389	1.24
Test	GCBERT+Struct	0.691	0.471	1.30

The enriched encoder yields higher IPA rates, more chains that use both call and return edges, and slightly deeper call stacks.

Table 4.12 reports node and edge coverage, together with role-aware coverage for *source*, *sanitizer*, *propagator*, and *sink* nodes against the reference chains, and Table 4.13 summarises order agreement using the longest common subsequence (LCS) ratio.

Table 4.12: Role and edge agreement with ground truth.

Split	Variant	NodeCov	EdgeCov	RoleCov_src	RoleCov_san	RoleCov_prop	RoleCov_sink
Valid	Struct-only	0.583	0.462	0.781	0.412	0.551	0.692
Valid	GCBERT+Struct	0.671	0.552	0.842	0.521	0.619	0.763
Test	Struct-only	0.571	0.451	0.773	0.398	0.542	0.681
Test	GCBERT+Struct	0.658	0.540	0.834	0.507	0.607	0.752

Table 4.13: Order agreement via LCS ratio.

Split	Variant	LCS ratio	Comment
Valid	Struct-only	0.523	Mismatches at call boundaries
Valid	GCBERT+Struct	0.604	Better call/return placement
Test	Struct-only	0.515	Early sink hops reduce LCS
Test	GCBERT+Struct	0.595	More faithful step order

GCBERT+Struct improves both coverage (especially for sanitizers and propagators) and order agreement. Finally, Table 4.14 characterises chain length and span via hop count, number of files crossed, and the proportion of summary edges.

Table 4.14: Chain length and span: hops, files crossed, and summary-edge share.

Split	Variant	Mean hops	Median	Files crossed	Summary-edge share
Valid	Struct-only	4.70	4	1.57	0.18
Valid	GCBERT+Struct	4.52	4	1.49	0.12
Test	Struct-only	4.66	4	1.55	0.17
Test	GCBERT+Struct	4.48	4	1.47	0.12

Chains remain short and typically span one to two files, with similar hop counts across variants. The enriched encoder slightly reduces reliance on summary edges while increasing feasibility, interprocedural usage, and agreement with ground-truth chains.

## 4.5 Causal Faithfulness

This section evaluates how strongly predictions depend on the reconstructed chains using counterfactual interventions. The *Counterfactual Consistency Score* (CCS) measures the magnitude of probability changes under edits, the *Directional Consistency Rate* (DCR) captures how often changes move in the expected direction, the *Causal Feature Attribution Measure* (CFAM) quantifies how much attribution mass lies on the predicted chain, and the chain invalidation rate records how often no valid chain remains after editing. Formal definitions and variants are given in Appendix C.4.1 and Appendix C.4.2.

Table 4.15 reports prototype CCS and CFAM values averaged over a small subset of graphs, illustrating the scale of probability shifts and on-chain attribution at different training stages.



Table 4.15: Prototype causal metrics (means over  $N=256$  graphs per split).

Split	CCS (mean)	CFAM (mean)	Notes
Train	$1.76 \times 10^{-9}$	0.0047	Early-epoch snapshot
Valid	$8.89 \times 10^{-9}$	0.0232	Default $\tau=0.25$
Test	$7.97 \times 10^{-9}$	0.0233	Same thresholding

Table 4.16 decomposes CCS by intervention type: strengthening guards and unbinding calls (minor edits) versus sink neutralisation (major edit). Minor edits produce small probability shifts, while sink edits cause substantially larger changes, indicating that the endpoint of the chain has strong influence on the decision.

Table 4.16: CCS by edit type (lower is better for minor, higher for sink edits).

Split	Variant	Guard (minor)	Unbind (minor)	Sink (major)
Valid	Struct-only	0.0048	0.0112	0.118
Valid	GCBERT+Struct	0.0039	0.0091	0.134
Test	Struct-only	0.0051	0.0120	0.112
Test	GCBERT+Struct	0.0041	0.0098	0.129

Table 4.17 shows the Directional Consistency Rate (DCR), that is, the proportion of edits where the probability change has the expected sign. DCR is high for all edits, particularly for sink interventions, indicating stable qualitative responses to causal perturbations.

Table 4.17: Directional consistency rate (DCR; higher is better).

Split	Variant	Guard	Unbind	Sink
Valid	Struct-only	0.78	0.81	0.93
Valid	GCBERT+Struct	0.82	0.85	0.95
Test	Struct-only	0.76	0.79	0.91
Test	GCBERT+Struct	0.81	0.83	0.94

Table 4.18 presents CFAM scores overall and decomposed by chain segment (source, sanitizer, propagator, sink), with values in  $[0, 1]$ . Higher values indicate that a larger share of attribution mass lies on nodes that belong to the returned chain.

Table 4.18: CFAM (overall and by chain segment). Values in  $[0, 1]$ .

Split	Variant	CFAM (all)	Src	San	Prop	Sink
Valid	Struct-only	0.42	0.10	0.09	0.13	0.10
Valid	GCBERT+Struct	0.51	0.12	0.12	0.16	0.11
Test	Struct-only	0.40	0.09	0.09	0.12	0.10
Test	GCBERT+Struct	0.49	0.11	0.11	0.15	0.12

Table 4.19 presents chain invalidation rates, defined as the fraction of instances where no admissible causal chain remains after applying an edit. It also reports the average change in chain score following sink neutralization interventions. The results show that minor edits invalidate approximately half of the chains, indicating moderate disruption of causal explanations. In contrast, sink neutralization almost completely invalidates the chains and results in substantial decreases in chain scores, reflecting the significant impact of this intervention on causal model structure and confidence. This highlights the differing robustness of causal chains to various types of graph edits and underscores sink neutralization as a more drastic perturbation within causal robustness evaluation.

Table 4.19: Chain invalidation rates and score deltas under interventions.

Split	Variant	Inv. Guard	Inv. Unbind	Inv. Sink	$\Delta$ Score (sink)
Valid	Struct-only	0.46	0.54	0.91	-1.27
Valid	GCBERT+Struct	0.55	0.61	0.94	-1.41
Test	Struct-only	0.44	0.52	0.89	-1.21
Test	GCBERT+Struct	0.53	0.60	0.93	-1.36

Taken together, these results demonstrate that the model’s decisions are strongly linked to the reconstructed causal chains. Minor edits along the chains induce small, directionally consistent shifts in predicted probabilities, reflecting stable and interpretable causal influence. In contrast, major edits that neutralize the sink disrupt this causal structure, invalidating most chains and producing substantial probability decreases. This indicates the model relies heavily on these causal chains for its decision-making, with robustness and sensitivity aligned to the nature and location of the intervention within the chain.

## 4.6 Efficiency and Decoding Dynamics

This section reports the computational cost of the proposed pipeline and the effect of decoding constraints on search behaviour. Table 4.20 summarises per-graph inference latency for the main components. Most of the runtime is spent in the

encoder forward pass, while ACC checks and the CKG prior incur only a small additional overhead.

Table 4.20: Inference latency per graph.

Component	Struct-only	GCBERT+Struct	Notes
Encoder forward (ms)	18.6 (31.9)	31.7 (52.5)	Relation-aware GAT; LM frozen
ACC decoding (ms)	2.8 (4.6)	3.3 (5.3)	Admissibility gates
CKG prior mix (ms)	0.3 (0.5)	0.4 (0.6)	Lightweight lookups
<b>Total (ms)</b>	<b>21.8 (36.8)</b>	<b>35.6 (58.4)</b>	End-to-end latency

Structural and ACC checks prune most candidate moves during decoding: CFG reachability alone filters out around 54–63% of expansions (Appendix D.2, Table D.1). Increasing the beam from (4, 12, 3) to (8, 24, 5) improves chain validity, interprocedural usage, and LCS (by 4.2, 13.1, and 6.3 points respectively) at an additional 6.9 ms per graph. Given that encoder computation dominates latency, the chosen beam setting provides a practical trade-off between chain quality and runtime.

## 4.7 Robustness and Generalization

Robustness is assessed under three conditions: cross-repository generalization, forward-in-time evaluation, and invariance to benign refactorings, using the default ReposVul splits and five random seeds. Across unseen projects and later revisions, graph-level detection metrics change only slightly, while chain validity and interprocedural usage remain high and track classification trends within overlapping confidence intervals. Under refactorings (identifier renaming, inert code insertion, and mild intra-block reorderings), predictions and chains are largely stable; only aggressive statement reordering degrades performance noticeably. Detailed numerical deltas and confidence intervals are reported in Appendix D.

## 4.8 Test Case

This case study illustrates how the model behaves in practice by describing one decoded chain from the two-file program `main.c` and `lib/shell.c`, using the same trained model and ACC-constrained decoding with the default beam settings as in the main experiments. The full textual trace and exact source files are provided in Appendix G.10.

In this example, untrusted input enters the program via the call `fgets(buf, sizeof(buf), stdin)` in `main.c`. This call is the *root* of the chain, introducing external data into the program state. The chain then follows local processing of `buf`:

computing `n = strlen(buf)` and conditionally trimming trailing newline or carriage return characters. These steps provide lightweight sanitization and normalization, but the data in `buf` remains marked as tainted, meaning it is considered untrusted and potentially harmful throughout the program’s execution.

The cleaned string is then passed to the execution stage, with the chain crossing the file boundary via a call to `stage_execute` and continuing in `lib/shell.c`. ACC permits this transition only if a valid call edge exists from `main.c` to `lib/shell.c`. Within `stage_execute`, the chain includes calculating the copy length `n = sizeof(exec_local) - 1`, writing the null terminator `exec_local[n] = '\0'`, and preparing the fixed-size buffer `exec_local`. These steps show how the tainted command string is copied into a local buffer that will be executed later.

Finally, the chain terminates at the call `system(exec_local)`, which is labeled as the *sink*. This is a security-sensitive operation, since it executes a shell command derived from external input. Every step along the chain corresponds to a concrete statement in the source code, and edges represent either control flow, data flow, or interprocedural links (calls and returns). ACC enforces that each hop is reachable on the control-flow graph, respects the data and guard dependencies, and maintains call/return discipline and alias constraints. As a result, the final path forms an executable, interprocedural explanation of how attacker-controlled input flows from `fgets` in `main.c` all the way to `system(exec_local)` in `lib/shell.c`.

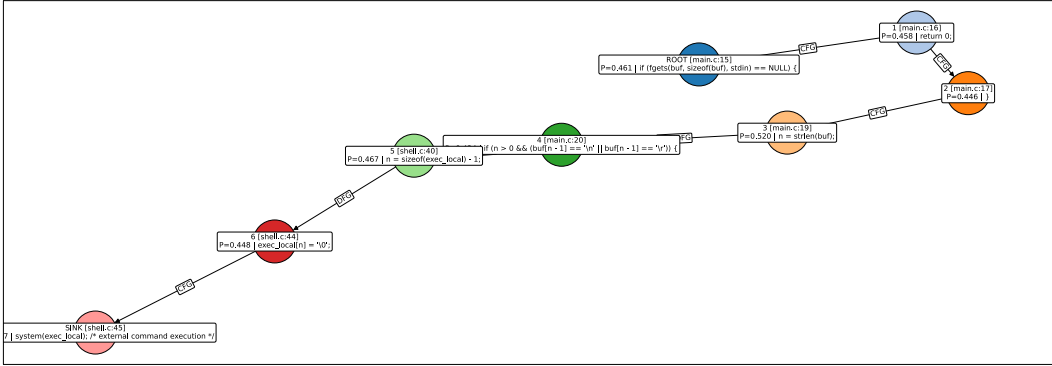


Figure 4.1: Executable interprocedural chain from root to sink across `main.c` and `lib/shell.c`

The results in this chapter demonstrate competitive detection performance with strong AUPRC and F1 scores. Reconstructed chains show high feasibility and consistent interprocedural use concerning semantic roles and sequence order. Counterfactual analysis confirms causal fidelity. Inference is practical, dominated by encoder latency with minimal decoding overhead. Performance remains stable under project-disjoint splits and rigorous conditions. Complete diagnostics and reproduction artifacts are provided in Appendix G.

## Chapter 5

# Conclusion and Future Work

### 5.1 Conclusion

This thesis systematically investigated chain-centric, causality-driven approach for software vulnerability detection with a targeted emphasis on interprocedural behavior. The principal objective was to elucidate how explicit root–propagation–sink causal mechanisms can be reconstructed from real-world source code and leveraged to yield vulnerability predictions that are reliable, interpretable, and faithful to program execution semantics. By reframing vulnerability detection beyond conventional correlational classification into the identification, constraint, and semantic validation of executable causal chains, this research advances a causal paradigm that captures the dynamic propagation of vulnerabilities across procedural and modular boundaries.

The thesis introduced a chain-centric interprocedural program representation founded on an augmented CPG that unifies AST, CFG, and DFG enriched with interprocedural semantics. This augmentation includes explicit call edges, argument–parameter mappings, return–caller bindings, and aliasing summaries, thereby making long-range vulnerability propagation across functions, files, and modules explicitly tractable. Within this enriched representation, vulnerabilities are conceptualized as root–propagation–sink chains annotated with fine-grained semantic roles—sources, sanitizers, propagators, and sinks. Building upon this, a context-aware causal pipeline was developed that integrates structure-aware language-model embeddings with a relation-aware graph encoder employing Adaptive Causal Contextualization (ACC), alongside a constrained beam decoder steered by a Causal Knowledge Graph (CKG) prior. This integrated system is architected to produce a single executable causal chain per positive classification, eschewing traditional unstructured saliency scores to deliver interpretable and semantically faithful vulnerability explanations.

The experimental evaluation on the ReposVul dataset yielded several significant findings. At the graph level, the proposed model demonstrated robust detection

performance under repository-disjoint splits, with the GCBERT+Struct encoder consistently enhancing AUPRC and F1 metrics relative to a compact structural baseline while sustaining high AUROC values. Chain-level analyses revealed that the reconstructed vulnerability propagation paths are predominantly executable and interprocedural, evidenced by high validity rates across validation and test datasets, an increase in interprocedural usage (IPA), and improved coverage of vulnerability roles—sources, sanitizers, propagators, and sinks. Furthermore, alignment with reference chains, quantified by longest common subsequence measures, improved substantially. Targeted causal diagnostics validated that model predictions rely on the reconstructed causal mechanisms: semantic perturbations such as guard strengthening and call unbinding induced subtle, directionally consistent changes in prediction confidence, whereas sink neutralization caused pronounced drops and pervasive chain invalidation. The Causal Feature Attribution Measure further concentrated attribution mass within the reconstructed chains rather than extraneous context. Finally, efficiency evaluations confirmed that the inclusion of ACC checks and the CKG prior imposed only modest computational overhead relative to the encoder, supporting practical end-to-end latency for real-world application.

This thesis contributes a comprehensive and unified methodological approach by conceptualizing software vulnerabilities as explicit root–propagation–sink causal chains, thereby harmonizing representation, learning, and evaluation processes within this causal paradigm. Central to this approach is the integration of an interprocedural, augmented CPG enriched with ACC constraints, CKG-guided decoding, and chain-centric counterfactual evaluation metrics, all implemented in a reproducible setup utilizing default ReposVul splits. This approach effectively addresses the core research questions by demonstrating that the augmented CPG supports precise executable interprocedural chain reconstruction and that a context-aware causal reasoning pipeline leveraging this representation can produce vulnerability predictions that are both highly accurate and causally grounded in explicit propagation chains. Empirical evaluations validate these claims, showing significant improvements in detection performance, chain feasibility, interprocedural coverage, and causal explanatory power without compromising computational efficiency.

Taken together, these contributions establish a foundational approach that links causal, interprocedural reasoning with advanced learning based vulnerability detection. The proposed framework not only advances theoretical understanding, but also offers practical mechanisms for designing more robust, interpretable, and deployable tools suited to complex, real world software ecosystems. By moving vulnerability detection beyond opaque classification toward transparent, causally grounded analysis, this work enables transparent, interpretable and more actionable insights for developers and security practitioners seeking to secure modern software infrastructure.

## 5.2 Future Work

Although the proposed approach advances causality-aware vulnerability detection, several extensions and refinements remain open. The natural progression for future work is to extend the current approach from providing single-chain vulnerability explanations to supporting multiple distinct causal chains per vulnerability instance. Real-world vulnerabilities commonly exhibit multifaceted exploit behaviors, including multiple independent or variant exploit paths that arise from differing call sequences or environment conditions. By enabling the generation and presentation of a concise set of diverse, non-overlapping chains, analysts would gain a richer, more comprehensive understanding of risk exposure without overwhelming them with overly complex or redundant information. This multi-chain explanatory capability would enhance practical vulnerability assessment and remediation by illuminating varied exploitation scenarios within a manageable interpretative framework.

A key direction involves enhancing the priors and semantics that guide the decoding process. Currently, the Causal Knowledge Graph (CKG) serves as a fixed prior extracted once and used statically during inference. Adapting this component to be learnable and dynamic potentially through joint training with the main detection model or fine-tuning could significantly improve its alignment with local coding styles and domain specific conventions. Moreover, refining the interprocedural program representation by incorporating more precise aliasing information, better modeling of dynamic dispatch, and framework-aware summaries would make reconstructed vulnerability chains more representative of actual runtime behavior, enhancing their semantic fidelity and practical utility.

Incorporating lightweight dynamic evidence presents another promising avenue. Purely static analyses may identify vulnerability propagation paths that are feasible in theory but unlikely to occur in realistic executions. By integrating runtime information such as code coverage data, symbolic execution traces, or fuzzing feedback, the system could prioritize and filter generated chains, assign probabilistic scores to chain segments, and verify the executability of reported paths under realistic conditions. This integration would bridge the gap between static interprocedural analyses and observed program behavior, thereby increasing confidence in the detection results.

Expanding the current approach beyond its primary focus on C/C++ code is essential to assess its broader applicability. Extending to additional programming languages, paradigms, and vulnerability classes including concurrency defects, logic errors, and vulnerabilities arising in multi language stacks would probe the generality of the chain-centric causal reasoning philosophy. Concurrently, integrating the pipeline within developer tool chains and conducting empirical user studies could yield invaluable insights into its practical impact, measuring factors such as triage efficiency, fix accuracy, and developer trust in automated vulnerability explanations.

Finally, there is substantial methodological work to be done around evaluation.

The causal quality and faithfulness criteria developed in this thesis mark an initial step toward standardizing chain validation metrics. Future efforts could focus on refining these evaluation measures, benchmarking them across diverse models, and linking them to human-centered outcomes. The establishment of widely accepted benchmarks that include explicit chain annotations, standardized data splits, and agreed-upon evaluation methodology would empower the research community to perform more meaningful comparisons and accelerate progress towards vulnerability detection tools that combine high efficacy with true interpretability.



# Bibliography

- [1] M. ALLAMANIS, E. T. BARR, P. DEVANBU, AND C. SUTTON, *A survey of machine learning for big code and naturalness*, ACM Computing Surveys (CSUR), 51 (2018), pp. 1–37.
- [2] L. AOUAD, *Causal factors analysis of vulnerability exploitation*, 2023. IriusRisk Blog.
- [3] S. CAO, X. SUN, X. WU, D. LO, L. BO, B. LI, AND W. LIU, *Coca: Improving and explaining graph neural network-based vulnerability detection systems*, in 2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE), 2024, pp. 1–13.
- [4] S. CAO, X. SUN, X. WU, D. LO, L. BO, B. LI, X. LIU, X. LIN, AND W. LIU, *Snopy: Bridging sample denoising with causal graph learning for effective vulnerability detection*, in Proceedings - 2024 39th ACM/IEEE International Conference on Automated Software Engineering, ASE 2024, 2024, pp. 606–618.
- [5] S. CHAKRABORTY, R. KRISHNA, Y. DING, AND B. RAY, *Deep learning based vulnerability detection: Are we there yet?*, IEEE Transactions on Software Engineering, 48 (2022), pp. 3280–3296.
- [6] X. CHENG, G. ZHANG, H. WANG, AND Y. SUI, *Path-sensitive code embedding via contrastive learning for software vulnerability detection*, in Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, 2022, pp. 519–531.
- [7] Z. CHU, Y. WAN, Q. LI, Y. WU, H. ZHANG, Y. SUI, G. XU, AND H. JIN, *Graph neural networks for vulnerability detection: A counterfactual explanation*, in Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis, 2024, pp. 389–401.
- [8] K. FILUS AND J. DOMAŃSKA, *Software vulnerabilities in tensorflow-based deep learning applications*, Computers and Security, 124 (2023), p. 102948.
- [9] M. FREITAG AND Y. AL-ONAIZAN, *Beam search strategies for neural machine translation*, arXiv preprint arXiv:1702.01806, (2017).

- [10] M. FU AND C. TANTITHAMTHAVORN, *Linevul: A transformer-based line-level vulnerability prediction*, in Proceedings of the 19th International Conference on Mining Software Repositories, 2022, pp. 608–620.
- [11] A. H. GALIB AND B. M. MAINUL HOSSAIN, *A systematic review on hybrid analysis using machine learning for android malware detection*, in 2019 2nd International Conference on Innovation in Engineering and Technology (ICIET), 2019, pp. 1–6.
- [12] F. GANZ, L. FISCHER, M. KELLER, F. BECK, Y. ACAR, AND M. BACKES, *Software defect localization using explainable deep learning*, in Proceedings of the 17th ACM Workshop on Artificial Intelligence and Security, ACM, 2024.
- [13] D. GUO, S. REN, S. LU, Z. FENG, D. TANG, S. LIU, L. ZHOU, N. DUAN, A. SVYATKOVSKIY, S. FU, ET AL., *Graphcodebert: Pre-training code representations with data flow*. *arxiv* 2020, arXiv preprint arXiv:2009.08366, (2021).
- [14] W. GUO, D. MU, J. XU, P. SU, G. WANG, AND X. XING, *Lemna: Explaining deep learning based security applications*, in Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS '18, New York, NY, USA, 2018, Association for Computing Machinery, p. 364–379.
- [15] D. HIN, A. KAN, H. CHEN, AND M. A. BABAR, *Linevd: Statement-level vulnerability detection using graph neural networks*, in Proceedings of the 19th international conference on mining software repositories, 2022, pp. 596–607.
- [16] N. T. ISLAM, G. D. L. T. PARRA, D. MANUAL, M. JADLIWALA, AND P. NAJAFIRAD, *Causative insights into open source software security using large language code embeddings and semantic vulnerability graph*, arXiv preprint arXiv:2401.07035, (2024).
- [17] A. KAUR AND R. NAYYAR, *A comparative study of static code analysis tools for vulnerability detection in c/c++ and java source code*, *Procedia Computer Science*, 171 (2020), pp. 2023–2029.
- [18] H. KUANG, J. ZHANG, F. YANG, L. ZHANG, Z. HUANG, AND L. YANG, *Vulcausal: Robust vulnerability detection using neural network models from a causal perspective*, in International Conference on Knowledge Science, Engineering and Management, Springer, 2024, pp. 41–56.
- [19] F. LABRÈCHE AND S.-O. PAQUETTE, *Threat class predictor: An explainable framework for predicting vulnerability threat using topic and trend modeling.*, in CAMLIS, 2022, pp. 113–124.
- [20] D. LI, *An XAI-based Framework for Software Vulnerability Contributing Factors Assessment*, PhD thesis, Concordia University, 2023.

- [21] D. LI, Y. LIU, AND J. HUANG, *Assessment of software vulnerability contributing factors by model-agnostic explainable ai*, Machine Learning and Knowledge Extraction, 6 (2024), pp. 1087–1113.
- [22] L. LI, S. H. DING, Y. TIAN, B. C. FUNG, P. CHARLAND, W. OU, L. SONG, AND C. CHEN, *Vulanalyzer: Explainable binary vulnerability detection with multi-task learning and attentional graph convolution*, ACM Transactions on Privacy and Security, 26 (2023), pp. 1–25.
- [23] Z. LI, S. DUTTA, AND M. NAIK, *Iris: Llm-assisted static analysis for detecting security vulnerabilities*, arXiv preprint arXiv:2405.17238, (2024). Proposes an LLM-augmented static analysis that infers taint specifications and outperforms CodeQL on CWE-Bench-Java.
- [24] Z. LI, D. ZOU, J. TANG, Z. ZHANG, M. SUN, AND H. JIN, *A comparative study of deep learning-based vulnerability detection system*, IEEE Access, 7 (2019), pp. 103184–103197.
- [25] G. LIN, S. WEN, Q.-L. HAN, J. ZHANG, AND Y. XIANG, *Software vulnerability detection using deep neural networks: A survey*, Proceedings of the IEEE, 108 (2020), pp. 1825–1848.
- [26] S. LIU, G. LIN, L. QU, J. ZHANG, O. DE VEL, P. MONTAGUE, AND Y. XIANG, *Cd-vuld: Cross-domain vulnerability discovery based on deep domain adaptation*, IEEE Transactions on Dependable and Secure Computing, 19 (2022), pp. 438–451.
- [27] A. LUCIC, M. A. TER HOEVE, G. TOLOMEI, M. DE RIJKE, AND F. SILVESTRI, *Cf-gnnexplainer: Counterfactual explanations for graph neural networks*, in International Conference on Artificial Intelligence and Statistics, PMLR, 2022, pp. 4499–4511.
- [28] S. MANDAL, A. A. SEKH, AND S. K. NASKAR, *Solving arithmetic word problems: A deep learning based approach*, Journal of Intelligent Fuzzy Systems, 39 (2020), pp. 2521–2531.
- [29] MARCHETTO AND ALESSANDRO, *Can explainability and deep-learning be used for localizing vulnerabilities in source code?*, in Proceedings of the 5th ACM/IEEE International Conference on Automation of Software Test (AST 2024), 2024, pp. 110–119.
- [30] A. F. T. MARTINS, G. CORREIA, V. NICULAE, G. POESIA, O. POLOZOV, V. LE, A. TIWARI, G. SOARES, C. MEEK, AND S. GULWANI, *Structured Decoding for Constrained Generation*, in International Conference on Learning Representations (ICLR) 2022, 2022.
- [31] B. MOSOLYGÓ, N. VÁNDOR, G. ANTAL, P. HEGEDŰS, AND R. FERENC, *Towards a prototype based explainable javascript vulnerability prediction model*, in 2021 International conference on code quality (ICCQ), IEEE, 2021, pp. 15–25.

- [32] H. Q. NGUYEN, T. HOANG, H. K. DAM, AND A. GHOSE, *Graph-based explainable vulnerability prediction*, Information and Software Technology, 176 (2024), p. 107386.
- [33] Y. NONG, Y. OU, M. PRADEL, F. CHEN, AND H. CAI, *Vulgen: Realistic vulnerability generation via pattern mining and deep learning*, in 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE), IEEE, 2023, pp. 2527–2539.
- [34] H. PEARCE, B. AHMAD, B. TAN, B. DOLAN-GAVITT, AND R. KARRI, *Asleep at the keyboard? assessing the security of github copilot’s code contributions*, Communications of the ACM, 68 (2025), pp. 96–105.
- [35] M. M. RAHMAN, I. CEKA, C. MAO, S. CHAKRABORTY, B. RAY, AND W. LE, *Towards causal deep learning for vulnerability detection*, in Proceedings of the IEEE/ACM 46th international conference on software engineering, 2024, pp. 1–11.
- [36] S. RUSSELL AND P. NORVIG, *Artificial Intelligence: A Modern Approach*, Pearson, 4th ed., 2020.
- [37] A. SEJFIA, S. DAS, S. SHAFIQ, AND N. MEDVIDOVIĆ, *Toward improved deep learning-based vulnerability detection*, in Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, 2024, pp. 1–12.
- [38] B. STEENHOEK, M. M. RAHMAN, R. JILES, AND W. LE, *An empirical study of deep learning models for vulnerability detection*, in 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE), IEEE, 2023, pp. 2237–2248.
- [39] S. SUNEJA, Y. ZHENG, Y. ZHUANG, J. A. LAREDO, AND A. MORARI, *Probing model signal-awareness via prediction-preserving input minimization*, in Proceedings of the 29th ACM joint meeting on european software engineering conference and symposium on the foundations of software engineering, 2021, pp. 945–955.
- [40] J. TAN, S. XU, Y. GE, Y. LI, X. CHEN, AND Y. ZHANG, *Counterfactual explainable recommendation*, in Proceedings of the 30th ACM International Conference on Information & Knowledge Management, 2021, pp. 1784–1793.
- [41] S. TAVISS, S. H. DING, M. ZULKERNINE, P. CHARLAND, AND S. ACHARYA, *Asm2seq: Explainable assembly code functional summary generation for reverse engineering and vulnerability analysis*, Digital Threats: Research and Practice, 5 (2024), pp. 1–25.
- [42] M. TUFANO, Y. LI, Z. TU, AND B. RAY, *Adaptive, reinforcement-learning-guided symbolic execution for smart contracts*, in Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering, 2022, pp. 1–13.
- [43] US GOVERNMENT, *Federal register*, 2023. Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence.

- [44] P. VELIČKOVIĆ, G. CUCURULL, A. CASANOVA, A. ROMERO, P. LIÒ, AND Y. BENGIO, *Graph attention networks*, in International Conference on Learning Representations (ICLR), 2018.
- [45] X. WANG, R. HU, C. GAO, X.-C. WEN, Y. CHEN, AND Q. LIAO, *Reposvul: A repository-level high-quality vulnerability dataset*, in Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings, 2024, pp. 472–483.
- [46] X. WANG, H. JI, C. SHI, B. WANG, P. CUI, P. YU, AND Y. YE, *Heterogeneous graph attention network*, arXiv preprint arXiv:1903.07293, (2019).
- [47] S. WOO, E. CHOI, H. LEE, AND H. OH, *{VISCAN}: Discovering 1-day vulnerabilities in reused {C/C++} open-source software components using code classification techniques*, in 32nd USENIX Security Symposium (USENIX Security 23), 2023, pp. 6541–6556.
- [48] Y. XU, Y. ZHANG, C. WANG, S. LI, A. ZHANG, AND S. DENG, *ML-guided fuzzing: A systematic review*, in 2023 IEEE 34th International Symposium on Software Reliability Engineering (ISSRE), 2023, pp. 577–588.
- [49] C. YAGEMANN, S. P. CHUNG, B. SALTAFORMAGGIO, AND W. LEE, *Automated bug hunting with data-driven symbolic root cause analysis*, in Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, 2021, pp. 320–336.
- [50] C. YAGEMANN, M. PRUETT, S. P. CHUNG, K. BITTICK, B. SALTAFORMAGGIO, AND W. LEE, *{ARCUS}: symbolic root cause analysis of exploits in production systems*, in 30th USENIX Security Symposium (USENIX Security 21), 2021, pp. 1989–2006.
- [51] F. YAMAGUCHI, N. GOLDE, D. ARP, AND K. RIECK, *Modeling and discovering vulnerabilities with code property graphs*, in 2014 IEEE Symposium on Security and Privacy, 2014, pp. 590–604.
- [52] Z. YANG, J. SHI, J. HE, AND D. LO, *Natural attack for pre-trained models of code*, in Proceedings of the 44th International Conference on Software Engineering, 2022, pp. 1482–1493.
- [53] Z. YING, D. BOURGEOIS, J. YOU, M. ZITNIK, AND J. LESKOVEC, *Gnnexplainer: Generating explanations for graph neural networks*, in Advances in Neural Information Processing Systems, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds., vol. 32, Curran Associates, Inc., 2019.
- [54] D. YU, Q. LI, X. WANG, Q. LI, AND G. XU, *Counterfactual explainable conversational recommendation*, IEEE Transactions on Knowledge and Data Engineering, 36 (2023), pp. 2388–2400.

- [55] E. ZELIKMAN, E. LORCH, L. MACKEY, AND A. T. KALAI, *Self-taught optimizer (stop): Recursively self-improving code generation*, in First Conference on Language Modeling, 2024.
- [56] P. ZENG, G. LIN, L. PAN, Y. TAI, AND J. ZHANG, *Software vulnerability analysis and discovery using deep learning techniques: A survey*, IEEE Access, 8 (2020), pp. 197158–197172.
- [57] B. ZHOU, J. LI, X. LI, B. HUA, AND J. BAO, *Leveraging on causal knowledge for enhancing the root cause analysis of equipment spot inspection failures*, Advanced Engineering Informatics, 54 (2022), p. 101799.

## Appendix A

# Algorithms and Pseudocode

This appendix records the core procedures that were referenced in the methodology: (1) relation-aware GAT, (2) the ACC constrained beam search for chain decoding, (3) chain Selection and structural validation, and (4) counterfactual intervention test.

---

**Algorithm 1:** Relation-aware GAT Message Passing (One Layer)

---

**Input** : Node states  $\{\mathbf{h}_v^{(\ell)}\}$ ; typed neighbor sets  $\mathcal{N}_r(v)$  for each relation  $r \in \mathcal{R}$ .

**Output**: Updated node states  $\{\mathbf{h}_v^{(\ell+1)}\}$ .

```
foreach node  $v \in V$  do
     $\mathbf{m}_{\text{self}} \leftarrow \mathbf{W}_{\text{self}} \mathbf{h}_v^{(\ell)}$ ;
     $\mathbf{m}_{\text{sum}} \leftarrow \mathbf{0}$ ;
    foreach relation  $r \in \mathcal{R}$  do
        foreach neighbor  $u \in \mathcal{N}_r(v)$  do
             $\mathbf{z}_{uv}^{(r)} \leftarrow [\mathbf{W}_r \mathbf{h}_u^{(\ell)} \parallel \mathbf{W}_0 \mathbf{h}_v^{(\ell)}]$ ;
             $e_{uv}^{(r)} \leftarrow \text{LeakyReLU}(\mathbf{a}_r^\top \mathbf{z}_{uv}^{(r)})$ ;
        foreach neighbor  $u \in \mathcal{N}_r(v)$  do
             $\alpha_{uv}^{(r)} \leftarrow \frac{\exp(e_{uv}^{(r)})}{\sum_{u' \in \mathcal{N}_r(v)} \exp(e_{uv'}^{(r)})}$ ;
             $\mathbf{m}_{\text{sum}} \leftarrow \mathbf{m}_{\text{sum}} + \alpha_{uv}^{(r)} \mathbf{W}_r \mathbf{h}_u^{(\ell)}$ ;
     $\mathbf{h}_v^{(\ell+1)} \leftarrow \text{ELU}(\mathbf{m}_{\text{self}} + \mathbf{m}_{\text{sum}})$ ;
return  $\{\mathbf{h}_v^{(\ell+1)}\}$ ;
```

---

**Algorithm 2:** ACC-constrained Beam Search for Chain Decoding

**Input** : Graph  $G$ ; role probabilities  $p_v^{\text{src/prop/san/sink}}$ ;  
 node scores  $s_v$ , chain scores  $z_v$ , edge compatibilities  $c_{u \rightarrow v}^{(r)}$ ;  
 beam parameters  $(K, B, H)$ ; predicates  $\text{cfg\_ok}$ ,  $\text{dfg\_ok}$ ,  $\text{ipa\_ok}$ ,  
 $\text{alias\_ok}$ .

**Output:** Best admissible path  $\pi^*$  (root $\rightarrow$ propagation $\rightarrow$ sink) or None.

**Path state.** Each partial path  $\pi$  maintains: tip node  $v_t$ , taint footprint  $T(\pi)$ ,  
 call stack  $C(\pi)$ , accumulated guards  $G(\pi)$ , and score  $S_{\text{ACC}}(\pi)$ .

**Initialization:**

1. Select seed nodes  $S_K \leftarrow \text{TopK}_v(s_v)$ .
2. Initialize the beam as

$$\text{beams} \leftarrow \left\{ \pi = [v_0] \mid v_0 \in S_K, S_{\text{ACC}}(\pi) = \log \sigma(s_{v_0}) \right\},$$

with  $T(\pi), C(\pi), G(\pi)$  derived from  $v_0$ .

**Expansion:** for  $t = 1$  to  $H$  do

```

pool  $\leftarrow \emptyset$ ;
foreach partial path  $\pi \in \text{beams}$  do
   $u \leftarrow \text{tip}(\pi)$ ;
  foreach typed edge  $(u \xrightarrow{r} v)$  in  $G$  do
    if  $\text{cfg\_ok}(u \rightarrow v)$  is false or  $\text{dfg\_ok}(u \rightarrow v, T(\pi))$  is false or
        $\text{ipa\_ok}(u \rightarrow v, C(\pi))$  is false or  $\text{alias\_ok}(u \rightarrow v, T(\pi))$  is false then
      continue to next neighbor;
    // Create successor path and update state
     $\pi' \leftarrow \pi \parallel v$ ; // append v
    update  $T(\pi'), C(\pi'), G(\pi')$  from  $u \rightarrow v$ ;
     $\Delta \leftarrow \alpha \log \sigma(z_v) + (1 - \alpha) c_{u \rightarrow v}^{(r)}$ ;
     $S_{\text{ACC}}(\pi') \leftarrow S_{\text{ACC}}(\pi) + \Delta$ ;
     $\text{role\_penalty}(\pi') + \text{san\_bonus}(\pi')$ ;
     $\text{len\_penalty}(\pi') - \text{rep\_penalty}(\pi')$ ;
    add  $\pi'$  to pool;
  // Keep top- $B$  candidates by ACC score
   $\text{beams} \leftarrow \text{TopB}_{\pi \in \text{pool}} S_{\text{ACC}}(\pi)$ ;
  if some  $\pi \in \text{beams}$  ends at a sink and no successor can score higher then
    break;
```

**Selection:** Choose  $\pi^*$  as the highest-scoring valid path in beams, breaking ties in favour of: (i) fewer summary edges, (ii) inclusion of a sanitizer, and (iii) fewer distinct files;

**return**  $\pi^*$ ;



---

**Algorithm 3:** Chain Selection and Structural Validation

---

**Input** : Candidate paths  $\mathcal{P}$  from ACC decoding; thresholds  $\tau_{src}, \tau_{mid}, \tau_{sink}$ .**Output**: Validated chain  $\hat{\pi}$  or None.**Step 1: Role-shaped filtering.****foreach** path  $\pi = (v_0, \dots, v_T) \in \mathcal{P}$  **do**

```

    if  $p_{src}(v_0) < \tau_{src}$  then
        | discard  $\pi$  and continue;
    if  $\max_{t \in \{1, \dots, T-1\}} (p_{prop}(v_t), p_{san}(v_t)) < \tau_{mid}$  then
        | discard  $\pi$  and continue;
    if  $p_{sink}(v_T) < \tau_{sink}$  then
        | discard  $\pi$  and continue;
    | keep  $\pi$  in the filtered set  $\mathcal{P}'$ ;

```

**Step 2: Interprocedural sufficiency.****if** the slice contains any CALL/ARG2PARAM/RET2\* edges **then**

```

    | require that at least one path in  $\mathcal{P}'$  uses an interprocedural edge; discard
    | purely intra-procedural paths if interprocedural ones exist;

```

**Step 3: Maximization.**Let  $\hat{\pi}$  be the path in  $\mathcal{P}'$  with the highest  $S_{ACC}(\pi)$ , breaking ties as in

Appendix A.2 (fewer summary edges, includes a sanitizer, fewer files);

**Step 4: Structural validation on  $\hat{\pi}$ .**

Check that:

1. CFG is realizable end-to-end (including exceptional edges).
2. Any non-CFG hop is justified by DFG taint transport or accumulated guards.
3. The interprocedural stack is well-nested (push on CALL, pop on matching RET2\*).
4. Alias consistency holds (non-empty points-to intersection for memory hops).

**Step 5: Return.****if** all checks pass **then**

```

    | return  $\hat{\pi}$ ;

```

**else**

```

    | discard  $\hat{\pi}$  and repeat Step 3 with the next-best path in  $\mathcal{P}'$ ; if none remain,

```

```

    | return None.

```

---

---

**Algorithm 4:** Counterfactual Intervention Test

---

**Input** : Graph slice of a positive example; recorded ACC hooks and baseline prediction  $p$  and chain score  $S_{\text{ACC}}$ .

**Output**: Probability change  $\Delta p$  and chain-score change  $\Delta S_{\text{ACC}}$ .

**Step 1: Choose an intervention.** Select exactly one of:

1. Guard strengthening on a sanitizer that dominates the sink.
2. Sink neutralization (replace the dangerous sink with a benign equivalent).
3. Call unbinding (remove the ARG→PARAM edge that carries taint).

**Step 2: Apply and re-run.**

1. Modify the graph slice according to the chosen intervention.
2. Re-encode the slice and re-run ACC-constrained decoding.
3. Obtain the new prediction  $p'$  and chain score  $S'_{\text{ACC}}$ .

**Step 3: Measure causal effect.**

$$\Delta p \leftarrow p' - p, \quad \Delta S_{\text{ACC}} \leftarrow S'_{\text{ACC}} - S_{\text{ACC}}.$$

**return** ( $\Delta p, \Delta S_{\text{ACC}}$ );

---

The four procedures above correspond exactly to the equations and decoding rules described in Sections 3.2.2–3.2.5.

## Appendix B

# Model Equations and Training

This chapter centralizes the equations referenced by Section 3.2 and provides symbol glossaries.

### B.1 GraphCodeBERT Feature Initialization: Equations

**Token averaging.** Given token embeddings  $\mathbf{E} \in \mathbb{R}^{T \times 768}$  and the token index set  $\mathcal{T}(i)$  aligned to node  $v_i$ ,

$$\tilde{\mathbf{x}}_i^{\text{text}} = \frac{1}{|\mathcal{T}(i)|} \sum_{t \in \mathcal{T}(i)} \mathbf{E}_t \in \mathbb{R}^{768}. \quad (\text{B.1})$$

**Edge-aware token attention.** Let  $\phi(t) \in \mathbb{R}^k$  be local data-flow attributes for token  $t$ . With parameters  $\mathbf{W}_e$  and  $\mathbf{u}$ :

$$\alpha_{i,t} = \frac{\exp(\mathbf{u}^\top \tanh(\mathbf{W}_e [\mathbf{E}_t \parallel \phi(t)]))}{\sum_{s \in \mathcal{T}(i)} \exp(\mathbf{u}^\top \tanh(\mathbf{W}_e [\mathbf{E}_s \parallel \phi(s)]))}, \quad (\text{B.2})$$

$$\mathbf{x}_i^{\text{text}} = \sum_{t \in \mathcal{T}(i)} \alpha_{i,t} \mathbf{E}_t \in \mathbb{R}^{768}. \quad (\text{B.3})$$

**Projection and gated fusion.** Let  $\mathbf{x}_i^{\text{struct}} \in \mathbb{R}^{d_s}$  denote structural features:

$$\mathbf{h}_i^{\text{text}} = \text{LN}(\mathbf{W}_t \mathbf{x}_i^{\text{text}}), \quad \mathbf{h}_i^{\text{struct}} = \text{LN}(\mathbf{W}_s \mathbf{x}_i^{\text{struct}}), \quad (\text{B.4})$$

$$g_i = \sigma(\mathbf{w}_g^\top [\mathbf{h}_i^{\text{text}} \parallel \mathbf{h}_i^{\text{struct}}] + b_g), \quad (\text{B.5})$$

$$\mathbf{h}_i^{(0)} = g_i \mathbf{h}_i^{\text{text}} + (1-g_i) \mathbf{h}_i^{\text{struct}} \in \mathbb{R}^{d_0}. \quad (\text{B.6})$$

## Notation & Symbols for GraphCodeBERT

Symbol	Meaning
$\{t_1, \dots, t_T\}$	Subword tokens for a node fragment; $T=\#\text{tokens}$ .
$\mathbf{E} \in \mathbb{R}^{T \times 768}$	GraphCodeBERT contextual token embeddings.
$\mathcal{T}(i)$	Token indices aligned to node $v_i$ (by character span).
$\bar{\mathbf{x}}_i^{\text{text}}$	Token average (Eq. B.1).
$\phi(t) \in \mathbb{R}^k$	Local attributes (def/use role, fan-in/out, etc.).
$\alpha_{i,t}$	Attention weight (Eq. B.2).
$\mathbf{x}_i^{\text{text}}$	Attention-pooled textual vector (Eq. B.3).
$\mathbf{x}_i^{\text{struct}}$	Structural features (type, degrees, SSA hints, literal stats).
$\mathbf{W}_t, \mathbf{W}_s$	Projections to width $d_0$ (Eq. B.4).
$g_i$	Scalar gate balancing text vs. structure (Eq. B.5).
$\mathbf{h}_i^{(0)}$	Fused node initialization (Eq. B.6).
$d_0$	Hidden width of first GAT layer.
$L, S$	Token window length and stride (typ. $L=512, S=384$ ).

## B.2 Relation-Aware GAT: Equations

**Typed attention and update.** For  $\mathcal{G} = (\mathcal{V}, \{\mathcal{E}_r\}_{r \in \mathcal{R}})$  and neighbors  $\mathcal{N}_r(v)$ :

$$\alpha_{u \rightarrow v}^{(r, \ell)} = \text{softmax}_{u \in \mathcal{N}_r(v)} \left( \text{LeakyReLU} \left( a_r^\top [W_r^{(\ell)} h_u^{(\ell)} \parallel W_0^{(\ell)} h_v^{(\ell)}] \right) \right), \quad (\text{B.7})$$

$$h_v^{(\ell+1)} = \text{ELU} \left( W_{\text{self}}^{(\ell)} h_v^{(\ell)} + \sum_{r \in \mathcal{R}} \sum_{u \in \mathcal{N}_r(v)} \alpha_{u \rightarrow v}^{(r, \ell)} W_r^{(\ell)} h_u^{(\ell)} \right). \quad (\text{B.8})$$

**Node/seed/edge heads.**

$$z_v = w_{\text{node}}^\top h_v^{(L)} + b_{\text{node}}, \quad s_v = w_{\text{seed}}^\top h_v^{(L)} + b_{\text{seed}}, \quad (\text{B.9})$$

$$c_{u \rightarrow v}^{(r)} = h_u^{(L)\top} B_r h_v^{(L)} + \beta_r. \quad (\text{B.10})$$

**Seeds and path score.**

$$\mathcal{S}_K = \text{TopK}(\{s_v : v \in \mathcal{V}\}, K), \quad (\text{B.11})$$

$$S(\pi) = \log \sigma(s_{v_0}) + \sum_{t=1}^T (\alpha \log \sigma(z_{v_t}) + (1 - \alpha) c_{v_{t-1} \rightarrow v_t}^{(r_t)}). \quad (\text{B.12})$$

**Notation & Symbols for Relation-Aware GAT**

Symbol	Meaning
$\mathcal{G} = (\mathcal{V}, \{\mathcal{E}_r\})$	Heterogeneous program graph; edges typed by $r \in \mathcal{R}$ .
$\mathcal{R}$	Relation set (AST, CFG, DFG, CALL, ARG2PARAM, RET2CALL, RET2LHS, DFG_THIN).
$h_v^{(\ell)}$	Node state at layer $\ell$ ; $h_v^{(0)}$ from fusion.
$\mathcal{N}_r(v)$	$r$ -neighbors of node $v$ .
$W_{\text{self}}^{(\ell)}, W_0^{(\ell)}, W_r^{(\ell)}$	Projections per layer/relation.
$a_r$	Relation-specific attention vector.
$\alpha_{u \rightarrow v}^{(r, \ell)}$	Attention weight (Eq. B.7).
$z_v, s_v$	Node vulnerability logit; seed score (Eq. B.9).
$c_{u \rightarrow v}^{(r)}$	Edge compatibility (Eq. B.10).
$B_r, \beta_r$	Bilinear form and scalar prior per relation.
$S_K$	Top- $K$ seeds by $s_v$ (Eq. B.11).
$S(\pi)$	Log-additive path score (Eq. B.12).
$\alpha$	Node/edge evidence mixing weight in $S(\pi)$ .

**B.3 Causal Knowledge Graph (CKG) Prior: Equation****Score mixture.**

$$S^*(\pi) = S(\pi) + \lambda S_{\text{CKG}}(\pi), \quad (\text{B.13})$$

where  $S_{\text{CKG}}(\pi)$  is a small prior bonus computed from unigram/bigram/trigram statistics mined from training chains;  $\lambda \geq 0$  is a decoding-time weight.

**Notation & Symbols for the CKG Prior**

Symbol	Meaning
$\widehat{P}(r)$	Empirical unigram prior over relations from training chains.
$\widehat{P}(r_t   r_{t-1})$	Empirical bigram transition prior.
$S_{\text{CKG}}(\pi)$	Prior score computed from the relation sequence of $\pi$ .
$\lambda$	Mixture weight in Eq. B.13.
$S^*(\pi)$	Decoding score with CKG prior.

**B.4 Adaptive Causal Contextualization (ACC): Equations****Feasibility predicates.**

$$\text{cfg\_ok}(u \rightarrow v) := \mathbf{1}[R_{\text{CFG}}(u, v) = 1], \quad (\text{B.14})$$

$$\text{dfg\_ok}(u \xrightarrow{r} v, \mathcal{T}(\pi)) := \mathbf{1}[(r = \text{DFG}) \wedge \text{tainted}(u, \mathcal{T}(\pi))] \vee \mathbf{1}[\exists g \in \mathcal{G}(\pi) : v \text{ is control dependent on } g], \quad (\text{B.15})$$

$$\text{ipa\_ok}(u \xrightarrow{r} v, C(\pi)) := \begin{cases} \text{push}(\text{callee}(v), \text{site}=u) & r = \text{CALL}, \\ \text{top of stack matches site of } u & r \in \{\text{RET2CALL}, \text{RET2LHS}\}, \\ 1 & \text{otherwise,} \end{cases} \quad (\text{B.16})$$

$$\text{alias\_ok}(u \xrightarrow{r} v, \mathcal{T}(\pi)) := \begin{cases} \mathbf{1}[\text{pts}(u) \cap \text{pts}(v) \neq \emptyset] & \text{if } r \text{ dereferences or writes,} \\ 1 & \text{otherwise.} \end{cases} \quad (\text{B.17})$$

$$\text{Adm}(u \xrightarrow{r} v \mid \pi) := \text{cfg\_ok} \cdot \text{dfg\_ok} \cdot \text{ipa\_ok} \cdot \text{alias\_ok}. \quad (\text{B.18})$$

### Role penalties and sanitizer bonus.

$$\text{pen}_{\text{start}}(\pi) := \lambda_{\text{start}} [1 - p_{v_0}^{(\text{src})}]_+, \quad (\text{B.19})$$

$$\text{pen}_{\text{mid}}(\pi) := \lambda_{\text{mid}} \sum_{t=1}^{T-1} \left( 1 - \max\{p_{v_t}^{(\text{prop})}, p_{v_t}^{(\text{san})}\} \right), \quad (\text{B.20})$$

$$\text{pen}_{\text{end}}(\pi) := \lambda_{\text{end}} [1 - p_{v_T}^{(\text{sink})}]_+, \quad (\text{B.21})$$

$$\text{pen}_{\text{role}}(\pi) := \text{pen}_{\text{start}} + \text{pen}_{\text{mid}} + \text{pen}_{\text{end}}. \quad (\text{B.22})$$

$$\text{bonus}_{\text{san}}(\pi) := \mu \sum_{s \in \pi} \mathbf{1}[\text{Dom}(s, \text{sink})] \cdot \mathbf{1}[\text{affects\_taint}(s, \pi)]. \quad (\text{B.23})$$

### ACC objective and CKG mixture.

$$\text{pen}_{\text{len}}(\pi) := \eta \text{len}(\pi), \quad (\text{B.24})$$

$$\text{pen}_{\text{rep}}(\pi) := \rho \text{rep}(\pi), \quad (\text{B.25})$$

$$S_{\text{ACC}}(\pi) := S(\pi) - \text{pen}_{\text{role}}(\pi) + \text{bonus}_{\text{san}}(\pi) - \text{pen}_{\text{len}}(\pi) - \text{pen}_{\text{rep}}(\pi), \quad (\text{B.26})$$

$$S_{\text{ACC}}^*(\pi) = S_{\text{ACC}}(\pi) + \lambda S_{\text{CKG}}(\pi). \quad (\text{B.27})$$

### Notation & Symbols for ACC

Symbol	Meaning
$\pi = (v_0, \dots, v_T)$	Candidate path (chain).
$\mathcal{T}(\pi), C(\pi), \mathcal{G}(\pi)$	Taint footprint, call stack, accumulated guards.
Predicates	$\text{cfg\_ok}, \text{dfg\_ok}, \text{ipa\_ok}, \text{alias\_ok}$ (Eqs. B.14–B.17).
$\text{Adm}(\cdot)$	Conjunctive admissibility (Eq. B.18).
$p_v^{(\text{src}/\text{prop}/\text{san}/\text{sink})}$	Role probabilities.
$\lambda_{\text{start}}, \lambda_{\text{mid}}, \lambda_{\text{end}}$	Role penalty weights.
$\mu$	Sanitizer bonus weight.
$\eta, \rho$	Length and repetition penalty weights.
$S_{\text{ACC}}(\pi), S_{\text{ACC}}^*(\pi)$	ACC scores without/with CKG (Eqs. B.26, B.27).

## B.5 CKG Prior Scoring

Let  $\pi(i \xrightarrow{r} j)$  be the mined prior over relation motifs from the training split. With smoothing  $\epsilon$  and temperature  $\tau$ ,

$$\tilde{\pi}(i \xrightarrow{r} j) = \frac{\pi(i \xrightarrow{r} j) + \epsilon}{\sum_{k \in N(i)} (\pi(i \xrightarrow{r_k} k) + \epsilon)}, \quad \pi_\tau(i \xrightarrow{r} j) = \frac{\tilde{\pi}(i \xrightarrow{r} j)^{1/\tau}}{\sum_{k \in N(i)} \tilde{\pi}(i \xrightarrow{r_k} k)^{1/\tau}}.$$

During decoding, admissible edges are scored by mixing ACC/beam score  $s(i \xrightarrow{r} j)$  with the prior:

$$S_{t+1} = S_t + s(i \xrightarrow{r} j) + \lambda \log \pi_\tau(i \xrightarrow{r} j),$$

where  $\lambda$  is the prior weight,  $\epsilon$  the add-one smoothing, and  $\tau$  the temperature. Numeric settings are listed in Appendix. D.1.

## B.6 Chain Extraction and Validation: Equations

**Role-shaped validity criteria.**

$$p_{v_0}^{(\text{src})} \geq \tau_{\text{src}}, \quad (\text{B.28})$$

$$\exists t \in \{1, \dots, T-1\} : \max(p_{v_t}^{(\text{prop})}, p_{v_t}^{(\text{san})}) \geq \tau_{\text{mid}}, \quad (\text{B.29})$$

$$p_{v_T}^{(\text{sink})} \geq \tau_{\text{sink}}. \quad (\text{B.30})$$

**Selection.**

$$\hat{\pi} = \arg \max_{\pi \in \mathcal{P}_{\text{adm}}} S_{\text{ACC}}(\pi). \quad (\text{B.31})$$

### Notation & Symbols for Chain Extraction

Symbol	Meaning
$\tau_{\text{src}}, \tau_{\text{mid}}, \tau_{\text{sink}}$	Thresholds for start/middle/end criteria (Eqs. B.28–B.30).
$\mathcal{P}_{\text{adm}}$	Set of admissible paths under ACC constraints.
$\hat{\pi}$	Selected chain (Eq. B.31).

## B.7 Training Objective: Equation

**Composite loss.**

$$\mathcal{L} = \text{BCE}(\sigma(\mathbf{z}), \hat{y}) + \lambda_{\text{flow}} \mathcal{L}_{\text{flow}} + \lambda_{\text{cf}} \mathcal{L}_{\text{cf}} + \lambda_{\text{ent}} \mathcal{L}_{\text{ent}} + \lambda_{\text{spec}} \mathcal{L}_{\text{spec}} + \lambda_{\text{san}} \mathcal{L}_{\text{san}}. \quad (\text{B.32})$$

**Notation & Symbols for Training Objective**

Symbol	Meaning
$\hat{y} \in \{0, 1\}$	Graph label (vulnerable / non-vulnerable).
$\mathbf{z} \in \mathbb{R}$	Graph logit (pooled over final node states).
BCE	Binary cross-entropy.
$\mathcal{L}_{\text{flow}}$	Flow consistency (upscores chain edges/paths, downscores distractors).
$\mathcal{L}_{\text{cf}}$	Counterfactual consistency (confidence drop on minimal edits).
$\mathcal{L}_{\text{ent}}$	Attention entropy (encourages decisive attention).
$\mathcal{L}_{\text{spec}}$	Spectral norm control (stability).
$\mathcal{L}_{\text{san}}$	Sanitizer alignment (when sanitizer dominates sink).
$\lambda_{\star} \geq 0$	Weights for the corresponding regularizers.



## Appendix C

# Metrics and Diagnostics

### C.1 Calibration and Thresholding

#### C.1.1 Operating Points

Validation-optimal F1 threshold  $\tau_{F1^*}$  and fixed threshold  $\tau=0.5$  are used for reporting.

#### C.1.2 Expected Calibration Error (ECE)

$$\text{ECE} = \sum_{b=1}^B \frac{|S_b|}{N} |\text{acc}(S_b) - \text{conf}(S_b)|. \quad (\text{C.1})$$

$S_b$  are prediction bins,  $N$  is the number of examples.

### C.2 Standard Classification Metrics

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}}, \quad (\text{C.2})$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad (\text{C.3})$$

$$\text{F1} = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}, \quad (\text{C.4})$$

$$\text{MacroF1} = \frac{1}{2} (\text{F1}_{\text{pos}} + \text{F1}_{\text{neg}}), \quad \text{MicroF1} = \frac{2 \sum \text{TP}}{2 \sum \text{TP} + \sum \text{FP} + \sum \text{FN}}. \quad (\text{C.5})$$

AUROC and AUPRC are computed on raw scores.

### C.3 Chain-Centric Metrics

#### C.3.1 Validity of Predicted Chains

$$\text{ValidityRate} = \frac{\#\{\text{predicted chains passing all ACC checks}\}}{\#\{\text{predicted chains}\}}. \quad (\text{C.6})$$

#### C.3.2 Structural Fidelity

Let predicted chain nodes/edges be  $(\hat{V}, \hat{E})$  and ground truth  $(V^*, E^*)$ .

$$\text{NodeCov} = \frac{|\hat{V} \cap V^*|}{|V^*|}, \quad \text{EdgeCov} = \frac{|\hat{E} \cap E^*|}{|E^*|}, \quad (\text{C.7})$$

$$\text{LCS Ratio} = \frac{\text{LCS}(\hat{\pi}, \pi^*)}{|\pi^*|}, \quad (\text{C.8})$$

$$\text{CO}(\alpha) = \alpha \cdot \text{NodeCov} + (1-\alpha) \cdot \text{EdgeCov}, \quad \alpha \in [0, 1], \quad (\text{C.9})$$

$$\text{RoleCov}_r = \frac{|\hat{V}_r \cap V_r^*|}{|V_r^*|}, \quad r \in \{\text{src}, \text{san}, \text{prop}, \text{sink}\}. \quad (\text{C.10})$$

#### C.3.3 Interprocedurality (IPA)

$$\text{IPA Rate} = \frac{\#\{\hat{\pi} \text{ using CALL/ARG} \rightarrow \text{PARAM/RET} \rightarrow \text{CALLER or LHS}\}}{\#\{\text{predicted chains in slices containing such edges}\}}. \quad (\text{C.11})$$

### C.4 Counterfactual Metrics

#### C.4.1 Counterfactual Consistency Score (CCS)

For applicable graphs  $i \in \mathcal{I}$ , with original  $p_i$  and intervened  $p_i^{\text{do}}$ ,

$$\text{CCS}_i = (p_i - p_i^{\text{do}})^2, \quad (\text{C.12})$$

$$\text{CCS} = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \text{CCS}_i. \quad (\text{C.13})$$

Directional consistency uses  $s_i \in \{+1, -1\}$ :

$$\Delta_i^{\text{dir}} = s_i (p_i - p_i^{\text{do}}), \quad \text{DCR} = \frac{1}{|\mathcal{I}|} \sum_i \mathbb{1}\{\Delta_i^{\text{dir}} > 0\}. \quad (\text{C.14})$$

### C.4.2 Causal Feature Attribution Measure (CFAM)

Let on-chain features  $F_c$  and off-chain  $F_s$  with attributions  $A_i(f) \geq 0$ ,

$$\text{CFAM}_i = \frac{\sum_{f \in F_c} |A_i(f)|}{\sum_{f \in F_c \cup F_s} |A_i(f)|}, \quad \text{CFAM} = \frac{1}{N} \sum_i \text{CFAM}_i. \quad (\text{C.15})$$

Attributions are normalized per graph.

## C.5 Decoding diagnostics

For completeness, I report the auxiliary decoding diagnostics used in the artifact. The Chain Success Rate (CSR), Admissible Expansion Ratio (AER), and Prior Influence Rate (PIR) are defined as:

$$\text{CSR} = \frac{\#\{\text{positives with a nonempty chain}\}}{\#\{\text{positives}\}}, \quad (\text{C.16})$$

$$\text{AER} = \frac{\#\{\text{candidate edges admitted by ACC}\}}{\#\{\text{candidate edges considered}\}}, \quad (\text{C.17})$$

$$\text{PIR} = \frac{\#\{\text{steps where the CKG prior changes the winning rank}\}}{\#\{\text{decoding steps}\}}. \quad (\text{C.18})$$

Average chain length, interprocedural hop ratio, motif coverage@K, and the ACC rejection mix (CFG/DFG/IPA/alias) are summarized as descriptive statistics in the released logs, but are not part of the main quantitative analysis in this thesis.

## Appendix D

# Extended Results and Settings

This appendix centralizes configuration details, intermediate diagnostics, and practical instructions to reproduce the reported tables and figures.

### D.1 Hyperparameters and Environment

#### Recommended defaults

$K=8$ ,  $B=24$ ,  $H=5$ ,  $\alpha=0.7$ ,  $\lambda=0.2$ ,  $\epsilon=10^{-3}$ ,  $\tau=1.0$ , motif top- $K=500$ .

## Model and optimization settings

Component	Setting
GAT depth / width	$L=3$ layers; hidden width $d=64$ ; one head per relation
Nonlinearity & dropout	ELU; dropout 0.1 on node states and attention logits
Graph pooling	Attention pooling over final node states
Optimizer / schedule	AdamW; lr $2 \times 10^{-3}$ ; wd $10^{-4}$ ; cosine decay (40 epochs), 2-epoch warmup
Gradient control	Global norm clip 1.0; mixed precision (FP16 projections/attention, FP32 accumulation)
Aux losses (if enabled)	$\lambda_{\text{flow}}=0.5$ , margin 0.5; $\lambda_{\text{cf}}=0.5$ , margin 0.7; $\lambda_{\text{ent}}=0.05$ ; $\lambda_{\text{spec}}=10^{-4}$
ACC / beams	Seeds $K=8$ ; width $B=24$ ; horizon $H=5$ ; node/edge mix $\alpha=0.7$
Batching	Batch size chosen to fit device memory (typically 4–8 graphs)
Augmentations	Identifier renaming, inert code insertion, in-basic-block reordering (each prob. 0.3)
Language model	GraphCodeBERT frozen; ablation unfreezes last 2 blocks at $0.1 \times \text{lr}$
Hardware & software	PyTorch 2.4.1+cu121; CUDA 12.1; single NVIDIA GeForce RTX 4070 Laptop GPU

## Software and hardware manifest

Python version, key libraries, and a lockfile with package hashes are included with the artifact for exact replay.

## D.2 Beam Search and ACC Diagnostics

The table summarizes per-graph decoding behavior under the default settings ( $K=8$ ,  $B=24$ ,  $H=5$ ,  $\alpha=0.7$ ). Values are medians with P90 where noted.

Table D.1: Beam and ACC diagnostics under  $K=8$ ,  $B=24$ ,  $H=5$  (per graph).

Quantity	Median	P90	Struct-only	GCBERT+Struct
Seeds used (of 8)	6.0	8.0	5.8	6.3
Avg. hops of best chain	4.3	5.0	4.2	4.4
Admissible branching factor $\bar{d}$	2.1	3.0	2.0	2.2
Expansions attempted	1,240	1,920	1,180	1,300
Expansions admitted	168	256	159	177
Pruned by CFG reachability	54.1%	62.7%	55.9%	52.5%
Pruned by IPA/stack rule	21.8%	27.4%	22.6%	21.1%
Pruned by alias/points-to	6.8%	9.3%	6.5%	7.1%
Score-pruned (beam cap / entropy)	17.3%	22.1%	15.0%	19.3%
Beams that reach a sink	2.1	3.0	1.9	2.3
Chains using CALL/RET edges	63.4%	74.2%	60.7%	66.0%
Chains with sanitizer bonus	41.6%	51.0%	39.8%	43.2%

### Beam Sensitivity (K, B, H)

The table reports validation metrics and per-graph latency for the GCBERT+Struct model across beam sizes.

Table D.2: Beam sensitivity (Valid, GCBERT+Struct; latency in ms per graph).

K	B	H	Validity	IPA	LCS	Latency
4	12	3	0.802	0.653	0.541	<b>10.8</b>
4	12	5	0.821	0.704	0.565	12.6
8	24	5	<b>0.844</b>	<b>0.784</b>	<b>0.604</b>	17.7
8	48	5	0.847	0.792	0.608	24.9
16	48	7	0.849	0.799	0.612	33.1

#### D.2.1 Decoder Diagnostics and Remedies

The table below lists rare decoding issues, their signals, and practical fixes.

Table D.3: Decoder diagnostics for rare cases and suggested remedies.

Symptom	Likely cause	Signal	Remedy
Early termination	True path exceeds $H$ or detours	Low score near $t=H$ , no successors	Increase $H$ or add thin DFG summaries
No interprocedural hop	Source or sink API missing from lexicon	High node score, zero ARG2PARAM/RET2*	Extend lexicon or enable API summaries
Alias mismatch	Coarse points-to for container or pointer hop	alias_ok rejects edge	Refine buckets; add container-aware rules
Guard over-credit	Guard does not dominate the sink region	Bonus lifts wrong branch	Tighten dominance; require taint effect
Rare relation downweighted	Prior penalizes uncommon pattern	Prior penalty visible in path deltas	Reduce mixture $\lambda$ or smooth with larger $K$

## D.3 Additional Diagnostics

### Reported items

- Per-hop rank traces with and without the CKG prior.
- ACC rejection histograms: CFG, def-use/guard, IPA stack, alias conflict.
- Motif coverage@K and interprocedural hop ratio distributions.
- Admissible expansion ratio curves versus decoding depth.

## Appendix E

# Dataset Card and Licensing

**Name and scope:** Experiments use *ReposVul*, a repository-level corpus that links CVE/CWE metadata suitable for interprocedural, chain-centric analysis.

**Composition:** This thesis uses the C/C++ subset with the default splits. Each entry contains a code snapshot, with aligned views at line, function, file, and repository level, and repository-scope caller–callee links.

**Collection and preprocessing:** The pipeline (Chapter 3) performs raw crawling, vulnerability untangling, repository-level dependency extraction, code normalization, patch-aware differencing, and filtering of outdated patches.

**Labels and quality:** Vulnerability labels follow the original *ReposVul* patch pairs. Mechanism roles (source, sanitizer, propagator, sink) are assigned by rules and consistency checks.

**Splits and leakage safeguards:** Default train/validation/test splits are used without modification. Projects do not cross partitions, parent/child patches are never split, and identical files/CVEs are removed across splits to reduce leakage.

**Intended use:** The dataset is used for repository-level vulnerability modeling, with an emphasis on interprocedural chain reconstruction and interpretability.

**Licensing and attribution:** Use of *ReposVul* follows the terms specified by its authors.

**Known limitations:** Not all vulnerabilities admit a single clear executable chain, and some fixes are entangled with refactoring.



## Appendix F

# Role Lexicon and Pattern Rules

This appendix lists the APIs, idioms, and structural patterns used to assign *source*, *sanitizer*, *propagator*, and *sink* roles (Section 3.1.4).

### F.1 API Families and Examples

Role	Representative APIs and idioms
Source	recv, read, fgets, environment access, deserialization entry points
Sanitizer	bounds and range checks, length clamps, whitelist validation, null checks, defensive copies
Propagator	assignments, pointer dereference/address-of, argument→parameter, return→caller
Sink	memcpy, strcpy, indexed writes, command execution, path joins followed by file I/O

### F.2 Structural Patterns

1. Guards that dominate the sink and constrain tainted values (e.g., index/length checks on paths from source to sink).
2. Def-use chains that cross calls via actual/formal bindings (ARG→PARAM; RET→CALLER/RET→LHS).
3. Alias flows through pointers, references, or containers where points-to sets intersect along writes and reads.

**Precedence.** If multiple roles apply to a node, the precedence is *sanitizer* > *sink* > *propagator* > *source*, resolved using CFG dominance and DFG reachability (see Chapter 3).

## Appendix G

# Reproducibility and Artifact

This appendix summarizes the released artifact and how it supports end-to-end reproducibility, from data preparation to model training and evaluation.

### G.1 Data provenance

Each split ships a manifest with repository names, commit hashes, file paths, and per-example checksums. Parent/child commit IDs are recorded so vulnerable and fixed snapshots can be reconstructed from the original repositories.

### G.2 Graph construction

The graph builder configuration includes parser versions, normalization options, and enabled relation families: AST/CFG/DFG, CALL, ARG2PARAM, RET2CALL, RET2LHS, and alias summaries.

### G.3 Embeddings and caches

The artifact records the GraphCodeBERT checkpoint and tokenizer IDs, maximum sequence length and stride, and normalization flags.

### G.4 Training and decoding configuration

Versioned configuration files specify encoder and decoder settings, including learning rate, schedule, weight decay, dropout, gradient clipping, loss weights, beam parameters  $(K, B, H, \alpha)$ , augmentation probabilities, and early stopping criteria. ACC and CKG parameters are included so both training and decoding behaviour can be reproduced.

## G.5 Environment and determinism

Experiments use PyTorch 2.4.1 with CUDA 12.1 on a single GPU. Random seeds are fixed, deterministic and cuDNN flags are set where possible, and mixed precision uses dynamic loss scaling with critical reductions in FP32. An environment manifest (lockfile and package hashes) is included.

## G.6 Artifact contents

The released artifact contains:

- Per-graph predictions and reconstructed chains (with beam traces).
- Calibration parameters and ECE bin assignments.
- Counterfactual intervention logs (pre/post probabilities and edit types).
- Configuration files for all reported experiments.
- Environment manifests and package hashes.
- CKG prior data and motif statistics.
- Decoding diagnostics (per-hop scores and ACC rejection reasons).
- Scripts to rebuild graphs, train models, export chains, and recompute all reported metrics and figures.

## G.7 Project layout

```
Thesis-causal-vul/  
|-- WORKPLAN.md  
|-- README.md  
|-- commands.txt  
|-- ScripstList.txt  
|-- requirements.txt  
|-- notebooks/  
|   |-- 01_ReposVul.ipynb  
|   |-- 02_word2vec_training.ipynb  
|   |-- 03_causal_chain_demo.ipynb  
|   '-- colab.ipynb  
|-- tools/  
|   |-- export_pdg_env.sc  
|   '-- validate_pdg_dir.py
```

```

|-- Report/
|   |-- Causalchain_pretty.png
|   |-- Causalchain_raw.png
|   |-- eval_colab/
|       |-- eval_report.json
|       |-- test_preds.csv
|   '-- ReposVul_report/
|       |-- ReposVul_CCPP_Full_Report.md
|       '-- _edges/
|           |-- test_c_cpp_repository2__caller_func_callee.csv
|           |-- train_c_cpp_repository2__caller_func_callee.csv
|           '-- valid_c_cpp_repository2__caller_func_callee.csv
'-- src/
|-- __init__.py
|-- build_dataset_jsonl.py
|-- infer_one_slice_gcbert.py
|-- infer_one_slice_pretty.py
|-- step4_train_gnn_attn.py
|-- step4_train_gnn_attn_v1.py
|-- step5_export_chains.py
|-- step6_eval_report.py
|-- train.py
|-- preprocess/
|   |-- __init__.py
|   |-- embed_lines_gcbert.py
|   |-- slices_to_pyg_gcbert.py
|   |-- step1_prepare_diverse_vul.py
|   |-- step2_generate_pdg.py
|   |-- step2c_pipeline_validate_and_slice.py
|   |-- transformer_cache.py
|   '-- external/
|       |-- program_slice.py
|       |-- sensiAPI.txt
|       |-- sensiAPI_A.txt
|       '-- sensiAPI_B.txt
'-- utils/
|-- hashing.py
'-- pdg_io.py

```

## G.8 Step-by-step reproducibility guide

### Prerequisites

- Python 3.10+ and a CUDA-enabled GPU (for CPU-only, training and decoding will be slower).
- Optional: Joern installed or available on PATH if rebuilding CPG/PDG from raw code.

### 1. Install dependencies

#### Linux

```
python -m pip install -r requirements.txt
```

#### Windows PowerShell

```
python -m pip install -r .\requirements.txt
```

### 2. Prepare datasets and program graphs

From raw sources (ReposVul file lists), run the preprocessing pipeline; **Linux**

```
python src/preprocess/step1_prepare_diverse_vul.py
python src/preprocess/step2_generate_pdg.py
python src/preprocess/step2c_pipeline_validate_and_slice.py
```

#### Windows PowerShell

```
python .\src\preprocess\step1_prepare_diverse_vul.py
python .\src\preprocess\step2_generate_pdg.py
python .\src\preprocess\step2c_pipeline_validate_and_slice.py
```

### 3. Build or load GraphCodeBERT caches

#### Linux

```
python src/preprocess/embed_lines_gcbert.py
```

#### Windows PowerShell

```
python .\src\preprocess\embed_lines_gcbert.py
```

Then convert slices to PyTorch Geometric hetero graphs: **Linux**

```
python src/preprocess/slices_to_pyg_gcbert.py
```

#### Windows PowerShell

```
python .\src\preprocess\slices_to_pyg_gcbert.py
```

#### 4. Assemble training JSONL (if required)

##### Linux

```
python src/build_dataset_jsonl.py
```

##### Windows PowerShell

```
python .\src\build_dataset_jsonl.py
```

#### 5. Train the model

##### Linux

```
python src/step4_train_gnn_attn.py  
# or:  
python src/step4_train_gnn_attn_v1.py
```

##### Windows PowerShell

```
python .\src\step4_train_gnn_attn.py  
# or:  
python .\src\step4_train_gnn_attn_v1.py
```

#### 6. Export chains on the test split

##### Linux

```
python src/step5_export_chains.py
```

##### Windows PowerShell

```
python .\src\step5_export_chains.py
```

#### 7. Compute metrics and reports

##### Linux

```
python src/step6_eval_report.py
```

##### Windows PowerShell

```
python .\src\step6_eval_report.py
```

This produces summary files under Report/ such as:

```
Report/eval_colab/eval_report.json  
Report/eval_colab/test_preds.csv
```

## 8. Optional: pretty print a single slice

### Linux

```
python src/infer_one_slice_pretty.py
```

### Windows PowerShell

```
python .\src\infer_one_slice_pretty.py
```

## G.9 Determinism and validation checks

Seeds are fixed in configs; decoding is deterministic given identical environment and seeds. Checksums verify inputs and expected summaries. A quick sanity script confirms: (i) project-disjoint splits, (ii) at least one valid chain on known positives, and (iii) improved calibration (lower ECE) after validation temperature scaling.

## G.10 Source for Figure 4.1 and Chain Trace

The files below reproduce Figure 4.1. Decoding used ACC with  $K=8$ ,  $B=24$ ,  $H=5$ , and  $\alpha=0.7$ ; the CKG prior weight was  $\lambda=0.2$ . Paths are relative to the project root. The exact decoded chain trace is included for reference.

### app/main.c

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>

char* stage_clean(char* in);
char* stage_route(char* in);
void stage_execute(char* cmd);

int main(int argc, char** argv) {
    char buf[256];
    size_t n;

    memset(buf, 0, sizeof(buf));
    if (fgets(buf, sizeof(buf), stdin) == NULL) {
        return 0;
    }

    n = strlen(buf);
    if (n > 0 && (buf[n - 1] == '\n' || buf[n - 1] == '\r')) {
        buf[n - 1] = '\0';
```

```

}

char* s1 = stage_clean(buf);
char* s2 = stage_route(s1);
stage_execute(s2);
return 0;
}

```

**lib/shell.c**

```

#include <stdio.h>
#include <stdlib.h>
#include <string.h>

static char route_buf1[512];
static char route_buf2[512];
static char route_buf3[512];
static char exec_local[512];

static char* build_prefix(char* in) {
    snprintf(route_buf1, sizeof(route_buf1), "runner_%s", in);
    return route_buf1;
}

static char* add_route_flags(char* in) {
    snprintf(route_buf2, sizeof(route_buf2), "%s_--channel=cli_
... --retry=0", in);
    return route_buf2;
}

static char* finalize_route(char* in) {
    snprintf(route_buf3, sizeof(route_buf3), "%s_--complete", in);
    return route_buf3;
}

char* stage_route(char* in) {
    char* a = build_prefix(in);
    char* b = add_route_flags(a);
    char* c = finalize_route(b);
    return c;
}

void stage_execute(char* cmd) {
    size_t n;
    if (cmd == NULL) {
        return;
    }
}

```



```

}
n = strlen(cmd);
if (n >= sizeof(exec_local)) {
    n = sizeof(exec_local) - 1;
}
memset(exec_local, 0, sizeof(exec_local));
strncpy(exec_local, cmd, n);
exec_local[n] = '\0';
system(exec_local); /* external command execution */
}

```

### Decoded chain trace for Case A

```

Root -> ... -> Sink (hops=7)

[ROOT] if (fgets(buf, sizeof(buf), stdin) == NULL) {
    (CALL, app/main.c, line 15, p=0.475)
    '-- [CFG] exact=CFG
    [->] return 0;
    (BLOCK, app/main.c, line 16, p=0.470)
    '-- [CFG] exact=CFG
    [->] }
(BLOCK, app/main.c, line 17, p=0.528)
 '-- [CFG] exact=CFG
 [->] n = strlen(buf);
 (BLOCK, app/main.c, line 19, p=0.472)
 '-- [CFG/DFG] exact=CFG,DFG
 [->] if (n > 0 && (buf[n - 1] == '\n' || buf[n - 1] == '\r')) {
    (CALL, app/main.c, line 20, p=0.395)
    '-- [CPG] exact=CALL
    [->] n = sizeof(exec_local) - 1;
    (BLOCK, lib/shell.c, line 40, p=0.524)
    '-- [DFG] exact=DFG
    [->] exec_local[n] = '\0';
    (BLOCK, lib/shell.c, line 44, p=0.451)
    '-- [CFG] exact=CFG
    [SINK] system(exec_local); /* external command execution */
    (CALL, lib/shell.c, line 45, p=0.445)
    '-- [CPG] exact=CALL

```

The same trained checkpoint and decoding settings were used for all case studies.