# CLNX Bridging Code and Natural Language for C/C++ Vulnerability-Contributing Commits Identification

Zeqing Qin, Yiwei Wu, Lansheng Han*

*School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China*

HUAZHONG UNIVERSITY OF SCIENCE AND TECHNOLOGY

## WHAT ARE Vulnerability-Contributing Commits (VCCs)?

**Vulnerability-Contributing Commits (VCCs)** are patch commits in open-source software development that unintentionally introduce security vulnerabilities into the codebase. These commits modify source code (e.g., configuration, logic) but contain flaws classified under Common Weakness Enumeration (CWE), regardless of exploit conditions.

For example, a VCC might alter file permissions in a way that enables privilege bypass. VCCs are critical for vulnerability analysis, as they reveal how and where security risks originate during code evolution.

## WHY CLNX IS NEEDED?

**BERT-based LLMs** show promise in vulnerability detection by learning code context and dependencies. However, they struggle with C/C++—the source of 52% of open-source vulnerabilities. This gap stems from insufficient C/C++-specific pre-training and inherent challenges like **complex syntax**, **low-level operations**, and **code redundancy**. While current solutions rely on further pre-training (e.g., CodeBERT-cpp), they yield minimal gains (+2.03% accuracy) at high computational costs, proving inefficient and limited.
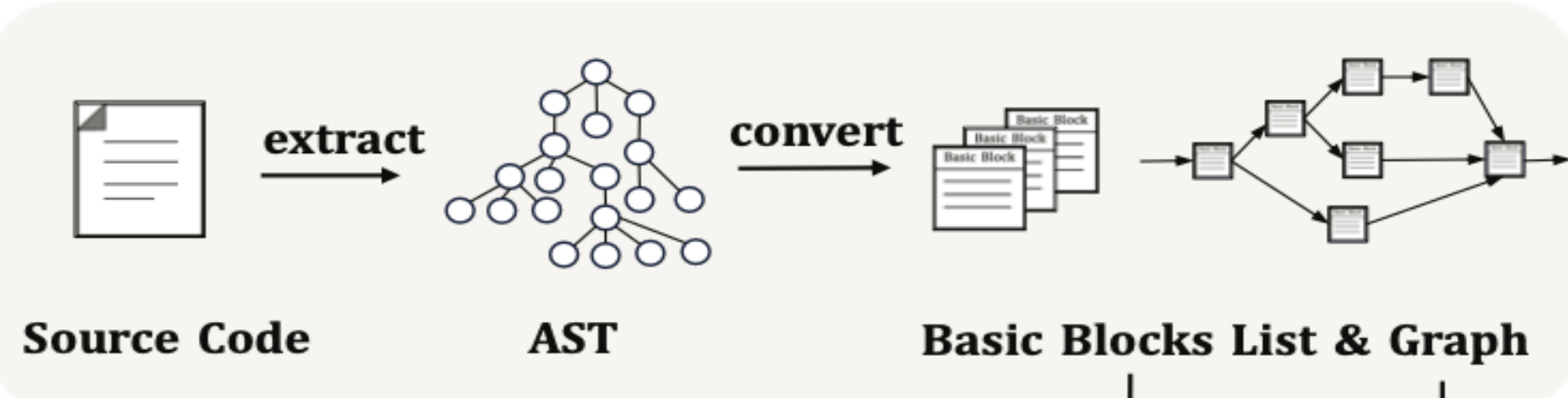
**CLNX aims to enhance the effectiveness of LLMs in identifying C/C++ VCCs while ensuring a lightweight implementation.**
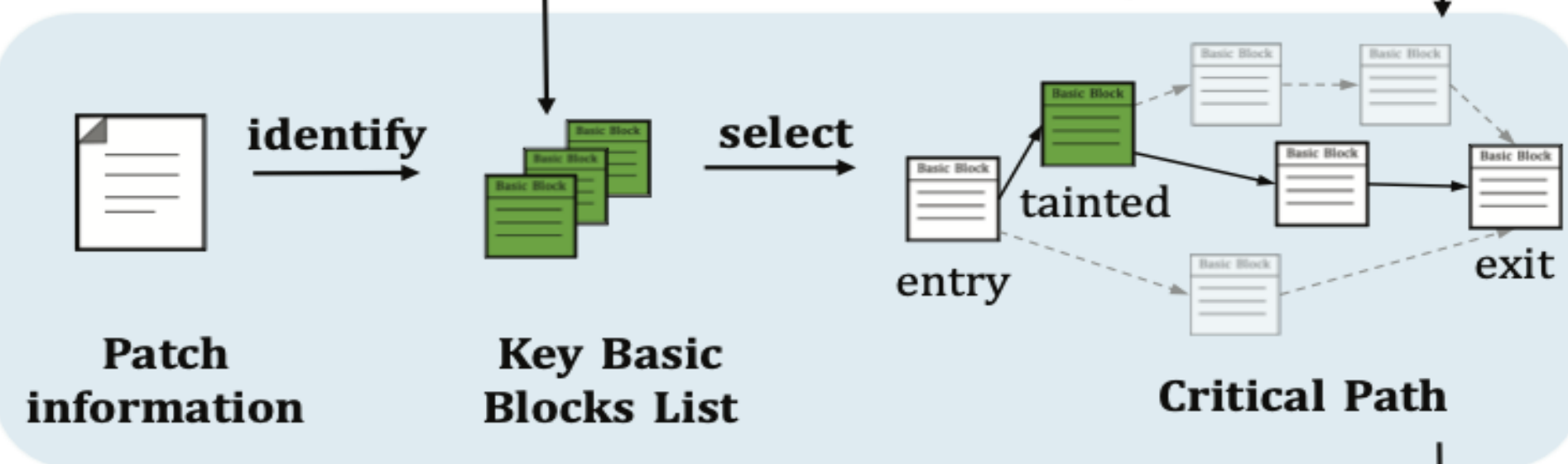
## HOW CLNX WORKS?

CLNX simplifies C/C++ code analysis for LLMs through two-step naturalization, bypassing heavy pre-training:

1. **Structural Naturalization:**
Transforms code into execution flow graphs and pinpoints critical paths using patch data.

2. **Token Naturalization:**
Converts critical symbols (e.g., variables, functions) into plain-language terms (e.g., "user_input")

3. **LLM-Friendly Output:**
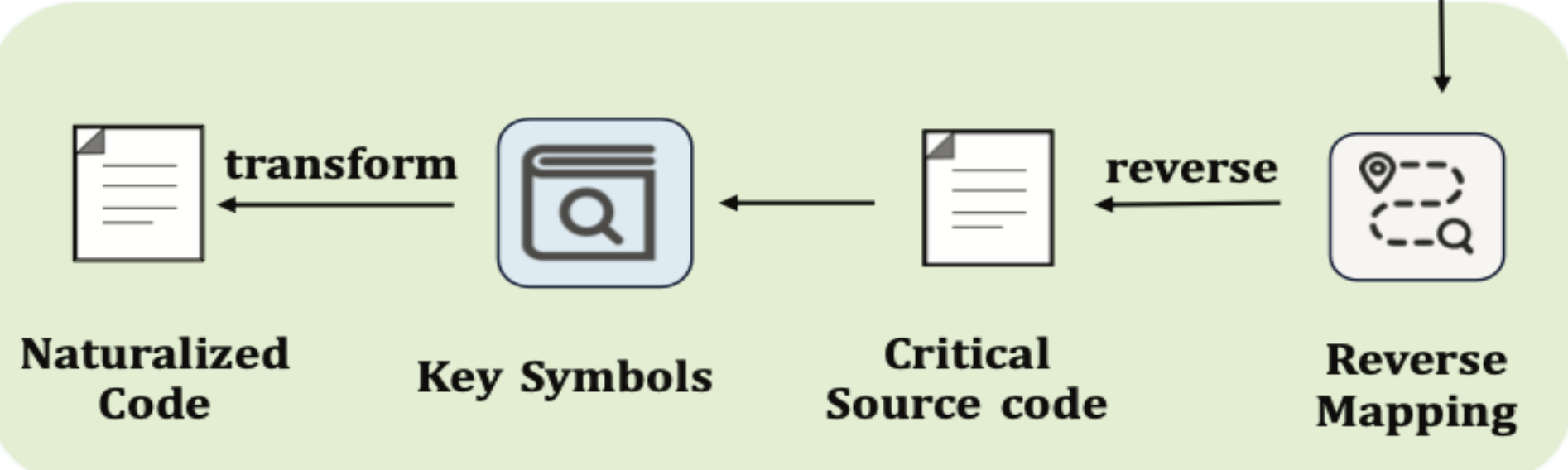Delivers simplified, naturalized code for LLMs to detect C/C++ VCCs efficiently.
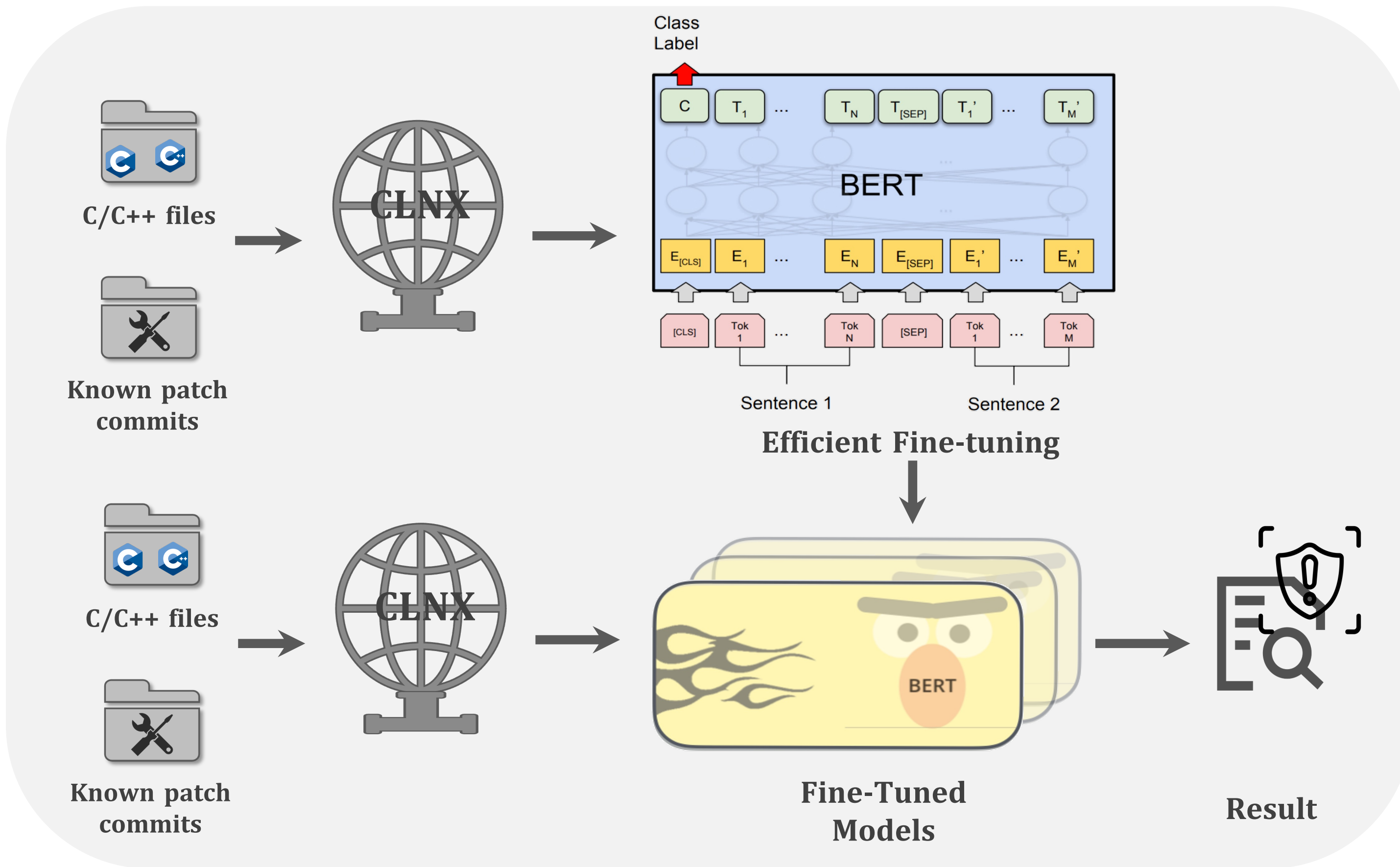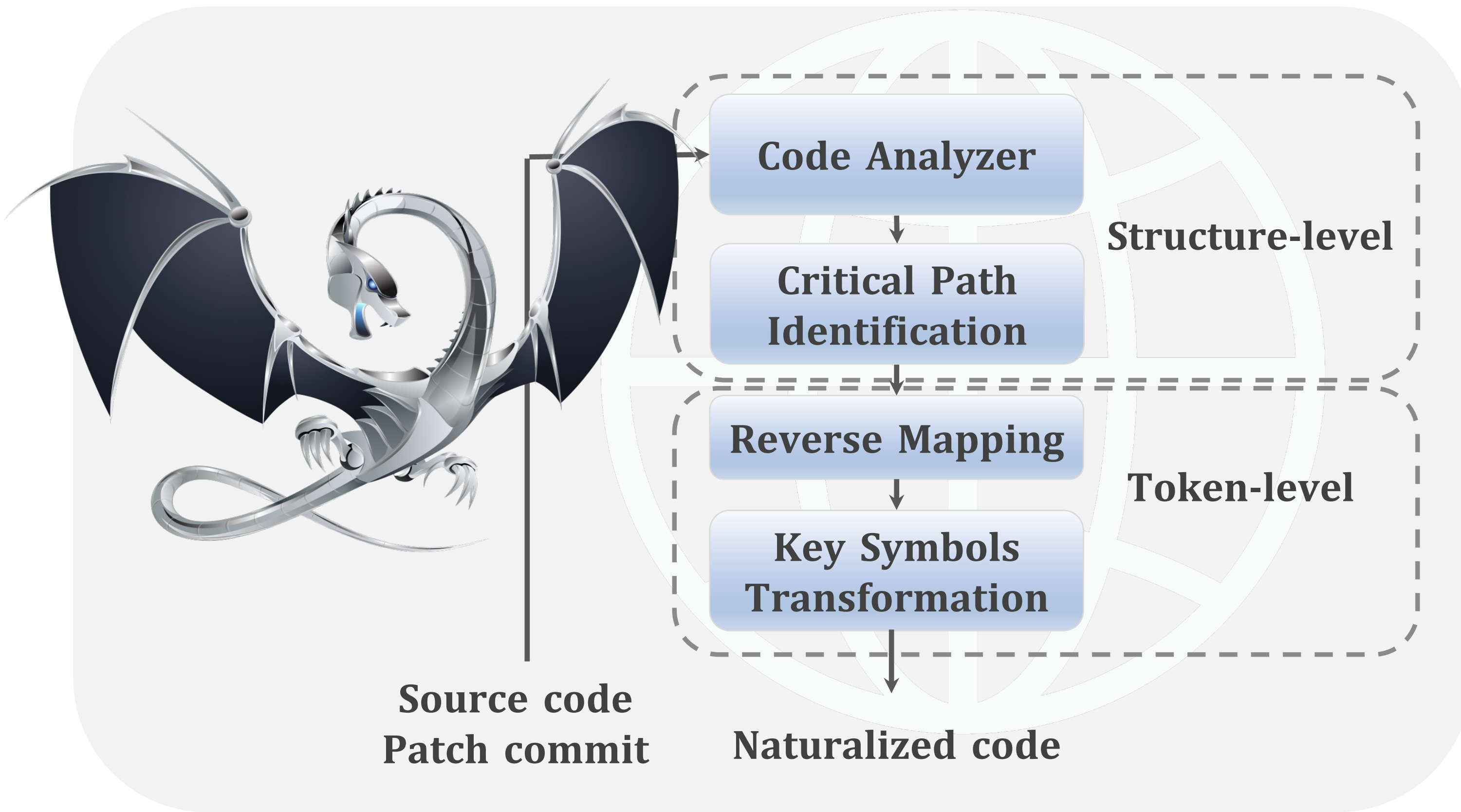
### Step 1: Code Analyzer



extract → AST → convert → Basic Blocks List & Graph

Source Code → AST → Basic Blocks List & Graph

### Step 2: Critical Path Identification

identify → select

Patch information → Key Basic Blocks List → tainted entry / exit → Critical Path

### Step 3: Reverse Mapping & Symbol Transformation

transform ← reverse

Naturalized Code ← Key Symbols ← Critical Source code ← Reverse Mapping

## CLNX ARCHITECTURE



Code Analyzer → Critical Path Identification — Structure-level

Reverse Mapping → Key Symbols Transformation — Token-level

Source code / Patch commit → Naturalized code

C/C++ files + Known patch commits → CLNX → BERT — Efficient Fine-tuning

Class Label → C T₁ ... Tₙ T[SEP] T₁' ... Tₘ'

E[CLS] E₁ ... Eₙ E[SEP] E₁' ... Eₘ'

[CLS] Tok 1 ... Tok N [SEP] Tok 1 ... Tok M

Sentence 1 / Sentence 2

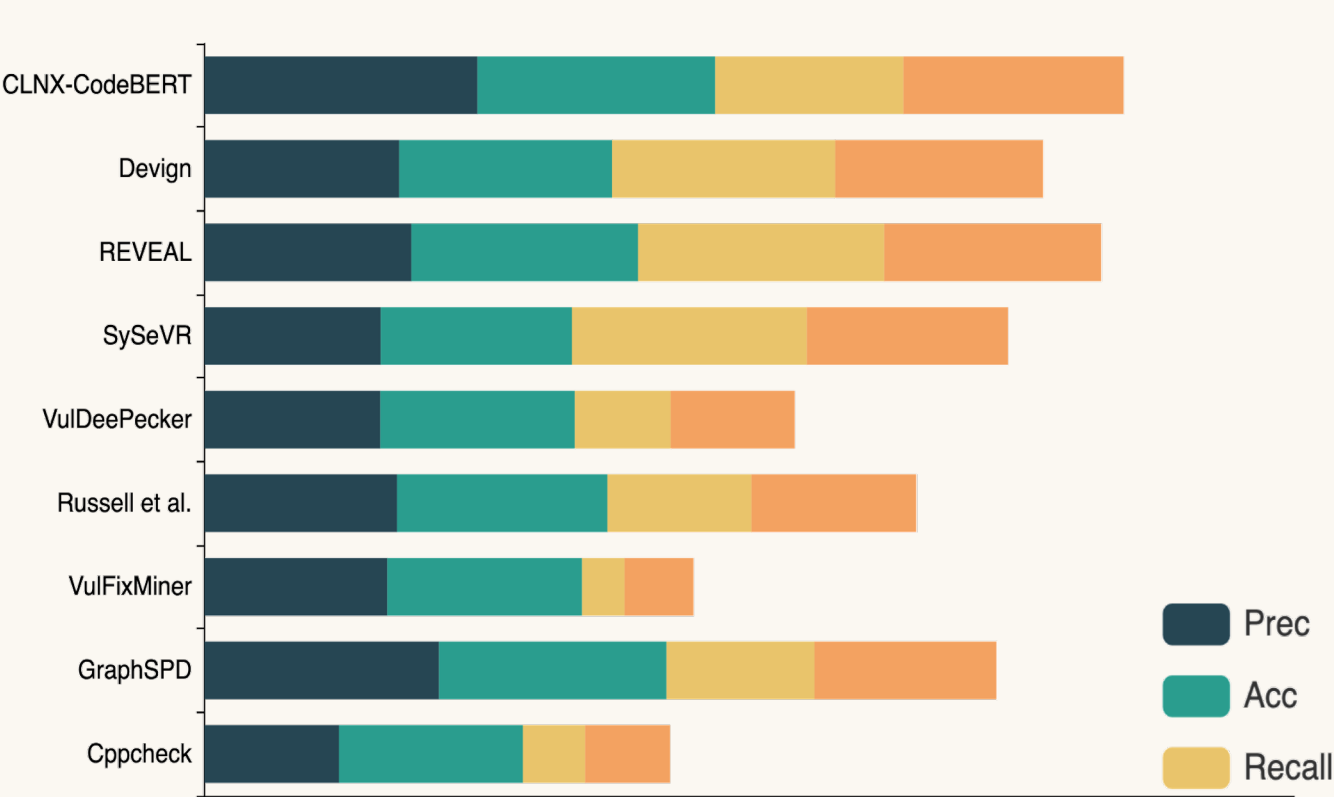C/C++ files + Known patch commits → CLNX → Fine-Tuned Models → Result

### OVERVIEW

*CLNX is a middleware designed to make C/C++ programs compatible with LLMs (especially BERT-based), thereby improving their ability to identify C/C++ VCCs. CLNX operates with minimal resource overhead, offering efficiency advantages over pre-training methods.*
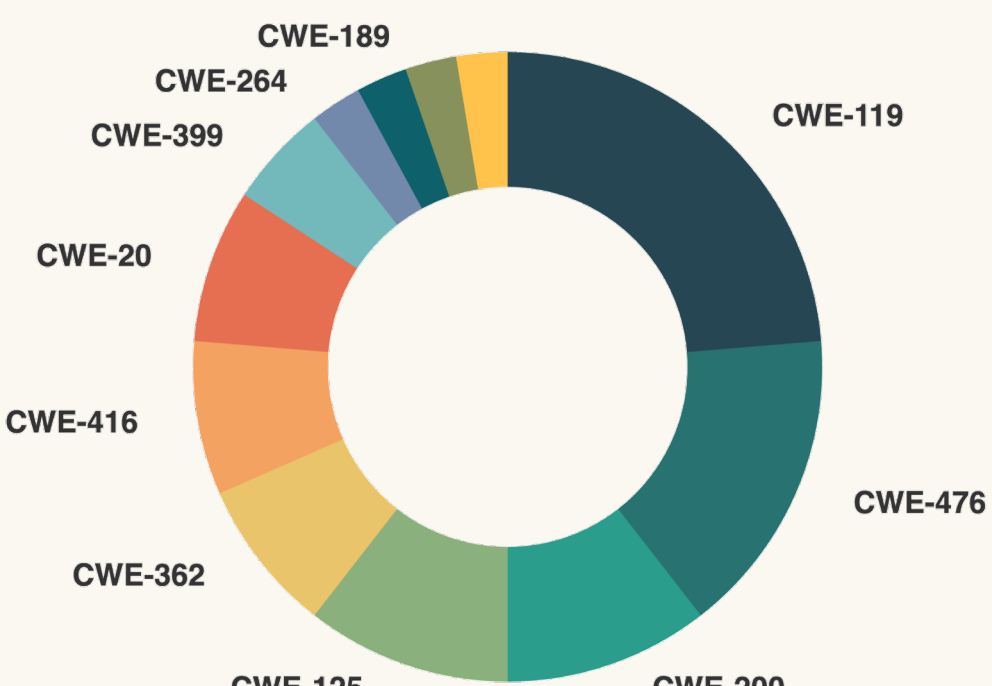
## EXPERIMENTS
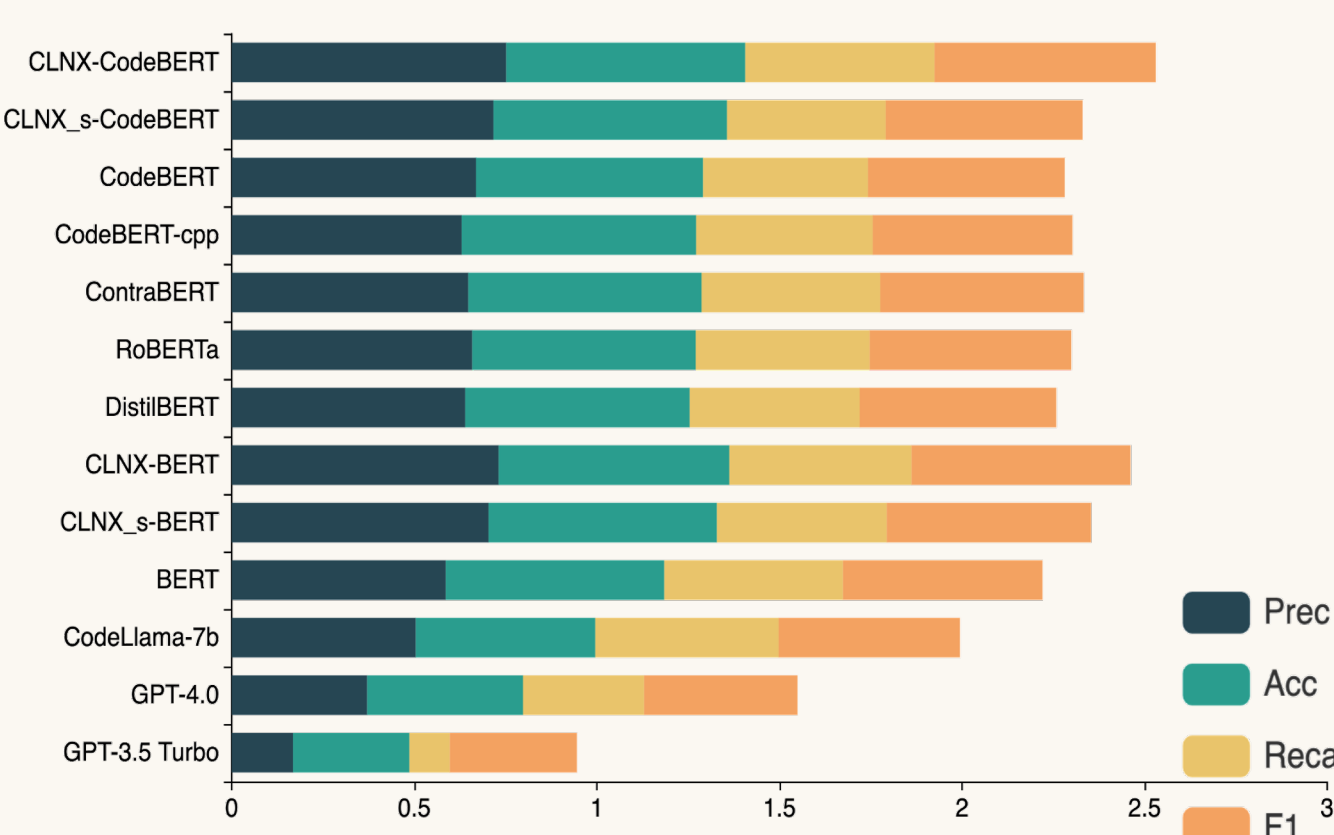
Results of comparative analysis



CLNX-CodeBERT, Devign, REVEAL, SySeVR, VulDeePecker, Russell et al., VulFixMiner, GraphSPD, Cppcheck

Legend: Prec, Acc, Recall, F1

| Items of Study |
|---|
| Effectiveness |
| Comparison |
| Read-World Application |

Results of LLMs in C/C++ VCCs identification



CLNX-CodeBERT, CLNX_s-CodeBERT, CodeBERT, CodeBERT-cpp, ContraBERT, RoBERTa, DistilBERT, CLNX-BERT, CLNX_s-BERT, BERT, CodeLlama-7b, GPT-4.0, GPT-3.5 Turbo

Legend: Prec, Acc, Recall, F1

Statistical Distribution of CWE Types



CWE-189, CWE-264, CWE-399, CWE-20, CWE-416, CWE-362, CWE-125, CWE-200, CWE-476, CWE-119

### EFFECTIVENESS: HOW DOES CLNX ENHANCE LLMs FOR C/C++ VCCs IDENTIFICATION TASK?

Both the Structure-level and Token-level Naturalization phases play crucial roles in CLNX's effectiveness. CLNX boosts LLMs' performance across all metrics (Precision: BERT: +14.48%, CodeBERT: +8.27%). CLNX outperforms resource-heavy pre-training methods. Despite minor recall dips in the early stages, full CLNX integration achieves the highest recall by retaining critical vulnerability patterns.

### COMPARISON: HOW DOES THE PERFORMANCE OF CLNX-LLMs COMPARE TO OTHER METHODS?

CLNX-CodeBERT surpasses graph-based models and traditional tools with 10.59% higher precision and the best F1 score, proving its lightweight code simplification outperforms complex embeddings. While recall lags slightly, its balanced efficiency and accuracy make it the top choice for C/C++ VCCs identification.

### REAL WORLD: HOW DOES CLNX-LLM PERFORMS IN IDENTIFY REAL-WORLD VCCs?

CLNX-CodeBERT detected 38 verified vulnerabilities across 35 C/C++ projects, excelling at identifying common patterns like Null Pointer Dereference (6) and Buffer Errors (9) by stripping code noise. While effective for logic-driven flaws, it struggles with complex, patternless vulnerabilities (e.g., Cryptographic Issues) due to inconsistent code logic.

*Experimantal Settings: Built on the public Devign dataset, our evaluation uses 25,872 real-world C/C++ function pairs (vulnerable/non-vulnerable) from FFmpeg and Qemu, split into 80% training, 10% validation, and 10% testing to ensure robust model assessment. For the real world one, we deploy the finetuned CLNX-CodeBERT to scan the repositories of 35 C/C++ open-source projects. Notably, VCCs are relatively rare in open-source software.*