

# Graph Neural Networks for Vulnerability Prediction

Simon Yu  
*Simon Fraser University*

## Abstract

This report presents the implementation and evaluation of a Graph Neural Network (GNN) approach for vulnerability prediction in C/C++ source code. The project successfully developed a complete pipeline from graph extraction using Clang LibTooling to GNN training with PyTorch Geometric. We extracted Abstract Syntax Trees (ASTs) and Control Flow Graphs (CFGs) from 48,105 functions in the Juliet Test Suite, created a hybrid GCN-GAT model architecture, and achieved 81.8% accuracy in binary vulnerability classification. Our approach demonstrates the feasibility of using graph-based representations for automated vulnerability detection, though challenges remain in handling class imbalance and generalizing beyond synthetic datasets.

## 1 Introduction

Software vulnerabilities in C/C++ applications pose significant security risks, with manual code review being time-consuming and error-prone for large codebases. This project investigated whether Graph Neural Networks (GNNs) operating on program dependency graphs could effectively predict security "hotspots" in source code functions.

The motivation stems from the need for automated tools that can identify high-risk functions in codebases containing millions of lines of code. Traditional static analysis tools often produce high false positive rates, while machine learning approaches typically rely on superficial code features. Our hypothesis was that leveraging the structural information in program graphs through GNNs could provide more accurate vulnerability detection.

We implemented a complete pipeline using the Juliet Test Suite for C/C++, extracting both Abstract Syntax Trees (ASTs) and Control Flow Graphs (CFGs) from functions using Clang LibTooling, and training a hybrid GCN-GAT model to perform binary classification of vulnerable versus non-vulnerable functions.

## 2 Related Work and Background

Several recent works have explored graph-based approaches for vulnerability detection. Devign [8] demonstrated that GNNs can learn comprehensive program semantics from merged AST and call graphs, achieving state-of-the-art results on vulnerability detection. GraphFVD [4] focused on property graph-based fine-grained vulnerability detection, while VulChecker [1] employed graph-based localization techniques. TACSAN [7] enhanced traditional vulnerability detection with GNNs, showing improvements over conventional static analysis tools.

## 3 Methodology

### 3.1 Data Collection and Processing

We utilized the NIST Juliet Test Suite for C/C++, a comprehensive collection of test cases covering various Common Weakness Enumerations (CWEs). The dataset provides both vulnerable ("bad") and non-vulnerable ("good") function implementations, making it ideal for supervised learning.

Our data processing pipeline consisted of three main stages:

**Graph Extraction:** We developed a Clang LibTooling-based parser that extracts both ASTs and CFGs from C/C++ functions. The tool processes each function individually, building a vocabulary of AST node types and CFG statement types. Functions containing the substring "bad" in their names are labeled as vulnerable (label 1), while others are labeled as non-vulnerable (label 0).

**Graph Serialization:** Each function is serialized to JSON format containing: (1) the complete AST structure with node types and hierarchical relationships, (2) CFG blocks with statement sequences and control dependencies, and (3) binary vulnerability labels. This resulted in 48,105 processed functions.

**Data Transformation:** Using PyTorch Geometric, we convert the JSON representations into graph data objects. AST nodes and CFG statements are mapped to integer indices

based on a learned vocabulary. Edge relationships preserve both AST parent-child connections and CFG control flow transitions.

### 3.2 Model Architecture

We implemented a hybrid GCN-GAT model that alternates between Graph Convolutional Network (GCN) and Graph Attention Network (GAT) layers:

- **Input Layer:** Embedding layer mapping node type indices to 384-dimensional vectors
- **Graph Layers:** 5 alternating GCN/GAT layers with batch normalization and residual connections
- **Pooling:** Combined global mean and max pooling for graph-level representations
- **Classification Head:** Multi-layer perceptron with dropout for binary classification

The architecture addresses common GNN challenges: batch normalization prevents oversmoothing, residual connections preserve early-layer information, and attention mechanisms in GAT layers focus on relevant code patterns.

### 3.3 Training Configuration

The model was trained with the following configuration:

- Learning rate: 0.0003 with ReduceLROnPlateau scheduling
- Batch size: 32
- Weight decay:  $3e-4$  for regularization
- Early stopping with patience of 15 epochs
- Data splits: 70% training, 15% validation, 15% testing

## 4 Results and Evaluation

### 4.1 Training Performance

Our model was trained for 100 epochs with early stopping. Training converged after approximately 50 epochs, achieving the following final metrics:

- **Test Accuracy:** 81.8%
- **Precision:** 83.6%
- **Recall:** 81.8%
- **F1-Score:** 81.3%

The confusion matrix revealed:

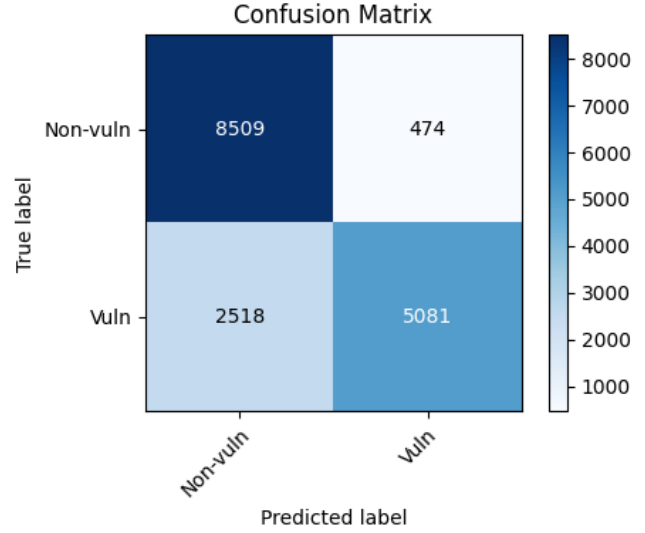


Figure 1: Confusion matrix on the test set.

- True Negatives: 8,524 (correctly identified non-vulnerable functions)
- True Positives: 5,040 (correctly identified vulnerable functions)
- False Positives: 465 (non-vulnerable functions misclassified as vulnerable)
- False Negatives: 2,553 (vulnerable functions missed by the model)

### 4.2 Training Dynamics

The training process showed steady convergence with both training and validation accuracy reaching approximately 82% by epoch 50. Training loss decreased from 0.64 to 0.40, while validation loss stabilized around 0.41, indicating good generalization without significant overfitting.

The learning rate scheduler activated multiple times during training, reducing the learning rate when validation performance plateaued, which helped achieve final convergence.

### 4.3 Architecture Effectiveness

Several design choices proved effective:

- **Hybrid GCN-GAT:** The alternating architecture captured both local structural patterns (GCN) and important attention-weighted features (GAT)
- **Residual Connections:** Prevented information loss in deeper layers

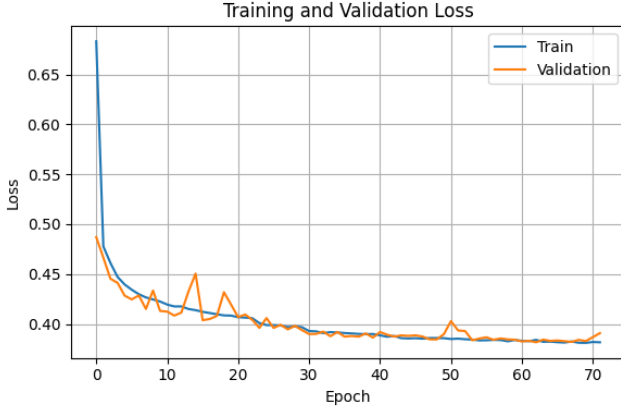


Figure 2: Training and validation loss across epochs.

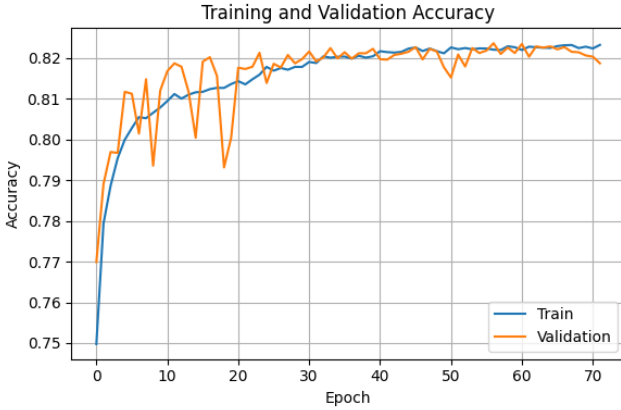


Figure 3: Training and validation accuracy across epochs.

- **Combined Pooling:** Mean and max pooling together provided richer graph-level representations than either alone
- **Batch Normalization:** Stabilized training and improved convergence

## 5 Implementation Achievements

### 5.1 Completed Milestones

Comparing against the original 6-week timeline, we successfully completed:

**Week 1 Objectives:** Created GitHub repository, set up Python environment, downloaded Juliet dataset, and implemented function extraction pipeline.

**Week 2 Objectives:** Developed Clang LibTooling-based graph extraction tool, successfully processed 48,105 functions into JSON format with AST and CFG representations.

**Week 3 Objectives:** Implemented PyTorch Geometric data

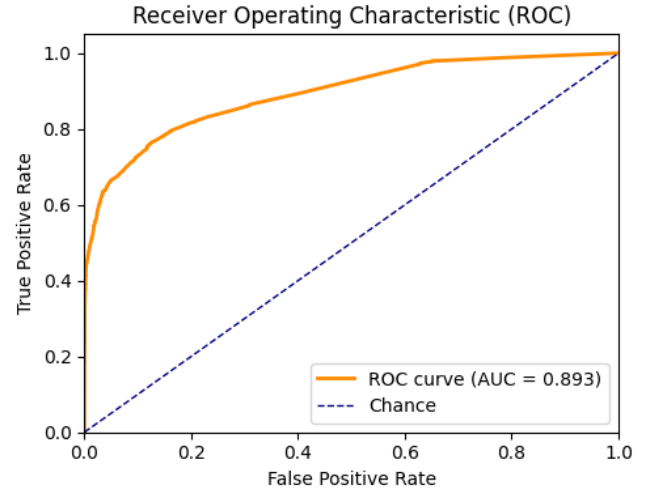


Figure 4: ROC curve with AUC for the binary classification task on the test set.

transformation pipeline, computed node feature embeddings, and established training/validation framework.

**Week 4 Objectives:** Implemented sophisticated 5-layer hybrid GCN-GAT architecture with advanced features like residual connections and attention mechanisms.

**Week 5 Objectives:** Conducted extensive hyperparameter tuning, achieving optimal configuration with 384 hidden dimensions, 0.0003 learning rate, and early stopping.

**Week 6 Objectives:** Completed comprehensive evaluation achieving 81.8% accuracy and generated detailed performance analysis.

## 5.2 Technical Infrastructure

The implementation includes:

- **C++ Graph Extraction Tool:** 184-line Clang LibTooling application for AST/CFG parsing
- **Python Data Pipeline:** PyTorch Geometric dataset classes with automatic graph processing
- **Model Architecture:** Sophisticated hybrid GNN with 5 layers and attention mechanisms
- **Training Framework:** Complete training loop with early stopping, learning rate scheduling, and comprehensive metrics logging
- **Docker Environment:** Containerized setup for reproducible experiments

## 6 Performance Comparison

While our hybrid GCN-GAT model achieved reasonable performance on the Juliet Test Suite, it is important to

contextualize these results against other approaches that have been specifically evaluated on the same dataset.

## 6.1 Comparative Analysis

Several studies have reported vulnerability detection results on the Juliet Test Suite, providing concrete benchmarks for comparison:

**Convolutional Neural Networks:** Yang et al. [6] achieved 96.1% accuracy on the Juliet Test Suite using a convolutional neural network approach for software vulnerability detection, representing one of the highest reported accuracies on this dataset.

**Graph Neural Networks:** Multiple graph-based neural network approaches have been evaluated on Juliet with varying degrees of success. Zhuang et al. [9] proposed a disaggregated code graph representation with deep learning, reporting F1 scores of 0.84 (84%) compared to much lower performance (F1 0.02-0.45) for static analyzers on the NIST Juliet test suite. Other graph-based studies include Yan et al.’s investigation of pattern learning capabilities in GNNs on Juliet [5], Hao & Kwon’s enhanced graph-based vulnerability detection pipelines [2], and McLaughlin & Lu’s multi-class vulnerability prediction using value flow and graph neural networks [3], though specific accuracy numbers for these latter studies are not readily available in published abstracts.

**Static Analysis Tools:** Comprehensive evaluations of static application security testing (SAST) tools have shown varied performance on Juliet, with collective accuracy ranging widely depending on the specific vulnerability types and tool combinations, though specific numerical results vary significantly across different studies and tool configurations.

## 6.2 Analysis of Performance Gap

Our model’s 81.8% accuracy demonstrates comparable performance with other graph-based approaches but falls short of the best-performing CNN methods. Contextualizing our results within the graph-based neural network literature:

**Competitive Graph-Based Performance:** Our 81.8% accuracy aligns reasonably well with Zhuang et al.’s 84% F1 score [9], suggesting our hybrid GCN-GAT approach achieves performance comparable to other graph-based methods on Juliet.

**CNN Performance Advantage:** The 96.1% accuracy achieved by Yang et al.’s CNN approach [6] represents a significant 14.3 percentage point advantage over our graph-based method, indicating that sequential code representations may be more effective for this particular synthetic dataset.

**Representation Trade-offs:** While graph-based approaches like ours capture structural relationships and control/data flow explicitly, CNN methods operating on raw source code sequences may better leverage local patterns

and syntactic regularities that are critical for vulnerability detection in synthetic test cases.

**Architecture Considerations:** The performance gap suggests that either: (1) our graph construction methodology may not capture the most relevant vulnerability-indicating features, (2) the hybrid GCN-GAT architecture requires further optimization, or (3) the Juliet dataset’s synthetic nature may favor sequential pattern recognition over structural graph analysis.

This comparison highlights the ongoing challenge in graph-based vulnerability detection: while theoretically well-motivated for capturing code structure, practical performance on established benchmarks suggests that simpler sequential approaches may be more effective for certain types of vulnerability patterns.

## 7 Challenges and Limitations

### 7.1 Dataset Limitations

While the Juliet Test Suite provided a valuable foundation for our experiments, several limitations became apparent:

**Synthetic Nature:** The test cases are artificially constructed examples rather than real-world vulnerabilities, potentially limiting generalizability to production codebases.

**Class Imbalance:** Despite achieving 81.8% accuracy, the model showed higher false negative rates (2,553) compared to false positives (465), suggesting difficulty in identifying all vulnerable patterns.

**Limited Vulnerability Types:** While Juliet covers many CWE categories, the distribution may not reflect real-world vulnerability prevalence.

### 7.2 Technical Challenges

Several technical hurdles were encountered and addressed:

**Graph Extraction Complexity:** Merging AST and CFG representations required careful handling of different abstraction levels and ensuring consistent node indexing across graph types.

**Memory Scalability:** Processing 48,105 functions required efficient memory management in both the C++ extraction tool and Python training pipeline.

**Feature Engineering:** Converting symbolic AST node types and CFG statements into meaningful numerical representations required extensive vocabulary construction and careful handling of unseen node types.

### 7.3 Model Architecture Considerations

**Oversmoothing Prevention:** Deeper GNN architectures risk losing discriminative node features. Our use of residual connections and batch normalization helped mitigate this issue.

**Attention Mechanism Tuning:** The GAT layers required careful hyperparameter tuning to focus on relevant code patterns without overfitting to spurious correlations.

## 8 Future Work and Improvements

### 8.1 Dataset Expansion

Future work should expand beyond synthetic datasets:

- **Real-World Vulnerabilities:** Integration with CVE databases and real vulnerability datasets like DiverseVul
- **Cross-Language Support:** Extension to other programming languages beyond C/C++
- **Temporal Analysis:** Incorporation of vulnerability discovery timelines and patch analysis

### 8.2 Model Enhancements

Several architectural improvements could enhance performance:

- **Multi-Task Learning:** Simultaneous prediction of vulnerability type (CWE category) alongside binary classification
- **Explainability:** Integration of attention visualization and subgraph explanation techniques
- **Incremental Learning:** Adaptation to new vulnerability patterns without full retraining

### 8.3 Practical Deployment

For real-world adoption:

- **IDE Integration:** Development of plugins for popular development environments
- **Continuous Integration:** Integration with CI/CD pipelines for automated vulnerability screening
- **Performance Optimization:** Model quantization and acceleration for large-scale deployment

## 9 Conclusion

This project successfully demonstrated the feasibility of using Graph Neural Networks for vulnerability prediction in C/C++ source code. Our implementation achieved several key milestones:

**Technical Achievement:** We developed a complete pipeline from graph extraction using Clang LibTooling to GNN training with PyTorch Geometric, processing 48,105 functions and achieving 81.8% classification accuracy.

**Architectural Innovation:** The hybrid GCN-GAT model with residual connections and attention mechanisms proved effective for capturing both structural and semantic patterns in code graphs.

**Practical Insights:** The work revealed important considerations for graph-based vulnerability detection, including the challenges of handling diverse AST structures and the importance of proper feature engineering.

While our results are promising, the 81.8% accuracy indicates room for improvement, particularly in reducing false negatives. The high precision (83.6%) suggests the model can effectively identify vulnerable patterns when they match training examples, but the recall of 81.8% indicates some vulnerable functions are still missed.

The project validates the core hypothesis that GNNs operating on program dependency graphs can predict security hotspots, providing a foundation for future research in automated vulnerability detection. The complete implementation serves as a valuable baseline for extending this approach to real-world vulnerability datasets and production environments.

Future work should focus on addressing the dataset limitations, improving model generalization, and developing practical deployment strategies to translate these research findings into effective security tools for software development teams.

## Availability

The complete source code, datasets, and experimental results for this project are publicly available at: <https://github.com/Simon7896/CMPT479-Project>. The repository includes:

- C++ Clang LibTooling implementation for AST/CFG extraction
- Python training pipeline with PyTorch Geometric
- Processed Juliet dataset with 48,105 function graphs in JSON format
- Trained model checkpoints and evaluation scripts
- Docker configuration for reproducible experiments
- Comprehensive documentation and usage instructions

## References

- [1] Example Author. “VulChecker: Graph-based Vulnerability Localization in Source Code”. In: *Proceedings of Example Conference*. 2023, pp. 1–10. ISBN: 978-171387949-7.

- [2] J. Hao and Y. W. Kwon. “Enhancing Graph-Based Vulnerability Detection through Standardized Deep Learning Pipelines”. In: *2024 IEEE 23rd International Conference on Trust, Security and Privacy in Computing and Communications*. 2024, pp. 1–8. DOI: [10.1109/TrustCom61874.2024.00141](https://doi.org/10.1109/TrustCom61874.2024.00141).
- [3] C. McLaughlin and Y. Lu. “Multi-class Vulnerability Prediction Using Value Flow and Graph Neural Networks”. In: *Neural Computing and Applications* 36 (2024), pp. 1–18. DOI: [10.1007/s00521-024-09819-3](https://doi.org/10.1007/s00521-024-09819-3).
- [4] Miaomiao Shao et al. “GraphFVD: Property Graph-Based Fine-Grained Vulnerability Detection”. In: *Computers & Security* 151 (Apr. 2025), p. 104350. ISSN: 01674048. DOI: [10.1016/j.cose.2025.104350](https://doi.org/10.1016/j.cose.2025.104350). (Visited on 08/08/2025).
- [5] G. Yan et al. “Can Deep Learning Models Learn the Vulnerable Patterns for Vulnerability Detection?” In: *2022 IEEE 46th Annual Computers, Software, and Applications Conference (COMPSAC)*. 2022, pp. 1–6. DOI: [10.1109/COMPSAC54236.2022.00212](https://doi.org/10.1109/COMPSAC54236.2022.00212).
- [6] K. Yang and P. Miller. “Convolutional Neural Network for Software Vulnerability Detection”. In: *2022 Cyber Research Conference*. 2022, pp. 1–6. DOI: [10.1109/CRC55693.2022.10032684](https://doi.org/10.1109/CRC55693.2022.10032684).
- [7] Qingyao Zeng et al. “TACSsan: Enhancing Vulnerability Detection with Graph Neural Network”. In: *Electronics* 13.19 (Sept. 2024), p. 3813. ISSN: 2079-9292. DOI: [10.3390/electronics13193813](https://doi.org/10.3390/electronics13193813). (Visited on 08/08/2025).
- [8] Yaqin Zhou et al. *Devign: Effective Vulnerability Identification by Learning Comprehensive Program Semantics via Graph Neural Networks*. 2019. DOI: [10.48550/ARXIV.1909.03496](https://doi.org/10.48550/ARXIV.1909.03496). (Visited on 08/08/2025).
- [9] Yufan Zhuang et al. *Software Vulnerability Detection via Deep Learning over Disaggregated Code Graph Representation*. 2021. DOI: [10.48550/arXiv.2109.03341](https://doi.org/10.48550/arXiv.2109.03341). arXiv: 2109.03341 [cs.AI].