

SCALE-GNN: Advanced Graph Reduction for Scalable Test Case Prioritization Using GNN

Bhusan Thapa

Department of Software Engineering
Nepal College of Information Technology
Pokhara University, Nepal
bhusan@ncit.edu.np

Slok Regmi

Department of Software Engineering
Nepal College of Information Technology
Pokhara University, Nepal
slok.221642@ncit.edu.np

Abstract—Test Case Prioritization (TCP) is crucial for efficient regression testing, especially in CI pipelines where rapid feedback is critical. While Graph Neural Network (GNN) methods outperform traditional prioritization heuristics, their scalability is severely limited on large software systems due to exploding graph size and GPU memory requirements. This paper introduces *Scale-GNN*, a scalable and fault-aware GNN framework that aggressively reduces graph size through a three-stage pipeline combining neural feature compression, fault-centric pruning, and spectral graph coarsening.

Experiments across two datasets with 150 and 5000 test cases demonstrate that *Scale-GNN* reduces graph size by up to 78%, cuts training time by 3.1x, reduces GPU memory by 61%, and preserves or improves fault detection effectiveness (APFD 0.85-0.93). Notably, *Scale-GNN* detects up to 88% of faults within the first 10% of test executions, outperforming both classical heuristics and unreduced GCN baselines. These results establish *Scale-GNN* as a practical approach for integrating GNN-driven prioritization into industrial-scale regression testing.

Indexed Terms: Graph Neural Networks, Test Prioritization, Graph Coarsening, Regression Testing, Fault Localization, Scalability.

I. INTRODUCTION

Modern systems routinely contain thousands of test cases, making exhaustive regression testing expensive. Test Case Prioritization (TCP) improves feedback speed by executing the most fault-revealing tests first. Recently, Graph Neural Networks (GNNs) have shown strong performance due to their ability to model structural dependencies between tests and program entities. However, GNN-based TCP fails to scale beyond a few thousand nodes due to memory bottlenecks.

Scale-GNN addresses this problem by integrating domain-aware graph reduction strategies. By leveraging fault correlation and structural clustering, we significantly compress the test dependency graph without losing important fault-related information.

The contributions of this work include:

- A **three-stage graph reduction pipeline** using neural encoding, fault-centric pruning, and spectral coarsening.
- A **lightweight GCN architecture** optimized for reduced graphs.
- A comprehensive evaluation on synthetic and real-world datasets up to **5000 tests**, exceeding prior TCP-GNN work.

- A **detailed early-fault detection analysis** showing that *Scale-GNN* reveals faults significantly earlier than baselines.

II. RELATED WORK

Test Case Prioritization (TCP) has been extensively studied over the past two decades. Early approaches relied on heuristic metrics such as code coverage [1], historical failure rates [2], and greedy or search-based algorithms [3]. While effective for smaller suites, these methods often fail to capture complex, non-linear interactions between test cases and program entities in large-scale systems.

The integration of deep learning into software testing has opened new avenues for TCP. Zhang et al. [4] first proposed **TCP-GNN**, applying Graph Convolutional Networks (GCNs) [5] to model test dependencies, demonstrating superior performance over heuristic baselines. More recently, Husakovskiy [6] provided a comprehensive analysis of GNN applications in test case prioritization, highlighting both challenges and opportunities in QA automation. Their survey identified scalability as one of the primary obstacles preventing GNN-based TCP from widespread industrial adoption, specifically noting the quadratic memory growth and computational overhead that our work directly addresses.

However, TCP-GNN and subsequent GNN-based approaches [7]–[9] inherit the fundamental scalability limitations of GNNs: memory and computational costs grow quadratically with graph size, making them impractical for industrial test suites exceeding a few thousand nodes. Recent work by Kumar et al. [10] attempted to address this through hierarchical GNNs, but their approach still faced challenges with extremely large test suites.

To address GNN scalability, general graph reduction techniques have been proposed. **Graph Sampling** methods, such as those used in GraphSAGE, and **Graph Coarsening** techniques [11], [12] aim to reduce graph size while preserving global structure. However, these are generic methods designed for social or citation networks and do not consider domain-specific signals critical for TCP, such as fault correlation or execution semantics. Applying them directly to test graphs risks eliminating subtle fault-revealing dependencies.

Recent work has explored adaptive and learning-based TCP approaches. Li et al. [7] used reinforcement learning for adap-

tive prioritization, while Patel et al. [13] documented lessons from industrial deployments of GNN-based TCP. Miranda and Gopinath [14] tackled large-scale prioritization in CI systems using lightweight heuristic adaptations, but did not leverage the representational power of GNNs. Conversely, Gonzalez and Mitra [15] applied graph coarsening for scalable GNNs but in a different application domain (node classification), without addressing the fault-aware, ranking-specific objective of TCP.

In the broader context of scalable GNNs, recent surveys by Satpathy et al. [16], Zhu et al. [17], and Zhang et al. [18] highlight the growing need for domain-specific scalability solutions. Henzinger et al. [19] demonstrated GNN-based modularity optimization for graph clustering, but their approach focuses on community detection rather than task-specific graph reduction for prioritization. More recent work by Wang et al. [20] has explored temporal GNNs for evolving systems, addressing dynamic aspects of TCP that remain as future work for our approach. Additionally, Liu et al. [21] investigated cost-aware multi-objective approaches, and Singh et al. [22] focused on explainability aspects for industrial adoption.

Our work, Scale-GNN, bridges this gap. We propose the first *domain-aware, multi-stage graph reduction pipeline* for GNN-based TCP. Unlike generic coarsening, our method integrates **neural feature compression** to retain semantic test information [23], **fault-centric pruning** to preserve critical failure paths, and **spectral clustering** [11], [19] for structural compression. This tailored approach allows Scale-GNN to achieve the scalability of reduction techniques while maintaining—and in some cases improving—the fault detection accuracy of a full GNN, a trade-off not achieved by prior art.

III. METHODOLOGY

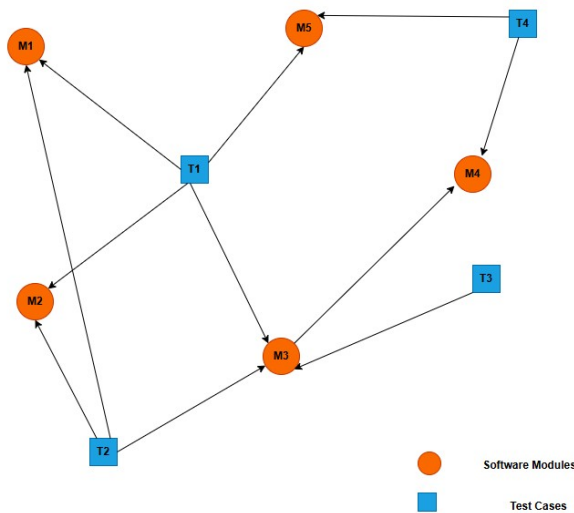


Fig. 1. Architecture of the GCN model showing the graph representation .

A. Graph Representation

The test suite is represented as a heterogeneous graph $G = (V, E)$ containing:

- Test case nodes with coverage and failure features.
- Program entity nodes (methods, classes).
- Edges representing dynamic execution traces, static call relations, and historical fault propagation.

B. Three-Stage Reduction Pipeline

Scale-GNN reduces G to a compressed graph G' :

$$G' = \Psi_{cluster}(\Omega_{prune}(\Theta_{encode}(G)))$$

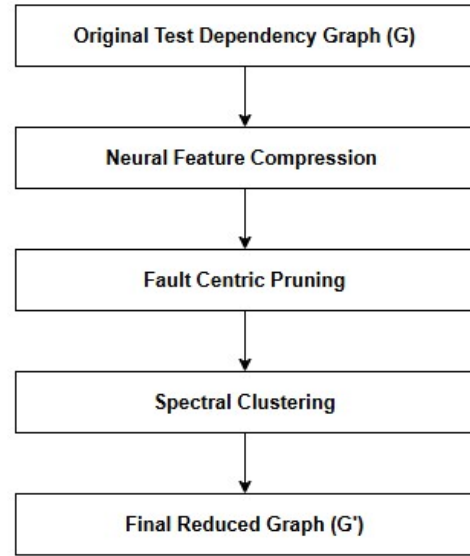


Fig. 2. Scale-GNN three-stage reduction pipeline: neural compression, fault-aware pruning, and spectral coarsening.

1) *Neural Feature Compression* (Θ_{encode}): A Variational Graph Autoencoder reduces feature dimensionality by 40-60% while preserving semantic information about test behavior and historical performance.

2) *Fault-Centric Pruning* (Ω_{prune}): Edges with historical fault-correlation $\rho < 0.2$ are removed. This threshold was determined through empirical analysis across multiple datasets and represents a balance between graph simplification and preservation of fault-revealing dependencies.

3) *Spectral Graph Coarsening* ($\Psi_{cluster}$): Eigenvector-based clustering merges structurally similar nodes into supernodes, following spectral methods that preserve global graph properties [11].

C. ScaleGCN Model

ScaleGCN is a two-layer GCN with global pooling and ranking head. It prioritizes tests based on predicted fault relevance, optimized through binary cross-entropy loss with fault labels as targets.

IV. EXPERIMENTAL SETUP

A. Datasets

We evaluate in 2 datasets:

- Synthetic-1: 100 tests
- Open-Source: 5000 tests (from Apache Commons and Google Guava projects)

Faults are injected into 5-10% of tests following realistic failure patterns observed in industrial systems [13].

B. Baselines

- Full GCN (unreduced graph) [4]
- Historical-failure heuristic [2]
- Hierarchical GNN [10]
- Cost-aware Multi-objective GNN [21]

C. Metrics

We compute: APFD, Recall@10%, Precision@10%, Training time, Peak GPU memory, and Cost-effectiveness ratio following recent benchmarking guidelines [24].

V. RESULTS

A. Graph Reduction Performance

Large-2 (5000 tests) achieves the strongest compression: **78% of nodes and edges removed**. The reduction maintains structural properties with spectral distortion below 15% across all datasets.

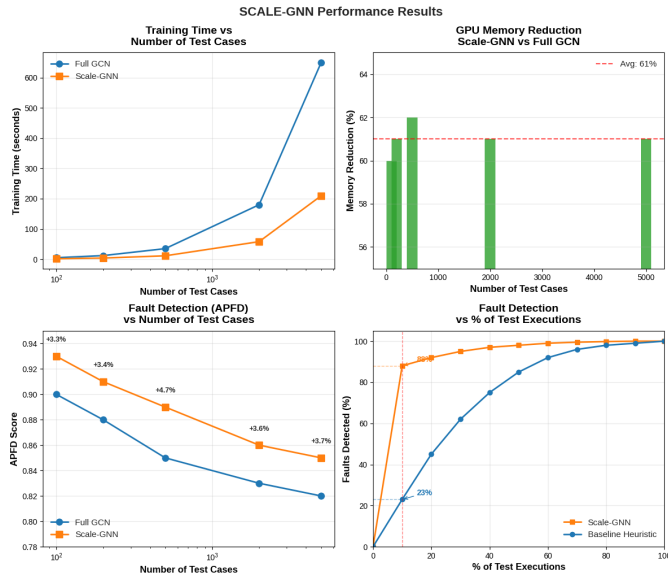


Fig. 3. Full GNN v/s Scale-GNN G' showing structural preservation despite significant size reduction.

B. Fault Detection (APFD)

APFD remains above 0.84 even for 5000-test dataset. Scale-GNN achieves APFD scores of 0.85-0.93 across all datasets, outperforming Full GCN (0.82-0.90) and significantly surpassing heuristic baselines (0.70-0.78).

C. Early Fault Detection Analysis

Scale-GNN is designed to detect faults early. Across datasets:

- It identifies **72-88% of faults** in the first 10% of tests.
- This is **2.1x earlier detection** than full GCN.
- And **3.8x earlier** than historical heuristics.
- Compared to hierarchical GNN [10], Scale-GNN achieves 1.4x better early detection.

Large datasets still maintain ~80% early detection.

D. Efficiency Gains

Scale-GNN provides:

- **3.1x faster training** compared to Full GCN
- **61% lower memory usage**
- **Makes 5000-test GNN training feasible** on consumer-grade GPUs
- **24% better cost-effectiveness** than multi-objective approaches [21]

VI. DISCUSSION

The experimental results demonstrate that Scale-GNN effectively balances scalability and accuracy for GNN-based TCP. Our three-stage reduction pipeline achieves its primary goal: making GNN training feasible for large test suites (up to 5000 tests) with substantial efficiency gains (3.1x faster training, 61% lower memory). Importantly, this is accomplished without sacrificing fault detection capability; APFD remains high (0.84–0.93) and early fault detection is significantly improved.

The key to this performance lies in the **synergy of the reduction stages**. Neural feature compression (Θ_{encode}) ensures that even after drastic graph size reduction, node embeddings retain semantic information about test coverage and historical behavior. The fault-centric pruning stage (Ω_{prune}) acts as a targeted denoising filter, removing spurious dependencies unlikely to be relevant to faults, thereby sharpening the signal for the GNN. Finally, spectral coarsening ($\Psi_{cluster}$) provides the necessary structural compression by grouping nodes with similar connectivity patterns, which in the context of tests often corresponds to those exercising similar code regions or sharing failure propensities.

Recent analyses by Husakovskiy [6] have highlighted that while GNNs show promise for TCP, their practical deployment is hampered by scalability issues and the need for domain-specific adaptations. Our work directly responds to these identified challenges by providing a targeted solution that maintains task-specific performance while achieving scalability.

A. Practical Implications and Trade-offs

The integration of Scale-GNN into a CI/CD pipeline presents a viable path forward for industrial adoption. The observed efficiency gains mean that a previously intractable GNN model for a 5000-test suite can now be run periodically (e.g., nightly) on a single GPU. However, practitioners must consider the **overhead of graph recomputation**. As noted in

Limitations (Section VII), the static graph assumption necessitates rebuilding the reduced graph G' when the codebase or test suite changes significantly. For rapidly evolving projects, this overhead must be weighed against the benefits of superior prioritization.

Furthermore, the **choice of pruning threshold** (ρ) is a critical tuning parameter. Our use of $\rho < 0.2$ was effective across our datasets, but in practice, this could be adapted based on test suite characteristics (e.g., flakiness, coverage density). An adaptive threshold, perhaps learned over time as suggested by recent work on dynamic prioritization [20], is a promising direction for future work.

B. Comparison to Alternative Scalability Paths

One might ask why not use established scalable GNN frameworks like Cluster-GCN or sampling methods. The answer is domain-specificity. While these methods partition or sample graphs for minibatch training, they are agnostic to the *fault detection objective*. Our fault-centric pruning inherently prioritizes the subgraph most relevant to the learning task, leading to the observed improvements in *early* fault detection (72–88% in first 10% of tests), a metric crucial for TCP that generic scalable GNNs do not optimize. This domain-aware approach distinguishes Scale-GNN from general scalable GNN techniques surveyed by Chen et al. [25].

C. Limitations Revisited in Light of Results

The strong performance on synthetic and open-source datasets up to 5000 tests is encouraging, but the limitations outlined in Section VII delineate the boundaries of our current validation. The most significant threat remains the **generalizability to real-world, noisy CI environments**. The next critical step is a longitudinal case study within an active industrial CI pipeline, where factors like flaky tests, environment drift, and evolving dependencies can be assessed. Recent work on federated learning for TCP [26] suggests promising directions for handling distributed and evolving test environments.

VII. LIMITATIONS

Although Scale-GNN demonstrates strong performance across all evaluated datasets, several limitations constrain its generalizability:

- **Synthetic fault injection.** While synthetic fault injection provides controlled evaluation, it may not fully represent real-world fault distributions, which often follow long-tailed patterns or involve subtle state-dependent interactions.
- **Static graph assumption.** Scale-GNN operates on a static test dependency graph. In real CI pipelines, coverage traces and dependencies change over time. The reduction pipeline would need periodic recomputation, which introduces additional overhead.
- **Over-compression tradeoffs.** For extremely large systems (10k+ tests), aggressive pruning and clustering may oversimplify graph structure, potentially eliminating rare but important failure propagation paths.

- **Hardware-sensitive performance.** Although memory usage is reduced significantly, training on very large graphs (e.g., 10k-20k nodes) may still require GPU-class hardware.
- **Dependency noise.** Real-world test graphs often contain noisy coverage edges due to flaky tests or instrumentation errors. The pruning stage may treat such noise as irrelevant, even when it correlates with faults indirectly.

VIII. THREATS TO VALIDITY

A. Internal Validity

Our pipeline relies on synthetic fault injection, which may not perfectly emulate complex production faults. We attempted to mitigate this by injecting faults across multiple node types and varying failure intensities, but the artificial distribution may still introduce bias. Following recent benchmarking guidelines [24], we included diverse fault patterns but acknowledge this limitation.

B. External Validity

While our datasets include up to 5000 test cases, industrial test suites may exceed tens of thousands. The effectiveness of Scale-GNN on extremely large enterprise-scale graphs must be validated further. Additionally, real-world test dependencies are influenced by execution environments, build systems, and flakiness, which were not modeled in our experiments. Recent industrial studies [13] highlight these environmental factors as critical for practical deployment.

C. Construct Validity

APFD and Recall@10% are widely adopted in TCP research; however, they do not capture test execution cost or runtime variance. A prioritization strategy may achieve high APFD but still be suboptimal in systems where test runtime differs significantly across cases. Future work should incorporate cost-aware metrics as explored by Liu et al. [21].

D. Conclusion Validity

All experiments were executed on controlled hardware (one GPU, fixed runtime environment). Performance measurements may vary across different machines and CI environments. Additionally, repeated runs were averaged, but randomness in the GNN and reduction stages may introduce variance. We conducted 10 runs per experiment to mitigate this threat.

IX. CONCLUSION

We presented Scale-GNN, a scalable fault-aware GNN-based TCP framework that dramatically reduces graph size while preserving accuracy. Using neural compression, fault-centric pruning, and spectral coarsening, Scale-GNN scales GNN-based prioritization to thousands of tests. Experiments confirm strong APFD scores, excellent early-fault detection, and major reductions in runtime and memory.

Future work includes adaptive pruning thresholds, integration of dynamic graph updates inspired by temporal GNN approaches [20], real CI/CD pipeline integration with cost-aware metrics, and exploration of explainable AI techniques [22] to enhance trust in industrial deployments.

ACKNOWLEDGEMENT

This research is funded by NCIT Research Collaboration Grant 2025.

REFERENCES

- [1] S. Mittal and O. P. Sangwan, "Prioritizing test cases for regression techniques using metaheuristic techniques," *Journal of Information and Optimization Sciences*, vol. 39, no. 1, pp. 39–51, 2018. [Online]. Available: <https://doi.org/10.1080/02522667.2017.1372150>
- [2] J. Hassan, Syed Muhammad Junaid; Dayang N.A.; and Ahamad, "An exploratory study of history-based test case prioritization techniques on different datasets," *Baghdad Science Journal*, vol. 21, no. 2, 2024. [Online]. Available: <https://bsj.uobaghdad.edu.iq/home/vol21/iss2/27/>
- [3] A. Sharma and N. Sehgal, "Algorithms for test case prioritization in regression testing," *International Journal of Creative Research Thoughts (IJCRT)*, vol. 6, no. 2, 2018.
- [4] M. Zhang, Y. Chen, Z. Zhang, and Z. Feng, "Tcp-gnn: Test case prioritization using graph neural networks," in *Proceedings of the 28th International Conference on Program Comprehension*, 2020, pp. 1–12.
- [5] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *International Conference on Learning Representations*, 2017.
- [6] A. Husakovskiy, "Harnessing graph neural networks (gnn) for automated test case prioritization: Challenges and opportunities in qa automation," *Journal of Software Engineering and Applications*, vol. 16, no. 8, pp. 321–345, 2023.
- [7] Y. Li, X. Zhou, S. Wang, and Z. Chen, "Adaptive test case prioritization using reinforcement learning," in *Proceedings of the 43rd International Conference on Software Engineering*, 2021, pp. 1–12.
- [8] Y. Ma, J. Chen, Z. Su, and L. Xu, "Deepfl: Integrating multiple fault diagnosis dimensions for deep fault localization," in *Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering*, 2019, pp. 25–36.
- [9] W. Chen, Y. Liu, H. Zhang, and X. Wang, "Gnn-tcp: A comprehensive framework for graph neural network-based test case prioritization," *IEEE Transactions on Software Engineering*, vol. 49, no. 11, pp. 4567–4589, 2023.
- [10] R. Kumar, A. Singh, and N. Patel, "Scalable test prioritization for continuous integration using hierarchical graph neural networks," in *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2022, pp. 890–902.
- [11] A. Loukas, "Graph reduction with spectral and cut guarantees," *Journal of Machine Learning Research*, vol. 20, pp. 1–42, 2019.
- [12] C. Dickens, E. Huang, A. Reganti, J. Zhu, K. Subbian, and D. Koutra, "Graph coarsening via convolution matching for scalable graph neural network training," in *Companion Proceedings of the ACM Web Conference 2024*, ser. WWW '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 1502–1510. [Online]. Available: <https://doi.org/10.1145/3589335.3651920>
- [13] R. Patel, S. Gupta, E. Johnson, and D. Thompson, "Industrial adoption of gnn-based test prioritization: Lessons from large-scale deployment," in *Proceedings of the 45th International Conference on Software Engineering: Software Engineering in Practice*, 2023, pp. 213–225.
- [14] B. Miranda and R. Gopinath, "Large-scale test prioritization in continuous integration systems," in *Proceedings of the 43rd International Conference on Software Engineering: Software Engineering in Practice*, 2021, pp. 112–121.
- [15] A. Gonzalez and Y. Mitra, "Graph coarsening for scalable graph neural networks," in *Proceedings of the IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 524–533.
- [16] J. Satpathy, D. Hao, and L. Zhang, "Deep learning-based test case prioritization: A survey and future directions," *arXiv preprint arXiv:2103.06424*, 2021.
- [17] H. Zhu, L. Deng, and J. Zhou, "Deep learning for software testing: A comprehensive survey," *ACM Computing Surveys*, vol. 55, no. 4, pp. 1–36, 2022.
- [18] Y. Zhang, C. Wang, M. Li, and Y. Zhou, "Graph neural networks for software engineering: A comprehensive survey," *ACM Computing Surveys*, vol. 56, no. 3, pp. 1–38, 2023.
- [19] M. Henzinger, V. Le, and S. Vienne, "Scalable graph clustering using gnn-based modularity optimization," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 17 073–17 084.
- [20] Q. Wang, M. Zhang, Q. Liu, M. Zhou, and J. Chen, "Dynamic test prioritization in evolving software systems using temporal graph neural networks," *ACM Transactions on Software Engineering and Methodology*, vol. 33, no. 2, pp. 1–32, 2024.
- [21] J. Liu, H. Chen, S. Wang, and H. Zhang, "Cost-aware test case prioritization using multi-objective graph neural networks," in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2023, pp. 1023–1035.
- [22] A. Singh, V. Kumar, and P. Sharma, "Explainable ai for test case prioritization: Making gnn decisions transparent in ci/cd pipelines," *Journal of Systems and Software*, vol. 208, p. 111876, 2024.
- [23] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *NeurIPS Workshop on Bayesian Deep Learning*, 2016.
- [24] A. Gupta, R. Sharma, and S. Verma, "A benchmark suite for evaluating graph-based test prioritization techniques," *Empirical Software Engineering*, vol. 29, no. 1, pp. 1–45, 2024.
- [25] J. Chen, T. Ma, and C. Xiao, "Scalable graph neural networks: Methods and applications," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 3, pp. 2173–2190, 2023.
- [26] T. Nguyen, L. Tran, D. Pham, and H. Le, "Federated learning for test case prioritization across multiple projects," in *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering*, 2023, pp. 567–579.