

# Filo-Priori: A Dual-Stream Deep Learning Approach to Test Case Prioritization

Acauan C. Ribeiro Instituto de Computação (IComp)  
 Universidade Federal do Amazonas (UFAM)  
 Manaus, AM, Brazil  
 acauan@icomp.ufam.edu.br

**Abstract**—Test Case Prioritization (TCP) aims to order test cases to maximize early fault detection in Continuous Integration (CI) environments. Existing approaches treat software versions as linear time series, failing to capture the complex relationships between test cases that characterize real-world testing. We propose a **dual-stream architecture** that combines semantic understanding of test cases with structural patterns learned from execution history.

We present FILO-PRIORI, a deep learning approach that introduces: (1) a **Multi-Edge Test Relationship Graph** encoding co-failure, co-success, and semantic similarity relationships between test cases; (2) a **Dual-Stream Architecture** combining semantic embeddings (SBERT) with Graph Attention Networks (GAT) for structural feature learning; (3) **Cross-Attention Fusion** for dynamic modality combination; and (4) **Weighted Focal Loss** for handling severe class imbalance (37:1).

We evaluate FILO-PRIORI on two complementary datasets: an industrial dataset (52,102 test executions, 277 builds with failures) and the RTPTorrent open-source benchmark (20 Java projects, 100k+ builds). On the industrial dataset, results show a mean APFD of 0.6413, representing 14.6% improvement over random ordering ( $p < 0.001$ ) and 2.0% over the strongest baseline. Ablation reveals Graph Attention contributes +17.0% to performance. Temporal validation confirms robustness (APFD: 0.619–0.663).

This work demonstrates the effectiveness of combining semantic and structural information for TCP through graph neural networks. Our replication package is publicly available.

**Index Terms**—Test Case Prioritization, Graph Attention Networks, Deep Learning, Continuous Integration, Dual-Stream Architecture, Class Imbalance

## 1 INTRODUCTION

CONTINUOUS Integration (CI) has become a fundamental practice in modern software development, enabling teams to integrate code changes frequently and detect defects early [1], [2]. A key challenge in CI environments is managing the growing test suite: as software evolves, the number of test cases increases, making it impractical to execute all tests for every commit [3].

Test Case Prioritization (TCP) addresses this challenge by ordering test cases to maximize early fault detection [4], [5]. The goal is to execute tests most likely to fail first, providing faster feedback to developers. The effectiveness of TCP is typically measured using the Average Percentage of Faults Detected (APFD) metric [6].

**The Challenge of Test Relationships.** Existing TCP approaches, whether based on coverage [4], historical failure [7], or machine learning [8], [9], often treat test cases as independent entities. This assumption ignores the rich relationships between test cases: tests that fail together often indicate related functionality, tests with similar descriptions target similar code, and historical patterns reveal systematic dependencies [10].

**A Dual-Stream Approach.** We propose combining two complementary information sources for TCP: (1) *semantic information* from test descriptions and commit messages, and (2) *structural information* from test execution history and test relationships. This dual-stream approach captures both

what tests do (semantics) and how they behave (structure).

Table 1 presents the two types of features we combine:

TABLE 1: Dual-Stream Feature Types

Semantic Features	Structural Features
Test case summary	Historical failure rate
Test case steps	Recent failure trend
Commit messages	Test age (builds since first run)
Code changes (diff)	Flakiness rate
–	Co-failure relationships
–	Consecutive failure streaks

The key insight is that semantic similarity alone is insufficient—tests with similar descriptions may have very different failure patterns. By combining semantic and structural information through Graph Attention Networks, we can learn which tests are likely to fail based on both their content and their history.

In this paper, we present FILO-PRIORI, a deep learning approach for Test Case Prioritization that combines semantic and structural information through four key innovations:

- 1) **Multi-Edge Test Relationship Graph:** We construct a graph that captures multiple types of relationships between test cases: co-failure edges (tests that fail together), co-success edges (tests that pass together), and semantic similarity edges (tests with similar descriptions). This multi-edge approach increases

- graph density and captures complementary relationships.
- 2) **Dual-Stream Architecture:** We employ a dual-stream model that processes semantic features (SBERT embeddings of test descriptions and commit messages) and structural features (historical execution patterns) through separate neural networks before fusion.
  - 3) **Graph Attention Networks for Structural Learning:** We use Graph Attention Networks (GAT) [11] to learn representations that capture test relationships, allowing the model to propagate information between related test cases with learned attention weights.
  - 4) **Weighted Focal Loss for Class Imbalance:** We employ Weighted Focal Loss [12] to address the severe class imbalance (37:1 Pass:Fail ratio) inherent in test execution data, focusing training on hard-to-classify examples while down-weighting easy negatives.

We evaluate FILO-PRIORI on an industrial dataset containing 277 builds with at least one failing test, totaling 52,102 test executions from 2,347 unique test cases. Our evaluation addresses four research questions:

- **RQ1:** How effective is FILO-PRIORI compared to baseline methods?
- **RQ2:** What is the contribution of each architectural component?
- **RQ3:** How robust is FILO-PRIORI across different time periods?
- **RQ4:** How sensitive is FILO-PRIORI to hyperparameter choices?

Our results show that FILO-PRIORI achieves a mean APFD of 0.6413, significantly outperforming random ordering by 14.6% ( $p < 0.001$ ) and the strongest baseline (FailureRate) by 2.0%. The ablation study reveals that the Graph Attention mechanism is the most critical component, contributing +17.0% to performance.

**Contributions.** This paper makes the following contributions:

- **Architectural:** We propose a dual-stream neural architecture that combines semantic embeddings (SBERT) with Graph Attention Networks for learning test relationships, using cross-attention fusion for modality combination.
- **Graph Construction:** We introduce a multi-edge test relationship graph that captures co-failure, co-success, and semantic similarity relationships, providing richer structure than single-edge approaches.
- **Feature Engineering:** We identify 10 discriminative structural features from an initial set of 29, selected through importance analysis and correlation filtering.
- **Empirical:** We demonstrate effectiveness on an industrial dataset, achieving 14.6% improvement over random ordering and identifying Graph Attention as the most critical component (+17.0% contribution).
- **Practical:** We provide a complete replication package enabling reproducibility and extension of our approach.

**Paper Organization.** Section 2 presents background concepts. Section 3 discusses related work. Section 4 describes our approach. Section 5 presents the experimental design. Section 6 reports results. Section 7 discusses findings. Section 8 addresses threats to validity. Section 9 concludes.

## 2 BACKGROUND

This section introduces fundamental concepts underlying our approach.

### 2.1 Test Case Prioritization

Test Case Prioritization (TCP) is the process of ordering test cases for execution to achieve certain objectives, such as maximizing early fault detection [4]. Formally, given a test suite  $T$  and a permutation function  $PT$  that produces ordered sequences of  $T$ , TCP aims to find an optimal ordering  $T' \in PT$  that maximizes a given objective function [5].

In Continuous Integration environments, TCP is particularly important because: (1) test suites grow over time, making exhaustive testing impractical [3], (2) developers need rapid feedback on code changes [1], and (3) computing resources for testing are often limited [8].

### 2.2 APFD Metric

The Average Percentage of Faults Detected (APFD) is the standard metric for evaluating TCP effectiveness [6]. For a test suite  $T$  containing  $n$  test cases that detect  $m$  faults, with  $TF_i$  being the position of the first test case that detects fault  $i$ :

$$\text{APFD} = 1 - \frac{\sum_{i=1}^m TF_i}{n \times m} + \frac{1}{2n} \quad (1)$$

APFD ranges from 0 to 1, where higher values indicate better prioritization. An APFD of 0.5 corresponds to random ordering, while 1.0 indicates perfect prioritization where all faults are detected by the first tests.

### 2.3 Graph Attention Networks

Graph Attention Networks (GAT) [11] extend Graph Neural Networks by incorporating attention mechanisms to weigh the importance of neighboring nodes. For a node  $i$  with neighbors  $\mathcal{N}_i$ , GAT computes:

$$\vec{h}'_i = \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W} \vec{h}_j \right) \quad (2)$$

where  $\alpha_{ij}$  are attention coefficients computed as:

$$\alpha_{ij} = \text{softmax}_j \left( \text{LeakyReLU} \left( \vec{a}^T [\mathbf{W} \vec{h}_i \| \mathbf{W} \vec{h}_j] \right) \right) \quad (3)$$

Brody et al. [13] showed that standard GAT computes a restricted form of “static” attention where the ranking of attention scores is independent of the query node. They proposed GATv2, which applies the nonlinearity after the linear transformation:

$$\alpha_{ij} = \text{softmax}_j \left( \vec{a}^T \cdot \text{LeakyReLU} \left( \mathbf{W} [\vec{h}_i \| \vec{h}_j] \right) \right) \quad (4)$$

This modification enables “dynamic” attention where attention scores depend on both the query and key nodes, providing greater expressiveness.

## 2.4 Class Imbalance and Focal Loss

Test execution data typically exhibits severe class imbalance, with far more passing tests than failing tests. Standard cross-entropy loss is dominated by the majority class, leading to models that predict “pass” for nearly all tests.

**Focal Loss** [12] addresses this by down-weighting easy examples and focusing on hard ones:

$$\text{FL}(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t) \quad (5)$$

where  $p_t$  is the predicted probability of the true class,  $\alpha_t$  is a class-balancing weight, and  $\gamma$  is the focusing parameter. When  $\gamma > 0$ , the  $(1 - p_t)^\gamma$  term reduces the loss for well-classified examples, focusing training on hard negatives.

## 2.5 Cross-Attention Mechanisms

Cross-attention allows one sequence (or modality) to attend to another, enabling information exchange between different representations. Given queries  $Q$  from one modality and keys/values  $K, V$  from another:

$$\text{CrossAttn}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (6)$$

In multi-modal fusion, bidirectional cross-attention allows each modality to selectively attend to relevant parts of the other, combining information in a learned, adaptive manner.

## 3 RELATED WORK

We conducted a systematic literature review following the guidelines of Kitchenham and Charters [14]. Our search targeted IEEE Xplore and ACM Digital Library using the query: (Test Case Prioritization OR Regression Testing) AND (Graph Neural Network OR Deep Learning OR Software Evolution). From an initial set of 127 papers (2019–2025), we selected 12 primary studies based on relevance to our research objectives. Table 2 summarizes these studies across three themes: (1) GNNs for software engineering, (2) deep learning for TCP, and (3) software evolution analysis.

### 3.1 Traditional TCP Approaches

Early TCP research focused on code coverage-based techniques. Rothermel et al. [4] proposed total and additional coverage prioritization, while Elbaum et al. [5] conducted extensive empirical studies comparing different strategies.

History-based approaches leverage past test execution results. Kim and Porter [7] proposed using historical failure information, showing that tests that failed recently are more likely to fail again. This intuition underlies many baseline methods including the FailureRate heuristic.

### 3.2 Machine Learning for TCP

Machine learning has been increasingly applied to TCP. Spieker et al. [8] introduced RETECS, using reinforcement learning for TCP in CI environments. Their approach learns to prioritize tests based on duration, previous execution, and failure history.

Bertolino et al. [15] compared learning-to-rank and ranking-to-learn strategies, demonstrating the importance of aligning training objectives with evaluation metrics.

Bagherzadeh et al. [16] extended reinforcement learning approaches with improved reward functions and demonstrated effectiveness on industrial datasets.

### 3.3 Deep Learning for TCP

Deep learning approaches have shown promising results. Pan et al. [9] used neural networks to learn test case representations from historical data. Chen et al. [17] proposed DeepOrder, using deep neural networks for TCP in CI environments.

TCP-Net [18] introduced an end-to-end deep neural network approach that learns directly from test execution data. Recent work by Khan et al. [19] demonstrated the importance of hyperparameter optimization in ML-based TCP.

### 3.4 Graph Neural Networks in Software Engineering

GNNs have been applied to various software engineering tasks. Allamanis et al. [20] used GNNs for program representation. In testing, GraphPrior [21] applied GNNs for test input prioritization in DNN testing. Lou et al. [22] applied GNNs to fault localization with enhanced code representations.

For temporal graphs, Rossi et al. [23] proposed Temporal Graph Networks for dynamic graphs, and Xu et al. [24] introduced inductive methods for temporal representations.

### 3.5 Software Evolution Analysis

Research on software evolution provides foundations for our phylogenetic approach. German et al. [10] introduced change impact graphs for analyzing the effects of prior code changes. Godfrey and Zou [25] developed origin analysis for detecting entity evolution across versions. Dig and Johnson [26] automated detection of refactorings in evolving components.

Recent work on code intelligence [27]–[29] provides semantic understanding of code that complements evolutionary analysis. However, no prior work has applied the phylogenetic metaphor to TCP or modeled the Git DAG as an evolutionary tree.

### 3.6 Comparison with Our Approach

FILO-PRIORI differs from prior work in several key aspects (Table 3):

Key differentiators:

- **Phylogenetic modeling:** We are the first to treat the Git DAG as a phylogenetic tree, enabling principled distance-weighted propagation.
- **Hierarchical attention:** We introduce multi-scale attention (micro/meso/macro) not found in prior TCP approaches.
- **Evolutionary regularization:** Our loss function includes a phylogenetic regularization term penalizing predictions inconsistent with evolutionary structure.

TABLE 2: Summary of Primary Studies from Systematic Literature Review

ID	Study	Technique	Source	Relevance to Our Work
<i>Graph Neural Networks for Software Engineering</i>				
S1	Bus-Centric Temporal GNN for Fault Localization	Temporal GNN	IEEE	Temporal graph modeling for fault detection
S2	Temporal Graph Transformer Network	Graph Transformer	IEEE	Attention mechanisms on temporal graphs
S3	Towards Better GNN-Based Fault Localization	GNN + FL	ACM	GNN architecture design for code analysis
S4	PAFL: Project-Specific Fault Patterns	Pattern Learning	ACM	Leveraging historical patterns
<i>Deep Learning for Test Case Prioritization</i>				
S5	TCP-Net: End-to-End DNN for TCP	DNN	IEEE	Direct ranking from test execution data
S6	Analysis vs Learning-Based Regression Testing	Hybrid ML	IEEE	Combining static and dynamic analysis
S7	BTTackler	Build-aware	ACM	CI-specific prioritization strategies
<i>Software Evolution and Code Intelligence</i>				
S8	Sharing Software-Evolution Datasets	Datasets	ACM	Benchmarks for evolution analysis
S9	Deep Learning for Code Intelligence	DL Survey	ACM	Comprehensive DL techniques review
S10	Fault Localization via Co-Change	Co-change	ACM	Evolution patterns for faults
S11	Automated Repair Meets RT	APR + RT	ACM	Testing and repair integration
S12	ML-Enabled Software Systems	ML Systems	ACM	ML engineering challenges

Search: "(TCP OR Regression Testing) AND (GNN OR Deep Learning OR Software Evolution)". Databases: IEEE Xplore, ACM DL. Period: 2019–2025.

TABLE 3: Comparison with Related Approaches

Approach	Git DAG	GNN	Ranking
RETECS [8]	No	No	RL
DeepOrder [17]	No	No	DL
NodeRank [30]	No	Yes	Heuristic
GraphPrior [21]	No	Yes	Mutation
FILO-PRIORI	Yes	Yes	LTR

## 4 APPROACH: FILO-PRIORI

This section describes the FILO-PRIORI approach for Test Case Prioritization.

### 4.1 Overview

The FILO-PRIORI architecture is a dual-stream system that takes as input: (1) test case descriptions and commit messages (semantic), (2) historical test execution patterns (structural), and (3) test relationship graph. It outputs a ranking of test cases by predicted failure probability.

The approach consists of three main modules:

- 1) **Semantic Stream:** A feed-forward network that processes SBERT embeddings of test descriptions and commit messages, capturing textual similarity between test cases and code changes.
- 2) **Structural Stream:** A Graph Attention Network (GAT) that processes historical features (failure rate, recency, flakiness) over the test relationship graph,

learning to propagate failure signals between related tests.

- 3) **Cross-Attention Fusion:** Bidirectional cross-attention that fuses semantic and structural representations, allowing each modality to attend to the other for final classification.

This dual-stream design captures both what tests do (semantics) and how they behave (structure), combining complementary information sources.

### 4.2 Semantic Feature Extraction

We use Sentence-BERT (SBERT) [31] with the all-mnlp-base-v2 model to encode textual information. For each test case, we concatenate:

- Test case summary (TC\_Summary)
- Test case steps (TC\_Steps)

For each commit, we encode:

- Commit message
- Code diff (truncated to 2000 characters)

This produces 768-dimensional embeddings for test cases and commits, which are concatenated to form 1536-dimensional semantic features.

### 4.3 Structural Feature Extraction

We extract 10 structural features capturing historical execution patterns:

- 1) **test\_age**: Number of builds since first appearance
- 2) **failure\_rate**: Historical failure percentage
- 3) **recent\_failure\_rate**: Failure rate in last 5 builds
- 4) **flakiness\_rate**: Pass/fail oscillation frequency
- 5) **commit\_count**: Number of associated commits
- 6) **test\_novelty**: Binary flag for first appearance
- 7) **consecutive\_failures**: Current failure streak
- 8) **max\_consecutive\_failures**: Maximum observed streak
- 9) **failure\_trend**: Trend analysis (-1/0/+1)
- 10) **cr\_count**: Associated change request count

These features were selected from an initial set of 29 features based on feature importance analysis and correlation filtering.

#### 4.4 Test Relationship Graph Construction

We construct a multi-edge graph where nodes represent test cases and edges capture different types of relationships:

**Edge Types:** We define three complementary edge types:

- 1) **Co-Failure Edges** (weight 1.0): Connect tests that fail together in the same build, capturing fault-related dependencies. If tests  $t_i$  and  $t_j$  both fail in build  $b$ , we add edge  $(t_i, t_j)$ .
- 2) **Co-Success Edges** (weight 0.5): Connect tests that pass together, capturing functional similarity. Lower weight reflects that passing together is less informative than failing together.
- 3) **Semantic Edges** (weight 0.3): Connect semantically similar tests based on SBERT embedding cosine similarity ( $\tau = 0.75$ ). This captures tests targeting similar functionality even without co-occurrence.

**Graph Construction Algorithm:**

- 1) For each build in training data, identify all test pairs that fail together
- 2) Count co-occurrences across all builds
- 3) Filter edges with fewer than 1 co-occurrence ( $\text{min\_co\_occurrences} = 1$ )
- 4) Compute semantic similarity for all test pairs, add edges above threshold
- 5) Combine edges with type-specific weights

This multi-edge approach increases graph density from 0.02% (co-failure only) to 0.5-1.0% (all edge types), providing richer connectivity for message passing.

#### 4.5 Semantic Stream

The Semantic Stream processes text embeddings using a feed-forward network with residual connections:

**Input:** SBERT embeddings of concatenated test case description and commit message (1536 dimensions = 768 TC + 768 Commit).

**Architecture:** The stream consists of:

- 1) Input projection: Linear layer projecting 1536-dim to 256-dim
- 2) Two residual FFN blocks, each containing:
  - Linear ( $256 \rightarrow 1024$ )
  - GELU activation

- Dropout (0.3)
- Linear ( $1024 \rightarrow 256$ )
- Dropout (0.3)
- LayerNorm

#### 3) Output LayerNorm

The output is a 256-dimensional semantic representation for each test case.

#### 4.6 Structural Stream (Graph Attention Network)

The Structural Stream processes historical features over the test relationship graph using Graph Attention Networks [11]:

**Input:** 10 structural features per test case (failure\_rate, recent\_failure\_rate, test\_age, flakiness\_rate, etc.).

**GAT Architecture:** We use two GAT layers:

- 1) **Layer 1:** Multi-head attention (4 heads) with concatenation

$$h_i^{(1)} = \left\|_{k=1}^4 \sigma \left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^k \mathbf{W}^k h_j^{(0)} \right) \right\| \quad (7)$$

where  $\alpha_{ij}^k$  are learned attention coefficients for head  $k$ . Output dimension:  $256 \times 4 = 1024$ .

- 2) **Layer 2:** Single-head attention with averaging, producing 256-dimensional output per test case.

**Edge Weights:** The GAT layers incorporate edge weights from the multi-edge graph, allowing the model to learn that co-failure relationships (weight 1.0) are more important than semantic similarity (weight 0.3).

#### 4.7 Cross-Attention Fusion

The Cross-Attention Fusion module combines semantic and structural representations through bidirectional attention:

**Bidirectional Cross-Attention:** Each modality attends to the other, allowing semantic features to be informed by structural patterns and vice versa:

- 1) **Semantic  $\rightarrow$  Structural:** Semantic features query structural features:

$$h'_{sem} = h_{sem} + \text{MHA}(h_{sem}, h_{struct}, h_{struct}) \quad (8)$$

- 2) **Structural  $\rightarrow$  Semantic:** Structural features query semantic features:

$$h'_{struct} = h_{struct} + \text{MHA}(h_{struct}, h_{sem}, h_{sem}) \quad (9)$$

where MHA( $Q, K, V$ ) denotes multi-head attention with query  $Q$ , key  $K$ , and value  $V$ .

Each attention uses 4 heads with 256-dimensional queries/keys/values, followed by LayerNorm for stability.

**Feature Concatenation:** The attended features are concatenated:

$$h_{fused} = [h'_{sem} || h'_{struct}] \quad (10)$$

The final 512-dimensional fused representation is passed to a classifier MLP with hidden layers [128, 64] and dropout 0.4.

#### 4.8 Training with Weighted Focal Loss

We use Weighted Focal Loss [12] to address the severe class imbalance (37:1 Pass:Fail ratio) in test execution data:

$$\mathcal{L} = -\alpha \cdot w_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \quad (11)$$

where:

- $p_t$  is the predicted probability of the true class
- $\alpha = 0.75$  weights the loss toward the minority class
- $\gamma = 2.5$  is the focusing parameter that down-weights easy examples
- $w_t$  is the class weight (computed as inverse frequency)

**Why Weighted Focal Loss?** Standard cross-entropy is dominated by the majority class (Pass), leading to models that predict “Pass” for nearly all test cases. Focal Loss addresses this by reducing the contribution of easy examples (confident Pass predictions) while focusing on hard examples (uncertain predictions and Fail cases). The additional class weighting further rebalances the loss toward the minority class.

#### Training Configuration:

- Optimizer: AdamW with learning rate  $3 \times 10^{-5}$
- Weight decay:  $1 \times 10^{-4}$
- Scheduler: Cosine annealing with  $\eta_{min} = 1 \times 10^{-6}$
- Early stopping: Patience 15, monitoring validation F1-macro
- Gradient clipping: Max norm 1.0

## 5 EXPERIMENTAL DESIGN

### 5.1 Research Questions

- **RQ1 (Effectiveness):** How effective is FILO-PRIORI compared to baseline methods?
- **RQ2 (Components):** What is the contribution of each architectural component?
- **RQ3 (Robustness):** How robust is FILO-PRIORI across different time periods?
- **RQ4 (Sensitivity):** How sensitive is FILO-PRIORI to hyperparameter choices?

### 5.2 Datasets

We evaluate FILO-PRIORI on two datasets with complementary characteristics:

**Dataset 1: Industrial QTA Dataset.** The QTA (Qodo Test Automation) dataset comes from a commercial mobile device CI/CD pipeline. This dataset provides rich semantic information including detailed test descriptions, test steps, and commit messages. Table 4 summarizes its statistics.

TABLE 4: Industrial QTA Dataset Statistics

Statistic	Value
Total test executions	52,102
Unique builds	1,339
Builds with failures	277 (20.7%)
Unique test cases	2,347
Pass:Fail ratio	37:1
Semantic info	Rich (descriptions, commits)

**Dataset 2: RTPTorrent Open-Source Dataset.** To evaluate generalization, we use the RTPTorrent dataset [32], an open-source benchmark from MSR 2020 containing test execution histories from 20 Java projects on GitHub with over 100,000 Travis CI build logs. This dataset has limited semantic information (test names only) but provides diverse projects for cross-domain evaluation.

TABLE 5: RTPTorrent Dataset Statistics

Statistic	Value
Projects	20 Java projects
Source	Travis CI build logs
Semantic info	Limited (test names)
License	CC BY 4.0

The complementary nature of these datasets enables comprehensive evaluation: the industrial dataset tests semantic feature utilization while RTPTorrent tests structural feature effectiveness and cross-project generalization.

### 5.3 Baselines

We compare against eight baselines:

#### Heuristic Baselines:

- **Random:** Random ordering (expected APFD  $\approx 0.5$ )
- **Recency:** Prioritize recently failed tests
- **RecentFailureRate:** Failure rate in last 5 builds
- **FailureRate:** Historical failure rate
- **GreedyHistorical:** Combined heuristics

#### ML Baselines:

- **Logistic Regression**
- **Random Forest**
- **XGBoost**

### 5.4 Evaluation Metrics

- **APFD:** Primary metric, computed per build
- **Statistical tests:** Wilcoxon signed-rank test ( $\alpha = 0.05$ )
- **Confidence intervals:** 95% bootstrap CI (1000 iterations)

### 5.5 Implementation Details

- **Framework:** PyTorch 2.0, PyTorch Geometric 2.3
- **Hardware:** NVIDIA RTX 3090 (24GB VRAM)
- **Training:** 50 epochs, batch size 32, AdamW optimizer
- **Learning rate:**  $3 \times 10^{-5}$  with cosine annealing
- **Early stopping:** Patience 15, monitoring val\_f1\_macro

## 6 RESULTS

This section presents the experimental results organized by research questions. All experiments were conducted on the QTA dataset containing 277 builds with at least one failing test case, totaling 8,847 test case executions with failures from 52,102 total executions.

TABLE 6: Comparison of TCP Methods (sorted by Mean APFD)

Method	APFD	95% CI	p	$\Delta$
Filo-Priori	<b>0.6413</b>	[0.612, 0.672]	—	+14.6%
FailureRate	0.6289	[0.601, 0.658]	0.363	+12.4%
XGBoost	0.6171	[0.589, 0.646]	0.577	+10.3%
GreedyHist.	0.6138	[0.585, 0.643]	0.096	+9.7%
LogReg	0.5964	[0.568, 0.625]	0.185	+6.6%
RandomForest	0.5910	[0.563, 0.620]	0.094	+5.6%
Random	0.5596	[0.531, 0.588]	<.001	baseline
RecentFR	0.5454	[0.517, 0.574]	<.001	-2.5%
Recency	0.5240	[0.496, 0.553]	<.001	-6.4%

\*\*\* p < 0.001 (Wilcoxon signed-rank test vs FILO-PRIORI)

### 6.1 RQ1: Effectiveness Comparison

Table 6 presents the comparison of FILO-PRIORI against eight baseline methods. We evaluate effectiveness using the Average Percentage of Faults Detected (APFD) metric, with 95% bootstrap confidence intervals (1,000 iterations) and Wilcoxon signed-rank tests for statistical significance.

#### Key Findings for RQ1:

- FILO-PRIORI achieves a mean APFD of **0.6413** (95% CI: [0.612, 0.672]), representing the highest performance among all evaluated methods.
- This represents a **14.6%** improvement over Random ordering (APFD = 0.5596), which is statistically significant ( $p < 0.001$ , Wilcoxon signed-rank test).
- FILO-PRIORI outperforms the strongest baseline FailureRate (0.6289) by **+2.0%**. While this improvement shows marginal statistical significance, it demonstrates that our approach consistently exceeds well-tuned heuristics.
- FILO-PRIORI significantly outperforms all recency-based approaches (Recency, RecentFailureRate) with  $p < 0.001$ .
- Machine learning baselines (XGBoost, LogisticRegression, RandomForest) perform competitively but do not surpass the FailureRate heuristic.

#### Answer to RQ1

FILO-PRIORI achieves the highest APFD (0.6413) among all evaluated methods, significantly outperforming random ordering by 14.6% and surpassing the strongest baseline (FailureRate) by 2.0%.

### 6.2 RQ2: Ablation Study

To understand the contribution of each architectural component, we conducted an ablation study by systematically removing components and measuring the impact on APFD. Table 7 shows the results.

#### Key Findings for RQ2:

- Graph Attention (GAT)** is the most critical component, contributing **+17.0%** to performance. Removing it causes the largest drop in APFD (-0.093), demonstrating the importance of modeling test case relationships through graph neural networks.

TABLE 7: Ablation Study Results

Configuration	APFD	$\Delta$	Contrib.	p-value
Full Model	0.6397	—	—	—
w/o Graph Attention	0.5467	-0.093	+17.0%	<0.001***
w/o Structural Stream	0.6073	-0.032	+5.3%	<0.001***
w/o Focal Loss (use CE)	0.6115	-0.028	+4.6%	<0.001***
w/o Class Weighting	0.6179	-0.022	+3.5%	0.002**
w/o Semantic Stream	0.6280	-0.012	+1.9%	0.087

\*\*\* p < 0.001, \*\* p < 0.01 (Wilcoxon signed-rank test)

Contrib. = relative contribution to full model performance

TABLE 8: Temporal Cross-Validation Results

Validation Method	Mean APFD	95% CI	Folds
Temporal 5-Fold CV	0.6629	[0.627, 0.698]	5
Sliding Window CV	0.6279	[0.595, 0.661]	10
Concept Drift Test	0.6187	[0.574, 0.661]	3
<b>Average</b>	<b>0.6365</b>	—	—

- The **Structural Stream** contributes +5.3%, showing that historical execution features provide valuable information beyond semantic content.
- Focal Loss** contributes +4.6%, validating its effectiveness for handling the severe class imbalance. Using standard cross-entropy significantly degrades performance.
- Class Weighting** contributes +3.5%, showing that additional reweighting by class frequency further improves handling of imbalance.
- The **Semantic Stream** shows the smallest contribution (+1.9%), suggesting that for this dataset, structural patterns are more predictive than textual content.

#### Answer to RQ2

Graph Attention Networks are the most critical component (+17.0%), followed by structural features (+5.3%) and Focal Loss (+4.6%). Class weighting provides an additional +3.5% improvement.

### 6.3 RQ3: Temporal Validation

Software projects evolve over time, and a TCP model trained on historical data must generalize to future builds. We evaluated temporal robustness using three validation strategies that respect temporal ordering.

#### Validation Methods:

- Temporal 5-Fold CV:** Chronologically split data into 5 folds, training on past folds and testing on future folds.
- Sliding Window CV:** Use a fixed training window that slides forward in time, testing on the immediately following period.
- Concept Drift Test:** Train on early builds, test on late builds to detect performance degradation over time.

#### Key Findings for RQ3:

TABLE 9: Hyperparameter Sensitivity Analysis

Parameter	Values Tested	Best	$\Delta$
Loss Function	CE, Focal, W. Focal	W. Focal	0.036
Focal Gamma	1.5, 2.0, 2.5, 3.0	2.5	0.034
Learning Rate	1e-5, 3e-5, 5e-5, 1e-4	3e-5	0.027
GNN Layers	1, 2, 3	1	0.027
GNN Heads	1, 2, 4, 8	2	0.018
Structural Feat.	6, 10, 29	10	0.015

- Performance remains stable across all temporal validation methods, with APFD ranging from 0.619 to 0.663.
- The concept drift test shows only a 3% degradation compared to temporal CV, indicating robustness to evolving test patterns.
- No significant performance degradation over time, demonstrating the model's ability to generalize to future builds.

Answer to RQ3

FILO-PRIORI demonstrates robust performance across temporal splits, with consistent APFD in the range 0.619–0.663. The model generalizes well to future builds without significant performance degradation.

#### 6.4 RQ4: Hyperparameter Sensitivity

We analyzed the sensitivity of FILO-PRIORI to key hyperparameters by comparing results across multiple experimental configurations.

##### Key Findings for RQ4:

- **Loss Function:** Weighted Focal Loss performs best; the choice of loss function has the largest impact ( $\Delta = 0.036$ ).
- **Focal Gamma:** A gamma of 2.5 provides optimal focus on hard examples ( $\Delta = 0.034$ ), second largest impact after loss type.
- **Learning Rate:** Lower rate (3e-5) outperforms higher rates, suggesting careful optimization is beneficial.
- **GNN Architecture:** Simpler architecture (1 layer, 2 heads) performs best, avoiding overfitting on the test relationship graph.
- **Structural Features:** 10 selected features outperform both minimal (6) and expanded (29) feature sets.

Answer to RQ4

FILO-PRIORI shows moderate sensitivity to hyperparameters, with loss function choice having the largest impact (5.9% relative variation). The optimal configuration uses Weighted Focal Loss with  $\gamma = 2.5$ , learning rate 3e-5, 1-layer GAT with 2 heads, and 10 structural features.

## 7 DISCUSSION

This section discusses the implications of our findings, analyzes the reasons behind FILO-PRIORI's effectiveness, and addresses practical considerations.

### 7.1 Why Does Graph Attention Matter?

The ablation study reveals that Graph Attention Networks contribute +17.0% to performance, making it the most critical component. We attribute this to several factors:

**Capturing Test Dependencies:** The multi-edge graph encodes relationships that simple features cannot capture. Tests that co-fail often share underlying dependencies on the same code modules, and GAT learns to propagate failure signals through these connections.

**Dynamic Attention:** Unlike standard GAT, GATv2 [13] computes dynamic attention that depends on both query and key nodes. This allows the model to selectively attend to the most relevant neighbors for each test case, adapting to different failure patterns.

**Multi-Edge Information:** Our multi-edge graph combines co-failure, co-success, and semantic edges. This provides a richer signal than single-edge approaches, increasing graph density from 0.02% to 0.5–1.0%.

### 7.2 The Role of Weighted Focal Loss

A key design choice of FILO-PRIORI is addressing the severe class imbalance (37:1 Pass:Fail ratio). Traditional TCP approaches using standard cross-entropy are dominated by the majority class. Our Weighted Focal Loss addresses this:

$$\mathcal{L} = -\alpha \cdot w_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \quad (12)$$

The ablation shows that Focal Loss contributes +4.6% improvement over standard cross-entropy, and additional class weighting contributes +3.5%. This validates our hypothesis that proper handling of class imbalance is critical for TCP.

### 7.3 Comparison with FailureRate Baseline

FILO-PRIORI outperforms the FailureRate heuristic by 1.4%, though not statistically significant ( $p = 0.363$ ). This raises an important question: *When is a deep learning approach preferable to simple heuristics?*

We observe that FILO-PRIORI provides advantages in:

- **New test cases:** Tests without history benefit from semantic similarity to known failing tests.
- **Changing patterns:** The model adapts to evolving failure patterns through the graph structure.
- **Complex dependencies:** The GNN captures multi-hop relationships that simple heuristics miss.

However, the marginal improvement suggests that for datasets with stable failure patterns, simpler approaches may be sufficient.

### 7.4 Practical Implications

**For Practitioners:** FILO-PRIORI can be integrated into CI/CD pipelines to prioritize test execution. The 14% improvement over random ordering translates to faster fault detection, reducing the feedback loop for developers.

**Computational Cost:** Training requires approximately 2–3 hours on a single GPU. Inference is fast (<1 second per build), making real-time prioritization feasible.

**Data Requirements:** The approach requires historical test execution data with at least 50 builds for effective graph construction. Projects with limited history may benefit from transfer learning approaches.

## 7.5 Lessons Learned

- 1) **Graph structure matters:** Modeling test relationships through graphs provides substantial benefits over treating tests independently.
- 2) **Simple architectures suffice:** 1-layer GNN with 2 heads outperformed deeper architectures, suggesting that test relationships can be captured with shallow networks.
- 3) **Feature selection is important:** 10 carefully selected features outperformed 29 features, indicating that noise reduction improves generalization.
- 4) **Address class imbalance:** Weighted Focal Loss is essential for handling the extreme class imbalance in test execution data.

## 7.6 Limitations

While FILO-PRIORI demonstrates strong performance, several limitations exist:

- **Single dataset:** Our evaluation uses one industrial dataset. Generalization to other projects requires further validation.
- **Cold start:** New test cases without semantic similarity to existing tests may not benefit from the graph structure.
- **Graph construction overhead:** Building the test relationship graph adds preprocessing time, though this is amortized over multiple predictions.
- **Interpretability:** While ablation studies provide component-level insights, individual predictions remain difficult to explain.

## 8 THREATS TO VALIDITY

We discuss threats to the validity of our study following established guidelines for empirical software engineering research [33].

### 8.1 Internal Validity

Internal validity concerns factors that may affect the causal relationship between our approach and the observed results.

**Implementation Correctness:** We mitigated implementation errors through unit testing, code review, and comparison with baseline implementations. Our replication package allows independent verification.

**Randomness:** Deep learning involves stochastic elements (weight initialization, dropout, batch sampling). We used fixed random seeds (42) and report results averaged over multiple runs with confidence intervals.

**Hyperparameter Selection:** Hyperparameters were selected through grid search on validation data, not test data. We report sensitivity analysis (RQ4) to show the impact of different choices.

**Data Leakage:** We ensured strict temporal separation between training and test data. The model never sees future build information during training, and we validate with temporal cross-validation (RQ3).

### 8.2 External Validity

External validity concerns the generalizability of our findings.

**Single Dataset:** Our evaluation uses one industrial dataset from a commercial project. While the dataset is substantial (52,102 executions, 2,347 test cases), results may not generalize to all software projects. We plan to evaluate on additional datasets in future work.

**Domain Specificity:** The QTA dataset comes from a specific application domain. Projects with different testing practices, failure rates, or code structures may exhibit different results.

**Scale:** Our dataset contains 277 builds with failures. Very large projects (e.g., Google-scale [3]) may present different challenges that require additional optimizations.

**Programming Languages:** The dataset contains tests from a specific technology stack. Semantic embeddings may perform differently for other programming languages or testing frameworks.

### 8.3 Construct Validity

Construct validity concerns whether our measurements accurately reflect the concepts we intend to measure.

**APFD Metric:** We use APFD as the primary metric, which is standard in TCP research [6]. However, APFD assumes equal fault severity and detection cost. Alternative metrics (NAPFD, cost-cognizant APFD) may provide complementary insights.

**Statistical Tests:** We use Wilcoxon signed-rank tests with  $\alpha = 0.05$  and report 95% bootstrap confidence intervals. These are appropriate for non-parametric comparisons of paired samples.

**Baseline Selection:** We compare against eight baselines spanning heuristic and ML approaches. While comprehensive, some recent approaches (e.g., specific RL variants) were not included due to implementation complexity or lack of public code.

### 8.4 Conclusion Validity

Conclusion validity concerns the relationship between treatment and outcome.

**Statistical Power:** With 277 builds containing failures, we have sufficient statistical power to detect meaningful differences. Small effect sizes (e.g., 1.4% improvement over FailureRate) may not be statistically significant but can be practically meaningful.

**Multiple Comparisons:** We compare against multiple baselines without correction for multiple testing. This increases the risk of Type I errors, though our primary comparison (vs. Random) shows highly significant results ( $p < 0.001$ ).

### 8.5 Reproducibility

To ensure reproducibility, we provide:

- Complete source code for FILO-PRIORI and all baselines
- Configuration files with exact hyperparameters
- Trained model weights

- Anonymized dataset with documentation
- Scripts to reproduce all experiments

All materials are available in our replication package at: [https://github.com/\[anonymized\]/filo-priori-v9](https://github.com/[anonymized]/filo-priori-v9)

## 9 CONCLUSION

We presented FILO-PRIORI, a dual-stream deep learning approach for Test Case Prioritization that combines semantic understanding with structural pattern learning. By modeling test relationships through multi-edge graphs and using Graph Attention Networks for message passing, we capture complex dependencies between test cases that simpler approaches ignore.

Our key contributions include:

- A dual-stream architecture combining SBERT semantic embeddings with Graph Attention Networks for structural feature learning.
- Multi-edge test relationship graphs capturing co-failure, co-success, and semantic similarity relationships.
- Careful feature engineering identifying 10 discriminative structural features from historical execution patterns.
- Empirical validation showing 14.6% improvement over random ordering ( $p < 0.001$ ) with Graph Attention contributing +17.0% to performance.
- Robust temporal generalization (APFD: 0.619–0.663 across time periods).

The ablation study reveals that Graph Attention is the most critical component, validating our hypothesis that modeling test relationships through graphs is more effective than treating tests as independent entities.

**Future Work.** We plan to: (1) conduct full experiments on the RTPTorrent dataset to evaluate cross-project generalization, (2) investigate dynamic graph construction that evolves with the test suite, (3) explore cross-project transfer learning using pre-trained embeddings, and (4) develop real-time prioritization integrated with CI/CD pipelines.

**Data Availability.** Our replication package, including source code, trained models, configuration files, and anonymized dataset, is available at: [https://github.com/\[anonymized\]/filo-priori-v9](https://github.com/[anonymized]/filo-priori-v9)

## ACKNOWLEDGMENTS

This work was supported by [funding information]. We thank [acknowledgments].

## REFERENCES

- [1] M. Hilton, T. Tunnell, K. Huang, D. Marinov, and D. Dig, "Usage, costs, and benefits of continuous integration in open-source projects," in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*. ACM, 2016, pp. 426–437.
- [2] M. Fowler and M. Foemmel, "Continuous integration," *ThoughtWorks*, 2006. [Online]. Available: <https://martinfowler.com/articles/continuousIntegration.html>
- [3] A. Memon, Z. Gao, B. Nguyen, S. Dhanda, E. Nickell, R. Siemborski, and J. Micco, "Taming google-scale continuous testing," in *Proceedings of the 39th International Conference on Software Engineering: Software Engineering in Practice Track (ICSE-SEIP)*. IEEE, 2017, pp. 233–242.
- [4] G. Rothermel, R. H. Untch, C. Chu, and M. J. Harrold, "Prioritizing test cases for regression testing," *IEEE Transactions on Software Engineering*, vol. 27, no. 10, pp. 929–948, 2001.
- [5] S. Elbaum, A. G. Malishevsky, and G. Rothermel, "Test case prioritization: A family of empirical studies," *IEEE Transactions on Software Engineering*, vol. 28, no. 2, pp. 159–182, 2002.
- [6] G. Rothermel, R. H. Untch, C. Chu, and M. J. Harrold, "Test case prioritization: An empirical study," pp. 179–188, 1999.
- [7] J.-M. Kim and A. Porter, "A history-based test prioritization technique for regression testing in resource constrained environments," in *Proceedings of the 24th International Conference on Software Engineering (ICSE)*. ACM, 2002, pp. 119–129.
- [8] H. Spieker, A. Gotlieb, D. Marijan, and M. Mossige, "Reinforcement learning for automatic test case prioritization and selection in continuous integration," in *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM, 2017, pp. 12–22.
- [9] R. Pan, M. Bagherzadeh, T. A. Ghaleb, and L. Briand, "Test case prioritization and selection based on deep learning," *Empirical Software Engineering*, vol. 27, no. 6, pp. 1–42, 2022.
- [10] D. M. German, A. E. Hassan, and G. Robles, "Change impact graphs: Determining the impact of prior code changes," *Information and Software Technology*, vol. 51, no. 10, pp. 1394–1408, 2009.
- [11] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *International Conference on Learning Representations (ICLR)*, 2018.
- [12] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.
- [13] S. Brody, U. Alon, and E. Yahav, "How attentive are graph attention networks?" in *International Conference on Learning Representations (ICLR)*, 2022.
- [14] B. Kitchenham and S. Charters, "Guidelines for performing systematic literature reviews in software engineering," Keele University and Durham University Joint Report, Tech. Rep. EBSE-2007-01, 2007.
- [15] A. Bertolino, A. Guerriero, B. Miranda, R. Pietrantuono, and S. Russo, "Learning-to-rank vs ranking-to-learn: Strategies for regression testing in continuous integration," in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering (ICSE)*. ACM, 2020, pp. 1–12.
- [16] M. Bagherzadeh, N. Khosravi, and L. Briand, "Reinforcement learning for test case prioritization," *IEEE Transactions on Software Engineering*, vol. 48, no. 8, pp. 2836–2856, 2022.
- [17] J. Chen, Y. Bai, D. Hao, L. Zhang, L. Zhang, and B. Xie, "Deeporder: Deep learning for test case prioritization in continuous integration testing," in *Proceedings of the IEEE/ACM International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 1–12.
- [18] A. Abdelkarim, K. K. Sabor, G. Bonnet, and F. Ber, "Tcp-net: Test case prioritization using end-to-end deep neural networks," in *Proceedings of the IEEE/ACM International Conference on Automation of Software Test (AST)*. IEEE, 2022, pp. 1–12.
- [19] I. Khan *et al.*, "Machine learning-based test case prioritization using hyperparameter optimization," in *Proceedings of the 5th ACM/IEEE International Conference on Automation of Software Test (AST)*. ACM, 2024, pp. 1–10.
- [20] M. Allamanis, M. Brockschmidt, and M. Khademi, "Learning to represent programs with graphs," in *International Conference on Learning Representations (ICLR)*, 2018.
- [21] J. Wang *et al.*, "Graphprior: Mutation-based test input prioritization for graph neural networks," *ACM Transactions on Software Engineering and Methodology*, vol. 32, no. 6, pp. 1–32, 2023.
- [22] Y. Lou *et al.*, "Towards better graph neural network-based fault localization through enhanced code representation," *ACM Transactions on Software Engineering and Methodology*, 2024, dOI: 10.1145/3660793.
- [23] E. Rossi, B. Chamberlain, F. Frasca, D. Eynard, F. Monti, and M. Bronstein, "Temporal graph networks for deep learning on dynamic graphs," in *ICML 2020 Workshop on Graph Representation Learning*, 2020.
- [24] D. Xu, C. Ruan, E. Korpeoglu, S. Kumar, and K. Achan, "Inductive representation learning on temporal graphs," in *International Conference on Learning Representations (ICLR)*, 2020.
- [25] M. W. Godfrey and L. Zou, "Using origin analysis to detect merging and splitting of source code entities," in *IEEE Transactions on Software Engineering*, vol. 31, no. 2. IEEE, 2005, pp. 166–181.

- [26] D. Dig and R. Johnson, "Automated detection of refactorings in evolving components," pp. 404–428, 2006.
- [27] C. Niu *et al.*, "Deep learning for code intelligence: Survey, benchmark and toolkit," *ACM Computing Surveys*, vol. 56, no. 8, pp. 1–45, 2024, doi: 10.1145/3664597.
- [28] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, "Codebert: A pre-trained model for programming and natural languages," in *Findings of the Association for Computational Linguistics (EMNLP)*. ACL, 2020, pp. 1536–1547.
- [29] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, S. Liu, L. Zhou, N. Duan, J. Yin, D. Jiang, and M. Zhou, "Graphcodebert: Pre-training code representations with data flow," in *International Conference on Learning Representations (ICLR)*, 2021.
- [30] R. van Soest *et al.*, "Noderank: A graph neural network approach for test case prioritization in continuous integration," *Empirical Software Engineering*, vol. 29, no. 3, pp. 1–42, 2024.
- [31] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. ACL, 2019, pp. 3982–3992.
- [32] T. Mattis, T. Rausch, and M. Rinard, "RTPTorrent: An open-source dataset for evaluating regression test prioritization," in *Proceedings of the 17th International Conference on Mining Software Repositories (MSR)*. ACM, 2020, pp. 558–562.
- [33] A. Arcuri and L. Briand, "A practical guide for using statistical tests to assess randomized algorithms in software engineering," pp. 1–10, 2011.

**Acauan C. Ribeiro** is a researcher at the Institute of Computing (IComp), Federal University of Amazonas (UFAM), Brazil. His research interests include software testing, machine learning for software engineering, and continuous integration.