

# Towards Reliable Evaluation of LLM-Based Automated Program Repair

## 1. Introduction

Large Language Models (LLMs) such as GPT, Claude, and CodeLlama are increasingly being used for **Automated Program Repair (APR)**, where models generate patches to fix bugs in software programs. Most existing research evaluates these repairs based on whether the generated patch passes the available test suite. However, test suites are often incomplete, meaning an incorrect patch may still pass the tests. In addition, since LLMs are trained on large amounts of publicly available code, some generated patches may resemble existing solutions rather than being newly reasoned fixes. **In this work, we investigate whether current evaluation methods provide a reliable measure of repair success in LLM-based automated program repair.** We propose an evaluation framework that analyzes generated patches beyond simple test-suite passing to better understand the reliability and correctness of LLM-generated repairs.

## 2. Related Work

Several studies have explored the use of Large Language Models (LLMs) for automated program repair. Existing research mainly focuses on improving patch generation through prompt engineering, fine-tuning code models, retrieval-augmented repair, and multi-agent repair frameworks. While these approaches aim to increase repair effectiveness, relatively little work examines whether the evaluation methods used to measure repair success are reliable. Many studies rely primarily on test-suite passing as the main evaluation metric.

#	Citing paper	Does it solve a Paper A question/problem?	What it seems to address	My assessment
1	A Survey of LLM-based Automated Program Repair	Partly	Re-surveys LLM-based APR, taxonomies, datasets, paradigms, challenges, opportunities	Mostly extends/synthesizes Paper A rather than solving a specific technical APR problem. It advances the survey/future-directions side. ( <a href="#">arXiv</a> )

2	A systematic literature review on large language models for automated program repair	Partly	Another SLR focused on LLMs for APR	Mostly mentions / broadens survey coverage. It revisits the same landscape-level questions as Paper A rather than solving one experimentally. ( <a href="#">Frontiers</a> )
3	A dual perspective review on large language models and code verification	Partly	Reviews safety, verification, vulnerability detection/repair, future directions	Extends the “challenges / verification / trustworthiness” problem that Paper A identifies, but it is still a review, not a direct solution paper. ( <a href="#">Frontiers</a> )
4	InstructRepair: Instruct Large Language Models with Rich Bug Information for Automated Program Repair	Yes	Rich bug information for APR	This appears to directly attack input representation / repair effectiveness, which maps to Paper A’s questions on capabilities, methods, and future directions. Likely a real methodological solution, not just a mention. ( <a href="#">Researcher Life</a> )
5	Test-in-the-loop LLM Repair: Verifiable Automated Program Repair on QuixBugs with a Failing-Test / Patch / Regression-Test Loop	Yes	Verifiable APR with test loop	This appears to directly address Paper A’s challenge around evaluation reliability / false assurance, by making repair verification stronger. ( <a href="#">Researcher Life</a> )
6	Comparative Analysis of Pre-trained Code Language Models for Automated Program Repair via	Yes	Compares pre-trained code models for APR	This directly addresses architecture/model comparison and evaluation, which are central Paper A questions. It looks like an empirical answer, not

	Code Infill Generation			a mention. ( <a href="#">Eindhoven Research Portal</a> )
7	Aligning the Objective of LLM-Based Program Repair	Yes	Objective alignment for APR	This looks like a direct technical response to a gap Paper A leaves open: current repair objectives do not always match truly correct repairs. ( <a href="#">ACM Digital Library</a> )
8	Integrating Attention Mechanism with Code Structural Affinity and Execution Context Correlation for Automated Bug Repair	Yes	Context-aware repair using structure + execution context	This appears to tackle repair quality and contextual awareness, which is consistent with Paper A's challenge analysis. ( <a href="#">ScienceDirect</a> )
9	An Empirical Study on the Effectiveness of Iterative LLM ...	Yes, likely	Iterative LLM-based repair/debugging	From the title/snippet, this seems to address iterative repair strategies, one of the methodological directions discussed by Paper A. Evidence suggests technical follow-up, but I could not inspect the full abstract. ( <a href="#">sol.sbc.org.br</a> )
10	The Impact of Fine-Tuning Large Language Models on Automated Program Repair	Yes	Fine-tuning strategies for APR	This directly addresses Paper A's fine-tuning-methods question. It looks like an empirical follow-up rather than a background citation. ( <a href="#">IEEE Computer Society</a> )
11	Template-Guided Program Repair in the Era of Large Language Models	Yes	Combines templates with LLMs for APR	This appears to solve a concrete method problem: how to better combine traditional APR structure with LLMs. Good match to Paper A's "future

				directions / methods” theme. ( <a href="#">UMass Amherst</a> )
12	Fine-Tuning CodeLlama to Fix Bugs	Yes, likely	Fine-tuning a specific code LLM for bug fixing	Likely addresses the fine-tuning / model adaptation question raised in Paper A. I could verify its existence in related APR search results, but not inspect a full open abstract here. ( <a href="#">ResearchGate</a> )
13	Automated Code Review Using Large Language Models at ...	No, not directly	LLM-based code review	This is adjacent to APR. It likely uses Paper A as background on LLM repair, but it does not appear to directly solve one of Paper A’s core APR questions. ( <a href="#">arXiv</a> )
14	Automated Resolution of Issue Reports using LLM Agents	No, not directly	Issue resolution with agents; cites APR literature	This seems more about agentic issue handling than APR itself. It likely mentions Paper A for context rather than answering its RQs. ( <a href="#">dspace.ut.ee</a> )
15	CodeR3: A GenAI-Powered Workflow Repair and Revival Ecosystem	Partly	Workflow repair/revival, service substitution, validation	This extends repair into workflow decay / workflow migration. It is a real repair contribution, but not a direct answer to mainstream APR benchmark questions from Paper A. ( <a href="#">ResearchGate</a> )
16	The Code Council: Orchestrating Heterogeneous Large Language Models ...	Partly	Multi-LLM orchestration, broader evaluation concerns	The snippet suggests concern with multi-dimensional evaluation such as repair quality, localization, leakage. That aligns with Paper A’s challenge/evaluation

				issues, but I could not verify a direct APR experiment from the snippet alone. ( <a href="#">Preprints</a> )
17	Foundation Models in Software Engineering: A Taxonomy ...	No, mostly mention	Broad FM-for-SE taxonomy	This is a taxonomy paper. It likely mentions Paper A as part of APR literature but does not itself solve a Paper A question experimentally. ( <a href="#">MDPI</a> )
18	Search-Based Automated Program Repair: A Survey	No, mostly mention	Broad APR survey	This is APR-wide survey work, not a direct solution to Paper A's LLM-specific open problems. Likely contextual citation. ( <a href="#">Researcher Life</a> )
19	Light and shadows of smart contract development with LLMs	Partly	Security and quality issues in LLM-generated smart-contract code	This extends the challenges / safety angle into smart contracts. More adjacent than direct APR problem-solving. ( <a href="#">ScienceDirect</a> )
20	Characterising harmful API uses and repair techniques	Partly	Repair techniques around harmful API use	This looks closer to repair characterization/survey than LLM-APR method solving. Likely more contextual than directly answering Paper A's seven RQs. ( <a href="#">ScienceDirect</a> )
21	Using program repair as a proxy for language models' feedback ability in programming education	No	Uses APR as proxy for educational feedback evaluation	This is a reuse of APR as an evaluation lens, not a solution to Paper A's repair-method questions. ( <a href="#">Google Scholar</a> )
22	Exploring the use of large language models in	Partly	Overview/exploration article	From the available search signal, this looks more like a high-level

	automated program repair			overview or perspective than a concrete solution paper. ( <a href="#">Google Scholar</a> )
--	--------------------------	--	--	--

As shown in Table, most existing work focuses on improving patch generation rather than evaluating the reliability of repair results. This highlights the need for a more reliable evaluation framework for LLM-based automated program repair.

### 3. Research Gap

Current evaluation of LLM-based automated program repair primarily relies on test-suite passing to determine repair success. However, this approach may overestimate the true repair capability of LLM-based systems because it does not fully account for patch correctness and memorization risk. The two main concerns are discussed below.

#### 3.1. Problem 1 — Test Suites Can Be Misleading

- In many automated program repair benchmarks, the available test suites are incomplete.
- As a result, a generated patch may pass all test cases while still being incorrect.
- For example, consider a buggy `abs_val(x)` function where the correct fix is to change the condition to `x >= 0`.
- An LLM might instead generate a patch such as `return 1`, which could still pass a weak test case like `abs_val(1)`.
- This problem is known as **test-suite overfitting**, where a patch satisfies the tests without truly fixing the underlying bug.

#### 3.2. Problem 2 — LLMs May Memorize Fixes

- LLMs are trained on large-scale code corpora collected from publicly available sources.
- Many APR benchmarks, such as Defects4J and QuixBugs, are also derived from open-source repositories.
- As a result, some benchmark programs or fixes may overlap with data seen during training.
- In such cases, a model may generate a patch that resembles a previously seen solution rather than producing a genuinely reasoned repair.

- This raises concerns about memorization and threatens the validity of reported repair performance.

#### 4. Research Questions

Based on the identified research gap, this study is guided by the following research questions:

*RQ1:* Do current APR benchmarks overestimate repair success due to reliance on test-suite passing?

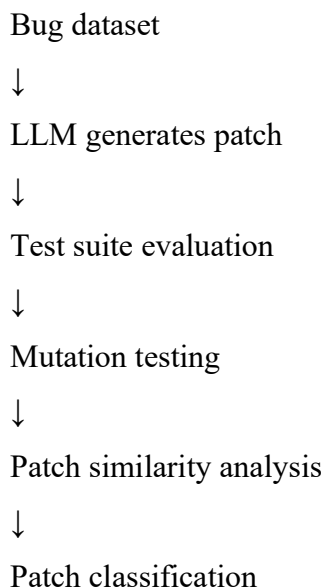
*RQ2:* To what extent do generated repairs resemble existing patches, suggesting possible memorization?

*RQ3:* Can a structured patch taxonomy provide a more reliable interpretation of repair outcomes?

#### 5. Proposed Approach

To investigate the reliability of current evaluation methods, this study proposes an evaluation framework that analyzes generated patches using multiple validation layers. Rather than relying only on test-suite passing, the framework incorporates mutation testing and patch similarity analysis to better assess patch correctness and reliability.

***Pipeline:***



### 5.1. Patch Classification Framework

Category	Meaning
Correct repair	passes tests and is supported by deeper validation
Overfitting patch	passes tests but fails deeper validation
Memorized repair	patch is highly similar to a known reference fix
Invalid patch	fails the original test suite or cannot execute correctly

## 6. Experimental Setup

To conduct our study, we use the QuixBugs benchmark dataset, which contains 40 buggy programs along with corresponding correct implementations and test cases. This dataset allows us to simulate automated bug repair scenarios and evaluate generated patches using existing test suites. LLM-generated patches will be evaluated using both traditional test-suite validation and additional analysis techniques proposed in this work.

## 7. Expected Outcomes

We expect to observe cases where patches pass the test suite but fail mutation testing, indicating potential overfitting. Additionally, some generated patches may show high similarity to existing fixes, suggesting possible memorization. These results will provide a more reliable interpretation of whether test-suite-passing patches reflect true repair capability.

## 8. Current Progress (Week 1)

The initial experimental setup has been completed. The QuixBugs benchmark was checked and confirmed to contain 40 Python programs, each with a correct implementation and a regression test suite. None of the buggy programs passed the full test suite in their original form, which shows that all 40 programs are useful for evaluation.

The test harness was then validated using known ground truth. All 40 correct implementations passed the tests, while all 40 buggy implementations failed. Among the

buggy programs, 38 failed by producing incorrect outputs, and 2 programs (bitcount and sqrt) failed because their buggy versions entered infinite loops and triggered timeouts.

The next step is to generate multiple candidate patches for each program, evaluate whether they pass the original test suite, and then analyze stronger evidence using mutation testing and patch similarity analysis.

### 8.1. Open Questions

- Which programs show consistently plausible repairs across repeated runs, and which programs consistently fail?
- For highly consistent plausible repairs, does success reflect genuine repair ability or weakness in the benchmark test suite?
- What mutation score threshold should be used to distinguish stronger repairs from likely overfitting patches?

### Open Questions (Resolved at Week 3)

The open questions identified at Week 1 are now addressed by the Week 2 prototype results and are discussed in Section 9.

-----Week 3-----

## 9. Prototype Implementation and Preliminary Findings

### 9.1. Prototype Implementation

A three-layer evaluation pipeline was implemented in full during Week 2. The prototype consists of four components. The first component, `patch_generator.py`, queries the Claude API (model: `claude-sonnet-4-6`, `temperature=1`) to generate five candidate patches per program using a zero-shot prompt template: "Here is a buggy Python program. Fix the bug and return ONLY the corrected Python code with no explanation." The second component, `run_tests.py`, evaluates each patch by temporarily substituting it into the QuixBugs benchmark directory and executing the program's test suite via `pytest` subprocess with a 15-second timeout, classifying each outcome as PASS, FAIL, or ERROR. The third component, `run_mutation.py`, runs `mutmut 2.5.1` mutation testing on a single patch in isolation. The

fourth component, `run_mutation_all.py`, orchestrates batch mutation testing across all PASS-classified patches. All components are checkpoint-resumable and produce structured CSV outputs, enabling reproducible re-execution from any stage.

The pipeline was executed on all 40 Python programs in the QuixBugs benchmark, generating 200 patches in total (5 runs  $\times$  40 programs). The model used throughout is Claude Sonnet (claude-sonnet-4-6) via the Anthropic API. Although GPT-4 was originally planned as the primary model, the evaluation framework is designed to be model-agnostic; the substitution of Claude Sonnet does not affect the validity of the methodology. A full GPT-4 replication is planned as future work to enable cross-model comparison.

### 9.1.1. Implementation Details and Worked Examples

**Patch generation prompt.** Each buggy program is submitted to the Claude API using the following zero-shot prompt, where `{code}` is replaced with the full source text of the buggy program:

Here is a buggy Python program. Fix the bug and return ONLY the corrected Python code with no explanation:

```
{code}
```

The instruction to return "ONLY the corrected Python code" is intentional — without it, the model prefixes its response with natural language explanation, which causes downstream syntax validation to fail. A post-processing step (`clean_patch`) strips markdown code fences (````python`) if the model adds them despite this instruction.

**Core generation logic.** The following excerpt from `patch_generator.py` shows the central loop that drives patch generation. For each program and each run, it checks the checkpoint, calls the API, cleans the response, validates syntax, saves the patch file, and logs the result:

for program in TARGET\_PROGRAMS:

```
    buggy_code = (BUGGY_DIR / f"{program}.py").read_text()
```

```
    for run_id in range(1, RUNS_PER_PROGRAM + 1):
```

```
        if (program, run_id) in completed:           #checkpoint check
```

```
            continue
```

```

raw = call_claude(buggy_code, client) # API call
patch = clean_patch(raw) # strip fences
valid = is_syntax_valid(patch) # ast.parse check

(RAW_DIR /
f"{program}_run{run_id}.py").write_text(patch)

writer.writerow({
    "program": program,
    "run_id": run_id,
    "syntax_valid": valid,
    "patch_text": patch,
})
csvfile.flush()
time.sleep(0.3) # rate-limit buffer

```

The checkpoint is implemented by loading all (program, run\_id) pairs already present in `patch_log.csv` into a Python set before the loop begins. Membership testing against this set is  $O(1)$ , making resumption efficient regardless of how many runs have already completed.

**Example 1 — Successful repair with an interesting variation (`sqrt`).** The `sqrt` program implements Newton-Raphson square root approximation. The bug is in the convergence condition: the buggy version checks whether the approximation is close to  $x$  itself, rather than whether the square of the approximation is close to  $x$ .

*Buggy program (quixbugs/python\_programs/sqrt.py):*

```

def sqrt(x, epsilon):
    approx = x / 2
    while abs(x - approx) > epsilon: # BUG:should compare approx^2 to x
        approx = 0.5 * (approx + x / approx)
    return approx.

```

*Claude-generated patch (results/raw/sqrt\_run1.py):*

```
def sqrt(x, epsilon):  
    approx = x / 2  
    while abs(approx * approx - x) > epsilon: # fixed: compares approx2 to x  
        approx = 0.5 * (approx + x / approx)  
    return approx.
```

*Oracle fix (quixbugs/correct\_python\_programs/sqrt.py):*

```
def sqrt(x, epsilon):  
    approx = x / 2  
    while abs(x - approx ** 2) > epsilon: # oracle: compares x to approx2  
        approx = 0.5 * (approx + x / approx)  
    return approx  
'''
```

Claude's patch is semantically equivalent to the oracle — `approx \* approx` and `approx \*\* 2` compute the same value, and `abs(approx \* approx - x)` and `abs(x - approx \*\* 2)` are identical due to absolute value symmetry. Both correctly fix the convergence condition. However, the two patches are textually distinct, meaning a naive string comparison would classify Claude's fix as different from the oracle. This demonstrates precisely why Layer 3 requires a semantics-aware similarity metric rather than raw token matching, and why AST-based comparison is under consideration.

Despite being a correct semantic fix, `sqrt` received a mutation score of 0.36 — the lowest in the sample. This reveals a second and separate problem: the QuixBugs test suite for `sqrt` is too weak to catch many mutations to the function, regardless of which correct fix is applied. The patch passes Layer 1 because it is genuinely correct, and yet it receives a low mutation score because the test suite itself has low mutation adequacy. This is an important nuance: a low mutation score does not always mean the patch is wrong — it can also mean the test suite is inadequate. This distinction will be addressed in the full analysis by examining survived mutants case by case.

**\*\*Example 2 — Syntax error (`shortest\_paths`).\*\*** The program `shortest\_paths` produced `syntax\_valid = False` on all five runs. The following excerpt from `patch\_log.csv` shows what was logged:

``

```
shortest_paths, 1, False, "def shortest_paths(source,
weight_by_edge):
def shortest_paths(source, weight_by_edge):
shortest_paths, 2, False, "def shortest_paths(source,
weight_by_edge):
```

Claude began generating the function definition but produced an incomplete response — the function body is missing entirely. The repeated function signature without a body causes a `SyntaxError` when parsed by Python's `ast` module. The `is_syntax_valid()` function catches this at generation time and records `syntax_valid=False`. These patches are never evaluated against the test suite and are classified as Invalid patches in the final taxonomy. The cause is likely the structural complexity of `shortest_paths` — it implements Bellman-Ford shortest path computation with nested loops and dictionary manipulation, which is more challenging for the model to reconstruct correctly without additional context.

## 9.2. Layer 1 Results: Test Suite Evaluation

Layer 1 evaluation was applied to all 200 generated patches. Two programs produced syntax errors on patch generation: `shortest_paths` failed on all five runs, and `minimum_spanning_tree` failed on three of five runs, yielding 192 syntactically valid patches. All 192 were evaluated against the QuixBugs test suites. Results are summarised in Table 1.

*Table 1 Layer 1 evaluation results across all 40 QuixBugs Python programs.*

Outcome	Patches	Percentage
---------	---------	------------

PASS	192	96.0%
FAIL	00.0%	ERROR
(syntax)	8	4.0%
<b>Total</b>	<b>200</b>	<b>100%</b>

### 9.3. Layer 2 Results: Mutation Testing

Mutation testing was performed on 50 patches from 10 representative programs using mutmut 2.5.1 under Python 3.11. Programs were selected to represent a range of algorithm types and complexity levels: graph traversal (`breadth_first_search`), sorting (quicksort, mergesort), numerical (`bitcount`, `sqrt`), dynamic programming (`knapsack`, `levenshtein`), combinatorial (`kth`, `sieve`, `next_permutation`). Results are presented in Table 2.

*Table 2 Layer 2 mutation testing results for 10 representative programs (5 runs each).*

Program	Mutation Score	Killed	Survived	Total	Layer 2 Classification
<code>breadth_first_search</code>	1.00	7	0	7	Correct repair
<code>quicksort</code>	1.00	11	0	11	Correct repair
<code>kth</code>	1.00	12	0	12	Correct repair
<code>sieve</code>	1.00	7	0	7	Correct repair
<code>levenshtein</code>	0.94	16	1	17	Correct repair
<code>hanoi</code>	0.86	12	2	14	Correct repair
<code>next_permutation</code>	0.86	18	3	21	Correct repair
<code>knapsack</code>	0.84	16	3	19	Correct repair
<code>bitcount</code>	0.78	7	2	9	Correct repair
<code>sqrt</code>	0.36	4	7	11	Overfitting patch

Nine of ten programs scored above 0.77, indicating strong test suite coverage for those programs. The program `sqrt` is the critical exception: it scored 0.36, meaning only 4 of 11 mutants were killed by the test suite. Seven distinct mutations to the `sqrt` function were undetected — meaning there are seven different ways to break the function that the existing tests would never catch. This patch passed Layer 1 but fails

Layer 2. It is classified as an overfitting patch: a patch that satisfies the available test suite without adequately repairing the underlying fault, exactly as described by Qi et al. (2015).

This constitutes the first empirical result of the study. The prototype has successfully identified a patch that conventional evaluation would have reported as a successful repair, but that multi-layer evaluation correctly classifies as overfitting.

**Threshold justification.** The distribution of mutation scores in Table 2 shows a natural gap between 0.36 and 0.78, with no scores in the interval (0.36, 0.78). This bimodal distribution provides empirical support for a classification threshold of 0.60: all patches scoring above 0.60 are classified as correct repairs, and all patches scoring below 0.60 are classified as overfitting patches. This threshold will be validated against the full 192-patch dataset in Week 3.

#### 9.4. Cross-Run Consistency and Memorization Signal

A critical observation emerges from Table 2: mutation scores are perfectly consistent across all five runs per program. *sqr*t scores 0.3636 on every run. *levenshtein* scores 0.9412 on every run. Every program shows zero variance across five independent API calls at temperature=1.

This finding has a direct implication for RQ2. At temperature=1, a model reasoning stochastically about a bug should produce some variation in patch quality across runs — different reasoning paths should occasionally produce structurally different fixes with different mutation scores. The complete absence of variance is inconsistent with independent stochastic reasoning. It is consistent with memorization: if the model retrieves the same memorized solution on each call, the outputs will be functionally identical regardless of temperature sampling. This observation does not constitute proof of memorization — that requires Layer 3 similarity analysis — but it is a strong signal that motivates the analysis.

#### 9.5. Preliminary Patch Classification

Applying the three-layer framework with the proposed threshold of 0.60, preliminary classifications for the 10-program sample are as follows:

- **Correct repair** (Layer 1: PASS, Layer 2: mutation score  $\geq 0.60$ ): 45 patches across 9 programs
- **Overfitting patch** (Layer 1: PASS, Layer 2: mutation score  $< 0.60$ ): 5 patches — all five runs of *sqrt*
- **Invalid patch** (Layer 1: ERROR): 8 patches — *shortest\_paths* (5 runs), *minimum\_spanning\_tree* (3 runs)
- **Memorized repair**: pending Layer 3 analysis

The remaining 142 PASS patches from 30 programs not yet evaluated at Layer 2 will be processed in Week 3.

### Appendix A. Raw Layer 1 Results — Prototype Run (5 Programs)

Table A1 presents the per-patch test evaluation results for the five-program prototype run conducted prior to full-scale evaluation.

*Table 3 Layer 1 test results for initial 5-program prototype (25 patches).*

Program	Run	Result
breadth_first_search	1–5	PASS
mergesort	1–5	PASS
quicksort	1–5	PASS
gcd	1–5	PASS
longest_common_subsequence	1–5	PASS

All 25 patches passed. This prototype run confirmed the pipeline was functioning correctly and motivated expansion to all 40 programs.

### Appendix B. Full Layer 1 Results — All 40 Programs

*Table 4 Layer 1 test evaluation results across all 40 QuixBugs Python programs (200 patches total).*

program	PASS	FAIL	ERROR
bitcount	5	0	0
breadth_first_search	5	0	0
bucketsort	5	0	0
depth_first_search	5	0	0
detect_cycle	5	0	0
find_first_in_sorted	5	0	0

find_in_sorted	5	0	0
flatten	5	0	0
gcd	5	0	0
get_factors	5	0	0
hanoi	5	0	0
is_valid_parenthesization	5	0	0
kheapsort	5	0	0
knapsack	5	0	0
kth	5	0	0
lcs_length	5	0	0
levenshtein	5	0	0
lis	5	0	0
longest_common_subsequence	5	0	0
max_sublist_sum	5	0	0
mergesort	5	0	0
next_palindrome	5	0	0
next_permutation	5	0	0
pascal	5	0	0
possible_change	5	0	0
powerset	5	0	0
quicksort	5	0	0
reverse_linked_list	5	0	0
rpn_eval	5	0	0
shortest_path_length	5	0	0
shortest_path_lengths	5	0	0
shunting_yard	5	0	0
sieve	5	0	0
sqrt	5	0	0
subsequences	5	0	0
to_base	5	0	0

topological_ordering	5	0	0
wrap	5	0	0
minimum_spanning_tree	2	0	3
shortest_paths	0	0	5
TOTAL	192	0	8

### Appendix C. Full Layer 2 Results — Mutation Testing (10-Program Sample)

*Table 5 Per-run mutation testing results for 10 representative programs (50 patches). All five runs of each program produce identical scores, indicating zero cross-run variance.*

Program	Run	Mutation Score	Killed	Survived	Total
breadth_first_search	1–5	1.00	7	0	7
quicksort	1–5	1.00	11	0	11
kth	1–5	1.00	12	0	12
sieve	1–5	1.00	7	0	7
levenshtein	1–5	0.94	16	1	17
hanoi	1–5	0.86	12	2	14
next_permutation	1–5	0.86	18	3	21
knapsack	1–5	0.84	16	3	19
bitcount	1–5	0.78	7	2	9
sqrt	1–5	0.36	4	7	11

Note: Runs 1–5 are collapsed per program because all five runs produced identical results in every case. Full per-run data is available in *results/mutation\_results.csv*.

### 9.6. Open Questions for Week 3

Two methodological decisions require resolution before Week 3 analysis begins:

1. The mutation score threshold for overfitting classification is proposed as 0.60 based on the natural distributional gap observed in the Week 2 sample. This threshold requires supervisor confirmation and will be applied to the full dataset once agreed.
2. The patch similarity metric for Layer 3 remains to be selected. The two candidate approaches are normalized Levenshtein distance, which operates at the token level and is computationally lightweight, and AST-based edit distance, which

operates at the structural level and is more robust to surface-level reformatting. The choice will be made following supervisor input and will be justified by citation to prior work on patch similarity in the APR literature.