

JaCoText and Integrated Gradients Performance in CodeGen Reasoning Explanation Benchmarks

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 8 peer-reviewed papers addressing the following research question: What is the relative performance of JaCoText vs. Integrated Gradients in explaining reasoning steps of CodeGen models on BigCode benchmarks, measured by both interpretability F1 score and downstream. 7 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Applied Explainability for Large Language Models: A Comparative Study. Research question: What is the relative performance of JaCoText vs. Integrated Gradients in explaining reasoning steps of CodeGen models on BigCode benchmarks, measured by both interpretability F1 score and downstream task accuracy?.

2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.8/10.

3 Results

8 papers retrieved. 7 claims extracted; 0 independently verified. Quality review score: 3.8/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
Integrated Gradients consistently highlighted sentiment-bearing tokens such as adjectives, negations, and intensifiers (×	0.03
Integrated Gradients attributions aligned well with human intuition and remained consistent across multiple examples.	×	0.05
Attention Rollout frequently emphasised syntactic or structural tokens, including stopwords, punctuation, and positional	×	0.02
In several cases, sentiment-relevant words received comparatively lower attention weights in Attention Rollout, reducing	×	0.04
SHAP explanations, when successfully computed, identified sentiment-relevant input components but often appeared noisy a	×	0.03
SHAP explanations were less visually stable than IG outputs and required careful preprocessing and configuration to inte	×	0.04
Integrated Gradients provides clearer and more intuitive explanations for sentiment classification compared to Attention	×	0.15

References

- <http://arxiv.org/abs/2402.04177v3>
- <http://arxiv.org/abs/2604.15371v1>
- <http://arxiv.org/abs/2505.24480v1>