

# Model Size and HumanEval Score Stability Across Evaluation Protocols

Assignee Research

May 31, 2026

## Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: What is the correlation between model size (1B–175B parameters) and HumanEval score stability across different evaluation protocols (e.g., deterministic vs. probabilistic sampling). Large language models (LLMs) achieve strong performance across many natural language processing tasks, yet their decision processes remain difficult to interpret. This lack of transparency creates challenges for trust, debugging, and deployment in real-world systems. 7 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.8/10. 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 correlation between model size (1B–175B parameters) and HumanEval score stability across different evaluation protocols (e.g., deterministic vs. probabilistic sampling)?.

## 2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.8/10.

## 3 Results

15 papers retrieved. 7 claims extracted; 1 independently verified. Quality review score: 4.8/10.

## 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	×	0.03
In several cases, sentiment-relevant words received comparatively lower attention weights in Attention Rollout, reducing	×	0.05
SHAP explanations, when successfully computed, identified sentiment-relevant input components but often appeared noisy a	×	0.04
SHAP explanations were less visually stable than IG outputs and required careful preprocessing and configuration to inte	×	0.05
Integrated Gradients provides clearer and more intuitive explanations for sentiment classification compared to Attention	✓	0.17

## References

- <http://arxiv.org/abs/2412.21199v2>
- <http://arxiv.org/abs/2410.12381v3>
- <http://arxiv.org/abs/2604.15371v1>