

Cross-Scale Alignment Consistency of the BNRM Framework Versus Bradley-Terry Models on the Chatbot Arena Leaderboard

Assignee Research

June 13, 2026

Abstract

Foundation models (FMs), particularly large language models (LLMs), have shown significant promise in various software engineering (SE) tasks, including code generation, debugging, and requirement refinement. Despite these advances, existing evaluation frameworks are insufficient for assessing model performance in iterative, context-rich workflows characteristic of SE activities. To address this limitation, we introduce *SWE-Arena*, an interactive platform designed to evaluate FMs in SE tasks. *SWE-Arena* provides a transparent, open-source leaderboard, supports multi-round conversational w

1 Introduction

This paper examines: *SWE-Arena: An Interactive Platform for Evaluating Foundation Models in Software Engineering*. Research question: Does the BNRM framework maintain consistent alignment performance across different language model scales (e.g., Qwen2.5-7B vs. Qwen2.5-72B) when evaluated using the Chatbot Arena leaderboard, and how does this compare to standard Bradley-Terry models?.

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 7.6/10.

3 Results

15 papers retrieved. 14 claims extracted; 11 independently verified. Quality review score: 7.6/10.

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
SWE-Arena is an interactive platform designed to evaluate Foundation Models in Software Engineering tasks.	✓	0.23
SWE-Arena provides a transparent, open-source leaderboard.	✓	0.20
SWE-Arena supports multi-round conversational workflows.	✓	0.19
SWE-Arena enables end-to-end model comparisons.	✓	0.20
SWE-Arena introduces a metric called 'model consistency score' that measures the consistency of model outputs through se	✓	0.25
SWE-Arena introduces a metric called 'conversation efficiency index' that evaluates model performance while accounting f	✓	0.29
SWE-Arena incorporates a feature called RepoChat.	×	0.14
RepoChat automatically injects repository-related context, such as issues, commits, and pull requests, into the conversa	✓	0.27
Foundation models have advanced software engineering tasks including code generation, debugging, and requirement refinem	✓	0.23
SWE-Arena plans to analyze user-submitted requests to identify common patterns and challenges in software engineering ta	✓	0.21
SWE-Arena intends to enable users to vote on model performance.	×	0.11
SWE-Arena aims to broaden coverage to include domain-specific models and multimodal foundation models.	✓	0.22
SWE-Arena aims to support infrastructure for web browsing and API integration.	×	0.14
SWE-Arena intends to incorporate advanced context compression methods such as LongRope or SelfExtend.	✓	0.21

References

- <http://arxiv.org/abs/2502.01860v5>

- <http://arxiv.org/abs/2306.09265v1>
- <http://arxiv.org/abs/2407.04065v4>