

Systematic Analysis of Protocol Factors Driving Extreme Rouge-L Variance in GPT-3.5 Across Divergent Evaluation Domains

Assignee Research

June 11, 2026

Abstract

Large language models (LLMs) are gaining increasing popularity in both academia and industry, owing to their unprecedented performance in various applications. As LLMs continue to play a vital role in both research and daily use, their evaluation becomes increasingly critical, not only at the task level, but also at the society level for better understanding of their potential risks. Over the past years, significant efforts have been made to examine LLMs from various perspectives. This paper presents a comprehensive review of these evaluation methods for LLMs, focusing on three key dimensions:

1 Introduction

This paper examines: A Survey on Evaluation of Large Language Models. Research question: Reproducibility meta-analysis: 4 independent publications report divergent GPT-3.5 performance on Rouge-L with a 78.2 percentage-point spread (range 1.8%–80.0%). Source papers: "ACE-RLHF: Automated Code Evaluation and Socratic Feedback Generation Tool using\ldots{" (2025, 1.8%); "Automated Literature Review Using NLP Techniques and LLM-Based Retrieval-Augmen\ldots{" (2024, 18.1%); "A Video Is Worth 4096 Tokens: Verbalize Videos To Understand Them In Zero Shot" (2023, 18.9%); "Comparison of Open-Source and Proprietary LLMs for Machine Reading Comprehensio\ldots{" (2024, 80.0%). Preliminary analysis suggests: The extreme variance likely stems from evaluating fundamentally different tasks under the same "Rouge-L" label, where the 80.0% score reflects a standard text-based Machine Reading Comprehension benchmark while the <20% scores originate from cross-modal video-to-text generation or specialized code evaluation domains w\ldots{" Systematically evaluate which evaluation protocol factors (model configuration, inference setup, quantization, tokenization, few-shot count, metric interpretation, or data-split selection) best explain the observed spread; identify the highest-confidence expla-

nation supported by each paper's stated methodology; and assess whether the highest-reported score is reproducible under the conditions described by the lowest-reporting paper..

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.0/10.

3 Results

13 papers retrieved. 10 claims extracted; 10 independently verified. Quality review score: 8.0/10.

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
Large language models (LLMs) are gaining increasing popularity in both academia and industry, owing to their unprecedented	✓	0.33
As LLMs continue to play a vital role in both research and daily use, their evaluation becomes increasingly critical, no	✓	0.39
Over the past years, significant efforts have been made to examine LLMs from various perspectives.	✓	0.25
This paper presents a comprehensive review of these evaluation methods for LLMs, focusing on three key dimensions: what	✓	0.33
We provide an overview from the perspective of evaluation tasks, encompassing general natural language processing tasks,	✓	0.41
We answer the 'where' and 'how' questions by diving into the evaluation methods and benchmarks, which serve as crucial c	✓	0.31
We summarize the success and failure cases of LLMs in different tasks.	✓	0.23
We shed light on several future challenges that lie ahead in LLMs evaluation.	✓	0.26
Our aim is to offer invaluable insights to researchers in the realm of LLMs evaluation, thereby aiding the development o	✓	0.32
Our key point is that evaluation should be treated as an essential discipline to better assist the development of LLM.	✓	0.29

References

- <https://doi.org/10.1093/jamia/ocae073>
- <https://doi.org/10.48550/arxiv.2402.06196>
- <https://doi.org/10.48550/arxiv.2307.03109>