

# Cross-Domain Adaptation of RAGalyst: Consistency in Human-Aligned Evaluation Metrics

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does cross-domain adaptation of RAGalyst (e.g., from religious texts to legal or medical domains) affect the consistency of human-aligned evaluation metrics compared to domain-specific fine-tuning. 18 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: RAGalyst: Automated Human-Aligned Agentic Evaluation for Domain-Specific RAG. Research question: How does cross-domain adaptation of RAGalyst (e.g., from religious texts to legal or medical domains) affect the consistency of human-aligned evaluation metrics compared to domain-specific fine-tuning?.

## 2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.0/10.

## 3 Results

14 papers retrieved. 18 claims extracted; 6 independently verified. Quality review score: 6.0/10.

## 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
RAGalyst is an automated, human-aligned agentic framework designed for the evaluation of domain-specific RAG systems.	✓	0.39
RAGalyst features an agentic pipeline that generates synthetic question-answering (QA) datasets from source documents.	✓	0.25
RAGalyst incorporates an agentic filtering step to ensure data fidelity in generated datasets.	✓	0.16
RAGAS is adopted as the primary baseline for evaluating the performance of RAGalyst because it is described as the only	×	0.05
Existing end-to-end RAG evaluation frameworks rely on some degree of manual validation for QA dataset generation.	×	0.08
Rule-based metrics in RAG evaluation often fail to capture subtle semantic nuances.	×	0.10
LLM-based metrics in existing works are rarely benchmarked for alignment with human judgment.	×	0.12
AutoCalibrate auto-tunes evaluation prompts to improve alignment with human judgment.	×	0.07
DSPy is a declarative framework that enables the programmatic creation and refinement of prompts for LLMs.	×	0.04
The PoLL evaluation framework advocates using a panel of various LLM evaluators to reduce bias and variance in generatio	×	0.06
High-quality QA datasets are often unavailable or insufficient in specialized domains.	×	0.09
Alberti et al. pioneered a round-trip consistency approach combining answer extraction, question generation, and answer	×	0.06
Shakeri et al. proposed an end-to-end transformer-based generator that outputs both the question and the answer from a g	×	0.05
RAGalyst was applied to evaluate RAG components across three domains: military operations, cybersecurity, and bridge eng	✓	0.19
No single embedding model, LLM, or hyperparameter configuration proved universally optimal across the tested domains.	✓	0.17
RAGalyst uses prompt optimization to achieve a strong correlation with human annotations for Correctness and Answerabili	✓	0.21
Modern Large Language Models suffer from hallucinations, defined as generating content that appears plausible but is fac	×	0.04
RAGalyst is available on Github.	×	0.05

## References

- <http://arxiv.org/abs/2402.12317v2>
- <http://arxiv.org/abs/2412.01496v2>
- <http://arxiv.org/abs/2511.04502v1>