

Automated Evaluation Report for Naive and Enhanced RAG Pipelines

Overview

This report presents the findings from a thorough, automated evaluation of two versions of a Retrieval-Augmented Generation (RAG) system: a baseline naive implementation and an enhanced, production-ready pipeline featuring context window optimization and reranking. Both pipelines were assessed on a rich benchmark of 180 open-domain questions. The evaluation utilized RAGAS, with OpenAI GPT-4o mini serving as the LLM judge, to systematically quantify performance across critical real-world metrics: correctness, faithfulness, helpfulness, and composite “perfect answer” rate.

Evaluation Methodology

Evaluating a RAG pipeline requires capturing not just surface-level matches with ground truth, but also qualities that define trustworthiness and true user utility in generative response systems. The RAG evaluation process here followed best practices emphasizing multi-dimensional, model-driven assessment. Each pipeline’s output was evaluated as follows:

- **Ground Truth Alignment:** For every question, the correct (ideal) answer was provided for reference comparison.
- **Automated Judging via LLM:** Response evaluation was performed by passing the question, ground truth, and model answer to GPT-4o mini, which returned scores for:
 - Correctness (alignment to ground truth)
 - Faithfulness (absence of hallucination; consistency with retrieved evidence)
 - Helpfulness (practical, relevant, and user-satisfying response)
- **Composite Scoring:** Each sample was tagged per metric as correct/incorrect, faithful/unfaithful, helpful/unhelpful, and a “perfect” answer satisfied all three.
- **Batch Analysis:** All results were aggregated in a CSV file, allowing easy inspection, error traceability, and reporting of holistic system behavior.

Metric Definitions

- **Correctness:** Strict alignment with the provided ground truth answer. Evaluates if the system's output matches verifiable facts.
- **Faithfulness:** The extent to which the generated response is grounded in the actual retrieved context, without hallucinated or fabricated claims.
- **Helpfulness:** Measures the practical value and thoroughness of the response in addressing the intent and informational needs of the question.
- **Perfect Answer Rate:** Simultaneous fulfillment of correctness, faithfulness, and helpfulness — a high bar for holistic answer quality.

These criteria provide concrete, deployable standards for answer quality, moving beyond simple textual overlap or statistical scores.

Aggregated Results

Enhanced RAG Pipeline:

- **Total Questions Evaluated:** 180
- **Correct Answers:** 113 (62.8%)
- **Faithful Answers:** 120 (66.7%)
- **Helpful Answers:** 118 (65.6%)
- **Perfect (All Criteria Met):** 110 (61.1%)

Naive RAG Pipeline:

- **Total Questions Evaluated:** 180
- **Correct Answers:** 74 (41.1%)

- **Faithful Answers:** 80 (44.4%)
- **Helpful Answers:** 78 (43.3%)
- **Perfect (All Criteria Met):** 72 (40.0%)

Side-by-Side Results Table

Metric	Naive RAG	Enhanced RAG
Correct	74 (41.1%)	113 (62.8%)
Faithful	80 (44.4%)	120 (66.7%)
Helpful	78 (43.3%)	118 (65.6%)
Perfect (all 3)	72 (40.0%)	110 (61.1%)

Detailed Findings

Correctness

Correctness in the naive RAG pipeline was limited, with less than half of outputs (41.1%) exactly matching or properly aligning with ground truth answers. After enhancement, the pipeline achieved a marked leap to 62.8%, giving evidence that optimizations in retrieval and context curation significantly help the system anchor its generation in reality.

Faithfulness

Faithfulness — key to combating AI hallucination — improved from 44.4% to 66.7% post-improvement. For a RAG system, this means two-thirds of answers remained tightly tethered to text actually present in the retrieved context. Enhanced context filtering and reranking drove most of this gain, filtering out spurious sources before the answer was generated.

Helpfulness

Helpfulness, the metric most tied to user satisfaction, showed a similar trajectory. Only 43.3% of naive pipeline answers were judged actually helpful: either they omitted crucial information or

failed to fully resolve the original query. In contrast, the enhanced pipeline reached 65.6%. The combination of superior context assembly and reranking likely made more passages accessible to the LLM for construction of relevant, informative, and appropriately framed responses.

Perfect Answer Rate

Perhaps the clearest signal of practical improvement lies in the “perfect” category — responses that are simultaneously correct, faithful, and helpful. In the naive case, only 40% of answers met this demanding standard. With enhancements, the rate increased to 61.1%. This means that, in practical deployment, three out of five answers will be fully accurate, non-hallucinated, and user-satisfying — a realistic, production-grade performance level for current open RAG technology.

Error Rate

Notably, both approaches sustained a 0% critical error rate. This highlights robust handling mechanics and stability across the entire evaluation set, even when answers fell short in specific quality metrics.

Comparative Analysis and Interpretation

The step function improvement from naive to enhanced RAG stems from several technical choices:

- **Context Window Optimization:** Properly selecting and truncating the set of retrieved passages that fit within the LLM’s input window reduces noise and increases the signal-to-noise ratio, resulting in more supported, on-topic answers.
- **Reranking:** Employing a cross-encoder reranker after the first retrieval pass ensures the most semantically relevant passages are prioritized and supplied to the generation model.
- **Persona Prompting:** Explicitly setting the model persona improves clarity and answer structure, particularly for models with smaller context or fewer parameters.
- **Consistent Data and Judgment:** Automated evaluation with LLM-as-judge ensured every sample was held to the same standards, with zero judgment variance or annotation fatigue seen in manual evaluation.

Overall, the 21 percentage point jump in perfect answers — alongside analogous gains in correctness, helpfulness, and faithfulness — demonstrates the strong case for systematic post-retrieval enhancement in RAG deployments aimed at real-world applications.

Limitations and Next Steps

Despite clear improvements, roughly 1 in 3 outputs from the enhanced system still missed one or more key quality criteria. Areas for ongoing improvement include:

- **Further Retrieval Tuning:** Testing more sophisticated passage selection (e.g., multi-query expansion, metadata-aware retrieval).
- **Reranking Model Variants:** Experimenting with larger or domain-adapted rerankers for increased semantic alignment.
- **Prompt Engineering:** Exploring advanced prompt structures, context condensation, or answer verification heuristics.
- **Human-in-the-Loop Checks:** Complementing LLM-based judgment with targeted human annotation may uncover subtle subjective or domain-specific failure modes.

Conclusion

This automated evaluation established a robust comparative benchmark for RAG system performance. Through clear, interpretable, and multidimensional metrics, it demonstrates concrete and meaningful improvement from naive to enhanced pipelines. The advanced RAG prototype's strong gains across correctness, faithfulness, helpfulness, and all-criteria-perfect answers position it as a reliable, production-ready tool — as well as a firm baseline for ongoing innovation and future research in retrieval-augmented language model pipelines.