



RAG Validation

By PK



RAG Validation components

Component	Goal	Potential Failures (and Root Cause)
1. Retrieval (Context)	Find all necessary and only relevant documents/chunks.	Missing Info: Retriever didn't find a relevant chunk (low Recall). Noise: Retriever brought back irrelevant chunks (low Precision).
2. Generation (Answer)	Use the provided context to generate an accurate, relevant, and fluent answer.	Hallucination/Inaccuracy: Answer contains facts not in the context (low Faithfulness/Groundedness). Irrelevance: Answer doesn't fully address the user's question (low Answer Relevance).



Step-by-Step Validation Workflow

Step-by-Step Validation Workflow

- Data Preparation: The Foundation
- Evaluation Criteria: What Makes a RAG Response "Good"?

1. Data Preparation: The Foundation

Validation is impossible without a good dataset.

- ❖ **Create a High-Quality Test Set:** Gather a set of **Questions**, **Expected Contexts** (the chunks that *should* be retrieved), and a **Ground Truth Answer** (the perfect response). **Establishing a "Golden" Reference Set**
- ✓ **Diversity is Key:** Include simple queries, complex multi-step questions, queries with misspellings, and **negative test cases** (questions that *cannot* be answered by the knowledge base).



Step-by-Step Validation Workflow

2. Evaluation Criteria: What Makes a RAG Response "Good"?

A. Retrieval Quality (Context)

Metric	What it Measures	How it's Used
Context Relevance/Precision	Is the retrieved context necessary to answer the question? (Signal-to-Noise Ratio)	Checks if a large portion of the retrieved chunks are noise/irrelevant.
Context Recall	Did the retriever find all the information required to generate the complete answer?	Crucial for multi-part or complex questions. (Requires a ground truth answer to verify if all facts are present in the context).
MRR (Mean Reciprocal Rank)	How quickly was the first most relevant document/chunk returned?	Important for systems where the highest-ranked result is most critical.



Step-by-Step Validation Workflow

2. Evaluation Criteria: What Makes a RAG Response "Good"?

B. Generation Quality (Answer)

Metric	What it Measures	How it's Used
Faithfulness (or Groundedness)	Is the generated answer factually consistent <i>with the retrieved context</i> ? (Anti-Hallucination)	The most important RAG metric. A faithful answer does not "hallucinate" new information.
Answer Relevance	Is the final answer <i>directly relevant</i> to the user's question?	Checks for tangential answers or incomplete responses.
Answer Correctness/Semantic Similarity	How semantically similar is the generated answer to the human-labeled Ground Truth Answer ?	A measure of overall accuracy, usually requiring a reference answer.



RAGAS Framework

RAGAS (Retrieval-Augmented Generation Assessment Suite) is a popular open-source framework that automates the calculation of RAG metrics, often using an LLM-as-a-judge approach to avoid heavy manual labeling.

RAGAS Metric	Component Focused	How it's Calculated (LLM-as-a-Judge)
1. Faithfulness	Generation	An LLM-Judge breaks the answer into statements and checks if each statement is logically supported by the retrieved context.
2. Answer Relevancy	Generation	An LLM-Judge generates a set of potential questions for the RAG answer. It then measures the similarity between these generated questions and the original user question.
3. Context Precision	Retrieval	An LLM-Judge rates the relevance of each retrieved context chunk to the question (a signal-to-noise ratio).
4. Context Recall	Retrieval	An LLM-Judge compares statements in the Ground Truth Answer and checks if all necessary facts are present in the Retrieved Contexts. (Note: This metric requires a Ground Truth Answer to work).

Examples

User query:

“What does Section 4.3 of the policy describe?”

Ground truth relevant chunks:

There are **2 chunks** in the whole dataset that contain Section 4.3 details:

- Chunk_12
- Chunk_14

Rank	Retrieved Chunk	Relevant?
1	Chunk_88	✗
2	Chunk_12	✓
3	Chunk_152	✗
4	Chunk_14	✓
5	Chunk_201	✗

Example 1: Recall@K Calculation

Assume your vector retriever returns top K=5 chunks as above :

Relevant items retrieved = **2** (Chunk_12, Chunk_14)

Total relevant items in dataset = **2**

➤ **Recall@5 = 2 / 2 = 1.0 (100%)**

This means your retriever found **all** relevant chunks within top-5 → excellent.

Example 2: Precision@K Calculation

Same top-5 retrieval:

Relevant items = **2**

Retrieved total K = **5**

➤ **Precision@5 = 2 / 5 = 0.4 (40%)**

Only **40%** of the retrieved chunks were relevant → retrieval is noisy.

Metric	Value	Interpretation
Recall@5	100%	Your retriever found all relevant chunks
Precision@5	40%	But it also retrieved many irrelevant chunks

❖ Interpretation for Example 1 + 2

Reranking or CRAG filtering is needed to improve precision.



RAG system's health

Metric Score	Interpretation	Action to Improve
Low Faithfulness	You have a Hallucination problem. The LLM is going off-script.	Improve your prompt template (be stricter), or try a smaller chunk size (to reduce noise in the context window).
Low Context Recall	You have a Comprehensiveness problem. Your retriever is missing key facts.	Check your chunking strategy (maybe overlap more, or use larger chunks), or refine your embedding model/retrieval parameters .
Low Context Precision	You have a Noise problem. Too many irrelevant chunks are retrieved.	Implement a C-RAG, re-ranker after initial retrieval, or use a different embedding model that captures relevance better.
Low Answer Relevancy	You have a Focus problem. The answer is vague or answers a different question.	Tune your prompt template to be more direct, or adjust your LLM's temperature (lower temperature for more focused answers).



Thank you