# Steps

1) RAG → DATA Ingestion, Retriever, Generation ✓

2) TEST DATA ⇒ Question ↔ Answer

3) Evaluation Metrics [LLM as a judge]

# Your Data



**Database**
Structured

**Documents**
Unstructured

**API**
Programmatic

Load →

## Index
Vector embeddings

Chicken

Wolf

Dog

Cat

[0.34, 2.35, 8.34, ... ]
300 dimensions

Banana

Apple

Query 💬

User

Prompt + query
+ relevant data →

response

## LLM 🤖
Model like GPT, llama

# **Correctness**: Response vs reference answer
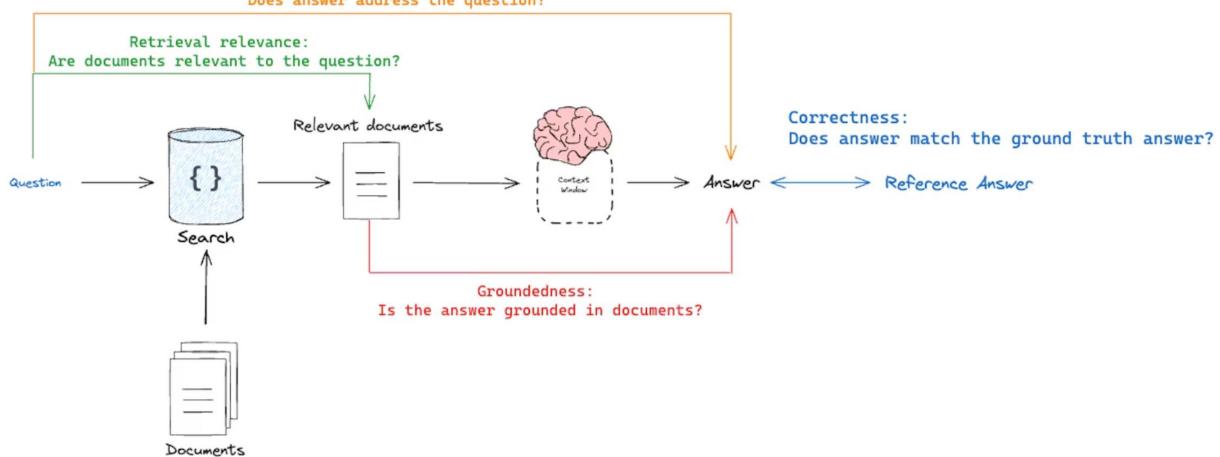
- Goal: Measure "*how similar/correct is the RAG chain answer, relative to a ground-truth answer*"

- Mode: Requires a ground truth (reference) answer supplied through a dataset

- Evaluator: Use LLM-as-judge to assess answer correctness.

# **Relevance**: Response vs input

- Goal: Measure "*how well does the generated response address the initial user input*"

- Mode: Does not require reference answer, because it will compare the answer to the input question

- Evaluator: Use LLM-as-judge to assess answer relevance, helpfulness, etc.

# **Groundedness**: Response vs retrieved docs

- Goal: Measure "*to what extent does the generated response agree with the retrieved context*"

- Mode: Does not require reference answer, because it will compare the answer to the retrieved context

- Evaluator: Use LLM-as-judge to assess faithfulness, hallucinations, etc.

# **Retrieval relevance**: Retrieved docs vs input

- Goal: Measure "*how relevant are my retrieved results for this query*"

- Mode: Does not require reference answer, because it will compare the question to the retrieved context

- Evaluator: Use LLM-as-judge to assess relevance