

Qdrant → Vector DB

Vector DB

Why? Metadata Filtering

Memory

In memory (RAM) (No Disk)  
Available in Runtime Only

{ col 1  
  { col 2  
    { col 3  
      { col 4

Collection :- Vector Dim (786)  
Distance Cosine Distance

Vec 1  
786

Vector 2

1536

RAG Evaluation Metrics

Metrics → Frameworks

RAGAS → RAG Assessment

\* DeepEval \*

## \* Evaluation Dataset \*

### Metrics

Faithfulness :- Generated content stick to the retrieved facts.

Check if the LLM is hallucinating

Faithful → Source Material.

### Examples

Input :- When was the first super bowl?

Response :- Jan 15, 1967

Actual

LLM :- The first super bowl hosted on Jan 15, 1967  
in Los Angeles.

Response

Retrieved Text :- The first AFL - NFL  
World Championship was  
played on Jan 15, 1967.  
in LA.

- 1) Extract all claims from generated answer
- 2) Check if each claim is supported by retrieved context.

3) Scoring Mechanism

$$\text{Score} = \frac{\text{Number of claims supported by context}}{\text{Total claims}}$$

	Claims	Score
1)	It was Jan 15, 1964 (supported)	
2)	It was in LA (supported)	
3)	It was a sunny day. (not supported) not in context	

$$\text{Score} = \frac{2}{3} = 0.67$$

Score Interpretations

1 : Perfect (All claims have R.C)

0 : Answer made up by LLM.

## RAG Metrics Mental Model

User Question



Retriever

→ Context Precision

→ Context Recall

→ Context Entity Recall



Retrieved  
Thunks



Generation  
LLM

→ Faithfulness

→ Noise Sensitivity



Response

→ Answer Relevancy

(address the question)

## Response Answer Relevancy

Measures :- Address whether the answer actually answering the question

Q1 What is the capital of France?

(Totally True/False  
but factual  
partially true)

R1 The Eiffel Tower is beautiful.

R2. Paris is a beautiful city

## Working :-

- 1) Generate hypothetical questions that the answer would be suitable.
- 2) Compare the generated questions with the original questions (Embedding)
- 3) \* Higher Similarity = more relevant answer  
\* Lower Similarity = less relevant answer

## Answer Relevancy

- 1 : Answers directly to the addressed question.
- 0 : Completely off-topic answers
- 0-1 : Partially off-topic.

## Good Relevancy

User Input :- When was the 1st S.B?

Response :- The first S.B was held on Jan 15, 1967.

R-1 :- The first AFL-NFL championship was played on 15<sup>th</sup> Jan, 1967 in LA

R2 :- The S.B is the annual championship game of NFL.

✓ What | Where X

Score :- (0.35)

## Context Precision

Where is Eiffel Tower?

The E.T located in Paris

R.C :- i)

The ET is located Paris

X ii) Paris is the capital of France

X iii) The ET was built in 1889.

X iv) Pizza originated Italy. \*

### Context Precision

Level :- Retrieval

Measure :- Are the retrieved chunks ranked properly?

(checks the top K highly relevant chunks)

$C_1 \rightarrow \frac{Q_1}{w_1}$  high  
 $C_2 \rightarrow w_2$  low  
 $C_3 \rightarrow w_3$  low  
 $C_4 \rightarrow w_4$  low

Working :-

- 1) For each retrieved chunk, determine whether it is relevant to the question
- 2) Calculate the precision at each position?  
(weightage by position)
- 3) Relevant chunks at the top = higher score  
Irrelevant chunk at the top = lower.

### Input

Where is ET located?

Reference: The ET is in Paris.

- R.C = [The ET is located in Paris] Highly ✓
- [Paris is the capital of France] Some relevant
- [The ET was built around 1889] IRRelevant
- [Pizza originated in Italy] Completely irrelevant

## Context Recall

Level:- Retrieval

Measures :- Did we retrieve all the necessary info to answer the question.

All chapters of book  $\rightarrow$  give an exam

Working :- individual

- 1) Breaking down of reference answer into claims
- 2) Check if the claim is attributed to the R.C
- 3) Score = 
$$\frac{\text{Claims found in the context}}{\text{Total claims in reference / response}}$$

## Example :-

Input :- Tell me about ET

Response :- The ET is in Paris, built in 1889.

R.C = [The ET is a landmark located in Paris, FR]

[The ET was completed in 1889]

Reference Claims :-

- 1) location ✓
- 2) Built ✓

$$Score = \frac{2}{2} = 1$$

### Score Interpretation

- 1 : - Retrieved contexts are provide the full answer
- 0 : R-C is completely useless.

### Context Recall

Measures :- Did we retrieve the context containing all necessary entities (people, place, org, date)

mentioned in the reference answer.

Analogy :-

Reference :- Paul was going to London on  
14 Jan, 2025.

3 Entities

Person :- Paul

Location :- London

Date :- 14 Jan, 2025

Extraction



Reference

Working

1) Extract Named Entity from Reference

2) Extract Name Entity from R-C

3) Score =  $\frac{\text{Entities in both}^{\text{common}}}{\text{Entities in reference}}$  (matching with R)

Example:-

Ref: Albert Einstein developed relativity in P.U  
in 1905.

R.C = [Albert Einstein was a famous physicist in P.U.]

A Entities in Ref :- [[A,B], [P,V], [Date]]

B Entities in R.C :- [[A,B], [P,V]]

$$\text{Score}:- \frac{2}{3} = 0.67$$

Noise sensitivity (unwanted)  
Jitter

Level :- Generation

Working:-

- ① Compare the reference answer with generated answer.
- ② Identify incorrect claims in the response.
- ③ Determine if the incorrect claims came from (irrelevant) unverified chunks noise
- ④ Score : 
$$\frac{\text{Incorrect Claims from noise}}{\text{Total Claims}}$$