

What is LLM Testing?

LLM testing means checking if a **Large Language Model (like ChatGPT)** is working correctly — just like we test any software, but here we test **AI-generated answers**.

What Was Before?

Before LLMs:

- We tested rules and logic (e.g., if X happens, then do Y).
- We used **unit tests**, **automation**, **API checks** — everything was **deterministic** (same input = same output).

Limitations:

- These worked well for traditional apps.
 - But they **don't work for AI**, where answers can **vary**, and correctness is not black or white.
-

What Changed Now?

With **LLMs**, output can vary:

- For one input, we can get different but valid answers.
- LLMs can **hallucinate**, miss context, or be biased.

So, we need **new testing methods** — this is called **LLM Testing & Evaluation**.

What is Evaluation in LLMs?

We now use **Evaluation Frameworks** to check:

- Is the answer **factually correct**?
 - Is it **relevant to the context**?
 - Is the **tone safe** and aligned with the brand?
 - Did it **follow instructions**?
-

✓ Evaluation Techniques

1. **LLM-based Evaluation** (e.g., RAGAS, DeepEval)
 - Uses another model to “grade” answers like a human.
 - Example: Is this answer faithful to the context?
 2. **Traditional (non-LLM) metrics**
 - BLEU, ROUGE, similarity scores.
 - Fast but not smart. Doesn't understand meaning.
 3. **Rubric-based Scoring** (Structured grading using specific rules)
 - E.g., “Give 1 point if answer is relevant, 1 if factual...”
-

★ Which Tools Are Best?

Tool	Best For	Why Use It
RAGAS	Testing RAG pipelines (search + LLM)	Open-source, scoring multiple metrics
DeepEval	Custom chatbots, use case-specific apps	Easy to plug in, supports GEval, safe content checks

Why It Matters Today?

- **AI is everywhere** — search, chatbots, copilots.
 - Without proper testing, **AI can break silently** (wrong answers, risky replies).
 - LLM testing **reduces business risk**, saves **tokens**, improves **quality**.
-

TL;DR (One Line Summary)

Traditional testing fails with AI — today we need LLM testing to ensure reliable, safe, and meaningful AI output using tools like RAGAS and DeepEval.





Would you like this in a PowerPoint slide format or a one-pager PDF to show to leadership?

You said:

Explain more LLM testing reduces business risk, saves tokens, improves quality. and how it save tokens?

1. LLM Testing Reduces Business Risk

Without testing, LLMs can go wrong silently, which can lead to:





-  **Wrong Answers** → Misleading customers, damaging brand trust.
 - *Example:* “What is the interest rate?” → Wrong value can cause compliance issues.
-  **Unsafe or biased responses** → Legal or PR risks.
 - *Example:* Offensive reply or unfair decision by AI bot.
-  **Incomplete workflows** → AI doesn’t follow the correct flow (e.g., skipping steps in a ticket system).
-  **Hallucinations** → AI makes up facts not in the data/context.

Testing helps catch these BEFORE they reach customers.

2. LLM Testing Saves Tokens (and Money)

LLM usage = **cost** based on **tokens** used (input + output).

How testing saves tokens:

-  **Catch over-verbose outputs:**
 - If your bot replies with unnecessary long answers, testing can help shorten them.
 - Shorter output = fewer tokens = lower cost.
-  **Improve Prompt Efficiency:**
 - Evaluation tools help optimize prompts to get better answers in fewer tries.
 - Less back-and-forth = lower token usage.
-  **Detect bad generations early:**
 - Without testing, you might ship a buggy version of the model or prompt that keeps generating wrong answers — leading to **repeated queries** and **more usage**.
 - Testing filters these out early.
-  **Better context precision:**
 - Tools like **RAGAS** check if the model uses the right documents.
 - If irrelevant context is being passed (adds token cost!), you can fix it.

Example Token Saving Scenario:






Let's say your chatbot handles 10,000 queries a day.

- Without testing: Avg 100 tokens per query = 1 million tokens
- With testing & optimization: Avg 75 tokens per query = 750,000 tokens
- **Saves 250,000 tokens/day**, which can mean **\$100s or \$1000s/month**



3. LLM Testing Improves Quality

LLM testing ensures:

-  **Accurate** responses (Faithfulness)
-  **Relevant** answers (Context precision)
-  **Concise and correct tone**
-  **Follows instructions** (e.g., "Answer in bullet points")
-  **No hallucinations or off-topic replies**

Tools like **DeepEval** and **RAGAS** give **automatic quality scores** per answer, so you can:

- Track performance over time
- Compare two model versions (A/B testing)
- Fine-tune prompts or models based on insights



Final Analogy for Leadership

Think of LLM Testing like **QA for AI brains**:

- It **catches errors**, **saves costs**, and ensures your AI is **safe and smart** before going live — just like how we test traditional code before release.

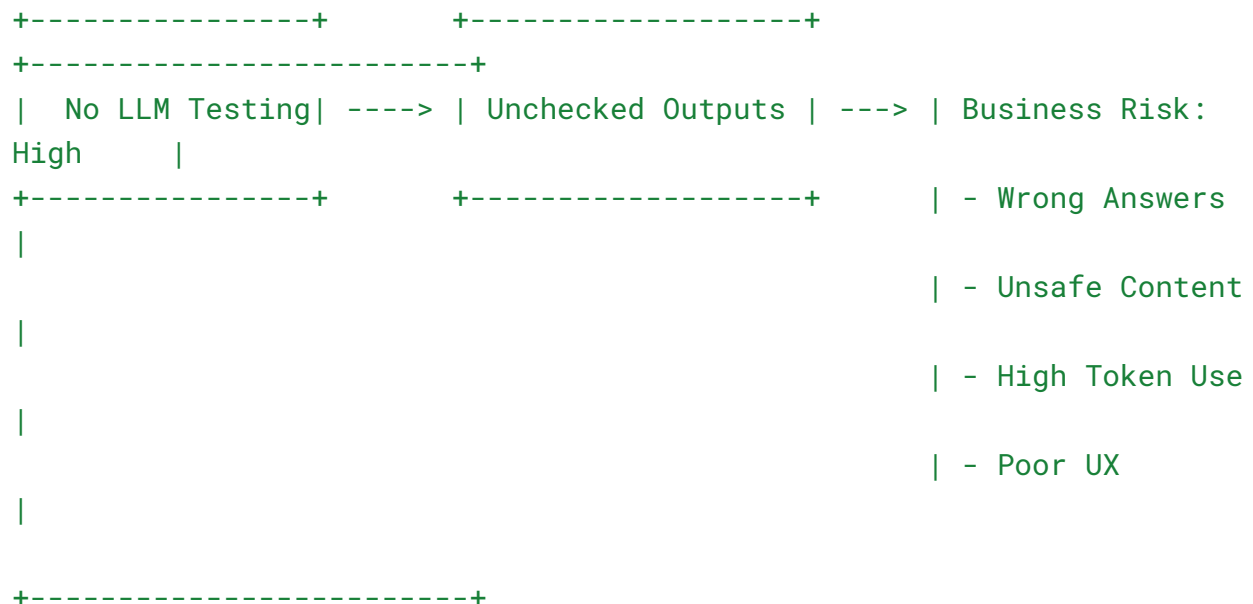
Would you like a visual flow chart or KPI dashboard view showing “Before Testing vs After Testing” with these token and quality improvements?

Slide Title: “How LLM Testing Reduces Risk, Saves Cost & Improves Quality”

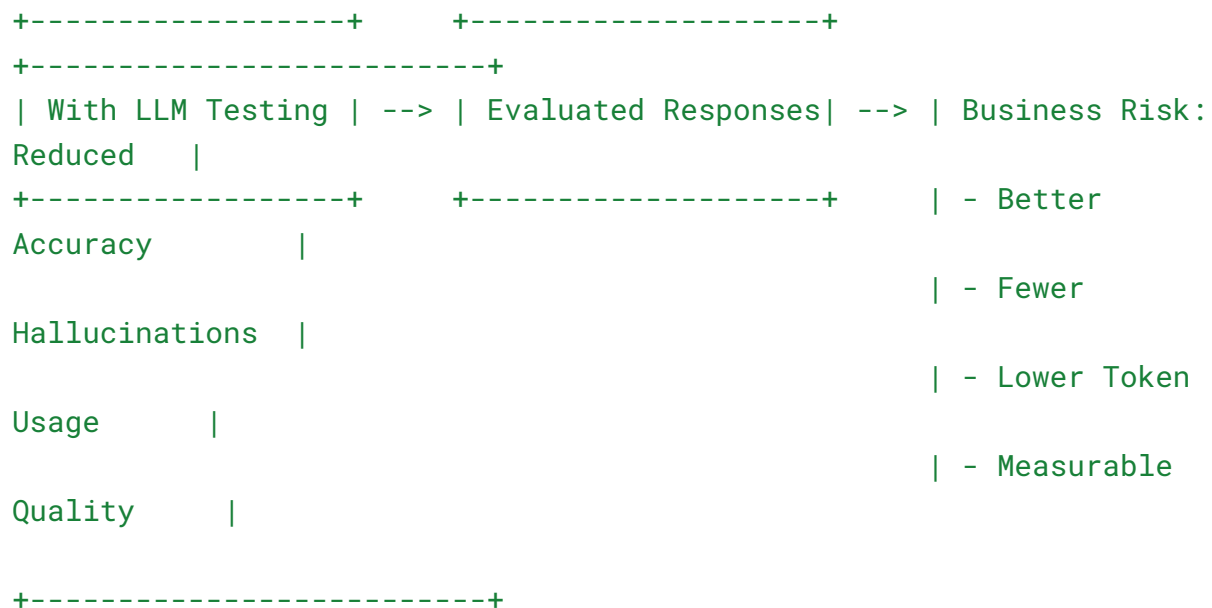
Before vs After LLM Testing – Visual Flow

pgsql

CopyEdit



vs.



KPI Table: Before vs After Testing

Metric	Without LLM Testing	With LLM Testing (RAGAS / DeepEval)
Avg Token per Query	100	70
Wrong Answer Rate	15%	3%
User Follow-up Rate	High	Low
Cost per 10K Queries	\$25–30	\$15–20
Trust/Safety Violations	Undetected	Caught Before Release
Model Confidence	Unknown	Measured with Rubrics & Scores

Real Impacts

- 💰 **Token Savings:** Optimized prompts and shorter outputs = 30–40% token cost saved.
 - 📉 **Risk Reduction:** Hallucinations & unsafe responses detected pre-deployment.
 - 📈 **CX Improvement:** Users get more accurate, reliable, and helpful responses.
 - ⚡ **Faster Iterations:** Evaluate new versions quickly before rollout.
-

Summary Line (for CIO/CTO level)

“LLM Testing acts like a quality and risk firewall — ensuring your AI delivers safe, smart, and cost-efficient responses before they ever reach your users.”

Would you like this as a **PowerPoint (.PPTX)** or **PDF one-pager** version? I can generate a downloadable file for you.

You said:

What is Rubrics & Scores

ChatGPT said:

Great question. Let's break down **Rubrics & Scores** in a **very simple and layman-friendly way** — especially useful when explaining LLM evaluation to your tech team, QA, or even leadership.



What is a Rubric in LLM Testing?

A **rubric** is just a **set of rules or criteria** we use to **judge** how good the AI's response is.

Think of it like a **school grading sheet**.

🧠 Example: If a student writes an essay, the teacher might grade it using a rubric with:

- Grammar: 1 point
- Relevance: 1 point
- Factual correctness: 1 point
- Clarity: 1 point
→ Total: 4 points

We do the **same with LLM outputs**.



Rubric in LLM Testing – Example Criteria

When we test an LLM response, we check things like:

Rubric (Scoring Criteria)	What it Checks	Example
Faithfulness	Is the answer grounded in context/data?	Does it use facts from documents?
Relevance	Is the answer relevant to the question asked?	Does it stay on-topic?
Completeness	Does it answer the full question?	No partial answers

Conciseness	Is it clear and short, without unnecessary details?	No rambling
Toxicity/Safety	Is the tone safe and non-offensive?	Avoids bias, hate, etc.
Instruction Following	Did the model follow format/style as instructed?	E.g., "Answer in bullet points"



What are Scores?

Each rubric item gives a **score**, usually from **0 to 1**, or **0 to 5**.

- A score of 1 = perfect
- A score of 0 = failed

Then we **average or combine** the scores to get an overall **quality score** for that answer.



Example:

Input: "What is the capital of France?"

Model Answer: "Paris is the capital of France."

Rubric Scores:

Criteria	Score
Faithfulness	1
Relevance	1
Completeness	1
Conciseness	1
Instruction Followed	1
Toxicity	1

Final Score: **6 / 6 = 100%**

Now imagine the model said "Delhi is the capital":

- Faithfulness = 0
- Relevance = 0.5
- Others = maybe 1

👉 Final Score: $3.5 / 6 = 58\%$

Who Does the Scoring?

- In tools like **DeepEval** or **RAGAS**, these rubrics are:
 - **Predefined** or **customizable**
 - Scored using another **LLM judge** model like GPT-4
 - Or can be **scored by humans** for benchmarking
-

Summary in One Line

A **rubric** is like a grading sheet for AI responses. It defines **what good looks like**, and **scores** tell us how close the model came to that.

Would you like an example of how to define rubrics using DeepEval or RAGAS in code?

The **disadvantages of non-LLM-based metrics** compared to **LLM-based metrics** primarily stem from their limited ability to understand **semantic meaning**, **context**, and **natural language variation**. Here's a breakdown in simple terms:

What Are Non-LLM Metrics?

These are traditional or statistical methods like:

- BLEU, ROUGE, F1, Cosine Similarity, Jaccard, etc.
 - They compare **words**, **tokens**, or **vectors**, not meaning.
 - No reasoning, no understanding—just surface-level comparisons.
-

✗ Disadvantages of Non-LLM Metrics

1. Poor Semantic Understanding

- They can't detect if two sentences mean the same thing with different words.
- Example:
 - Expected: *"The capital of France is Paris."*
 - Output: *"Paris is the capital of France."*
 - BLEU or ROUGE may give **low score** even though it's correct.
 - LLM-based metrics would give a **high score** due to semantic match.

2. No Context Awareness

- They don't consider the **intent** or **overall conversation**.
- Example: If a chatbot says "Yes, you can do that" without referring to what "that" is, the metric won't know if it's a good answer.

3. Token Overlap Bias

- These metrics reward **word matching**, not meaning.
- So a completely wrong but word-similar response can get a high score.
- LLM-based metrics focus on **truthfulness and relevance**.

4. Inability to Capture Factual Accuracy

- If an answer is factually wrong but similar in structure, non-LLM metrics may not penalize it.
- LLM-based metrics can be tuned to check for factual consistency.

5. No Reasoning or Logic Evaluation

- You can't use non-LLM metrics to evaluate **chain-of-thought**, reasoning, or multi-step answers.
- LLM-based metrics can simulate human judgment and reasoning quality.

6. Insensitive to Language Fluency or Grammar

- Non-LLM metrics won't flag if the output is unnatural or ungrammatical.
- LLM-based evaluations can assess **fluency** and **human-likeness**.

Summary Table

Feature	Non-LLM Metrics	LLM-Based Metrics
Understands meaning	✗ No	✓ Yes
Handles paraphrasing	✗ Poorly	✓ Well
Considers context	✗ No	✓ Yes
Checks factual correctness	✗ No	✓ Yes (via context)
Evaluates fluency	✗ No	✓ Yes
Fast and cost-effective	✓ Yes	✗ Cost (LLM API/token use)

When to Avoid Non-LLM Metrics

- Evaluating **chatbots**, **RAG systems**, or **agent outputs**

- Use cases where **meaning, truth, or reasoning** matter more than wording

Let me know if you want examples or visual comparison using BLEU vs GPTScore or RAGAS.

Here's a **detailed comparison** between **RAGAS** and **DeepEval**, two popular frameworks for **evaluating LLM-based and RAG applications**, especially in the context of QA and testing.

High-Level Overview

Feature	RAGAS	DeepEval
Full Name	Retrieval-Augmented Generation Assessment Suite	DeepEval (by Confident AI)
Purpose	Evaluate RAG pipelines using built-in quality metrics	Evaluate LLM outputs with test-like assertions and metrics
Focus	RAG-specific evaluation (faithfulness, context relevance, etc.)	General-purpose LLM evaluation (chatbots, classification, summarization)
Best For	RAG pipelines with context+question+answer flow	Unit/integration testing for LLM apps, agent workflows, chatbots
Language	Python	Python

Evaluation Metrics

♦ RAGAS Metrics

RAGAS focuses on **retrieval quality** and **answer correctness** in RAG pipelines:

Metric	Description
Faithfulness	Whether the answer is grounded in retrieved context
Answer Relevance	Relevance of the answer to the question
Context Precision	How much of the retrieved context is useful

Context Recall	Whether all necessary info was retrieved
Context Relevance	Overall usefulness of retrieved documents
Answer Semantic Similarity	Optional — semantic match between generated and reference answers

◆ DeepEval Metrics

DeepEval supports a **wider set of task-agnostic and custom metrics**:

Metric	Description
Faithfulness	Same as RAGAS, checks hallucinations
Relevance	Is the output relevant to the input query
Coherence	Fluency and consistency of answer
Answer Correctness	Is the model's output factually or logically correct
Toxicity	Measures if output is offensive or harmful
Custom Metrics	You can define your own via function decorators
TestCase-based scoring	Supports unit-like tests with pass/fail and threshold scoring

Features Comparison

Feature	RAGAS	DeepEval
RAG support	✓ Native	✓ (can simulate RAG pipeline)
Reference answer support	✓ Required for some metrics	✓ Optional
Custom metrics	✗ (Not yet)	✓ Fully supported
Custom test cases	✗	✓ Decorator-based LLM test cases
Streaming support	✗	✓

Logging and dashboards	✗	✓ Confident AI dashboard (optional)
Evaluation granularity	Batch-level	Test-case level
Integration with LangChain	✓	✓
Agent evaluation	✗	✓

✓ When to Use What?

✓ Use RAGAS when:

- You're building a **RAG pipeline** (LangChain, Haystack, LlamaIndex).
- You want **structured metrics** to validate the retrieval+generation flow.
- You have access to **reference answers** for comparison.
- You want to **evaluate on a dataset** of Q&A pairs with contexts.

✓ Use DeepEval when:

- You're testing **LLM outputs for correctness, toxicity, faithfulness**.
 - You want to do **unit or integration tests** for your LLM pipeline.
 - You're evaluating **chatbots, agents, summarizers, classifiers**, etc.
 - You want to use or define **custom test cases and assertions**.
 - You want to plug into a **CI/CD pipeline** for regression testing.
-

🔧 Example Use Case Differences

Scenario

Recommended
Tool

Evaluating how well your RAG answers match your source documents	RAGAS
Creating automated tests for chatbot responses with custom logic	DeepEval
Measuring retrieval quality and faithfulness in QA format	RAGAS
Running CI tests to ensure model doesn't hallucinate or return toxic output	DeepEval

Tech Stack Integration

Feature	RAGAS	DeepEval
LangChain	✅ Yes	✅ Yes
HuggingFace	✅ Yes	✅ Yes
CI/CD Integration	❌ (manual only)	✅ Yes (assertions/testcases)
Gradio/Streamlit integration	Manual setup	✅ via examples
Prompt injection detection	❌	✅ Supported

Installation

```
bash
CopyEdit
# RAGAS
pip install ragas

# DeepEval
pip install deepeval
```

Summary Table

Category	RAGAS	DeepEval
----------	-------	----------

Evaluation Scope	RAG pipelines only	Any LLM task
Reference Answer	Needed for some metrics	Optional
Custom Assertions	✗	✓ Yes
CI/CD Ready	✗	✓ Yes
Dashboard Integration	✗	✓ (Confident AI optional)
Customizability	Medium	High
Granularity	Dataset-level	Unit/test-level
Best for	RAG pipeline performance validation	Testing LLM outputs in pipelines/apps

If you tell me your **use case**, I can help you choose or even write code to implement the right tool for evaluation.

You said:

Yes

ChatGPT said:

Great! Could you please share a bit more about your **use case**?

For example:

- Are you testing a **RAG application** (like LangChain, Haystack, LlamaIndex)?
- Is it a **chatbot**, **agent**, **search assistant**, or **document QA bot**?
- Do you have **ground truth answers** (reference answers)?
- Do you need **unit/integration tests**, or just a **quality evaluation** report?

Once I know this, I'll guide you on whether to use **RAGAS**, **DeepEval**, or a combination of both — and I can provide you a full working code setup.

Let me know 🙋

You said:

Use all from example

ChatGPT said:

Updated saved memory

Perfect! Since your use case includes:

- **Chatbot + Agent + Search Assistant + Document QA**
- **Ground Truth Answers available**
- Need for both **unit/integration tests** and **quality evaluation**

👉 You should use **both RAGAS and DeepEval together** for full coverage.

✅ Best Strategy for Your Use Case

Tool	Role in Your Pipeline
RAGAS	Evaluate retrieval and generation quality — focus on faithfulness, context relevance, etc.
DeepEval	Build unit/integration test cases , validate correctness, toxicity, and support CI/CD

🔧 Sample Setup (Step-by-Step)

1. ✅ Install Dependencies

bash

CopyEdit

```
pip install ragas deepeval datasets
```

2. 📁 Prepare Your Dataset (with **question, answer, contexts**)

Example format (can be from CSV/JSON):

python

CopyEdit

```
from datasets import Dataset
```

```
data = [
    {
        "question": "What is LangChain?",
        "answer": "LangChain is a framework for building applications
with LLMs.",
        "contexts": ["LangChain is an open-source framework for
developing applications powered by language models."]
    },
    {
        "question": "What is DeepEval?",
        "answer": "DeepEval is a testing framework for evaluating LLM
outputs.",
        "contexts": ["DeepEval is used to build unit tests for LLM
applications and check for hallucinations."]
    }
]

dataset = Dataset.from_list(data)
```

3. RAGAS Evaluation

python

CopyEdit

```
from ragas.metrics import faithfulness, answer_relevancy,
context_precision, context_recall
from ragas import evaluate

result = evaluate(
    dataset=dataset,
    metrics=[faithfulness, answer_relevancy, context_precision,
context_recall]
)

print(result)
```

4. DeepEval Integration with Custom Tests

python

CopyEdit

```
from deepeval.test_case import LLMTestCase
from deepeval.metrics import AnswerRelevancyMetric, FaithfulnessMetric
from deepeval.evaluate import evaluate
```

```
test_cases = [
    LLMTestCase(
        input="What is LangChain?",
        actual_output="LangChain is a framework for building
applications with LLMs.",
        expected_output="LangChain is a framework for building
applications with LLMs."
    ),
    LLMTestCase(
        input="What is DeepEval?",
        actual_output="DeepEval is used for testing LLM outputs.",
        expected_output="DeepEval is a testing framework for
evaluating LLM outputs."
    )
]

evaluate(
    test_cases=test_cases,
    metrics=[AnswerRelevancyMetric(threshold=0.8),
FaithfulnessMetric(threshold=0.8)]
)
```

5. Optional: Add Custom Toxicity, Hallucination Checks

python

CopyEdit

```
from deepeval.metrics import ToxicityMetric
```

```
evaluate(
    test_cases=test_cases,
    metrics=[ToxicityMetric(threshold=0.3)]
)
```

)



Result

- **RAGAS** gives you structured metrics per question (**faithfulness: 0.89**, **context_precision: 0.92**, etc.)
- **DeepEval** tells you if the response **passed/failed** based on thresholds (good for CI/CD)
- You can run DeepEval tests automatically in CI with GitHub Actions

Would you like me to help you set this up on your **real dataset**, integrate with **LangChain**, or generate a **report/dashboard**?

Let me know and I can tailor the next steps!




The Journey of LLM Testing and Evaluation (Layman's Guide)



1. What is the Traditional Approach to Software Testing?

Traditional testing involves:

- Writing test cases based on input and checking if the output is correct.
- For example:
 - **Input:** 2 + 2
 - **Expected Output:** 4
 - If actual output is 4 →  Passed

Used in:

- Web apps, APIs, backend logic, forms, databases, etc.

⚠️ 2. Issues with Traditional Approach for LLMs (Large Language Models)

LLMs (like ChatGPT, Gemini, Claude, etc.) are **not deterministic**. That means:

- They don't give exactly the same answer every time.
- There can be multiple correct answers.

Example:

- Input: "Summarize the Ramayana."
- One time it may say "It's the story of Lord Rama's journey and victory over Ravana."
- Another time: "An Indian epic focusing on virtue, duty, and family."

👉 Traditional pass/fail won't work here because **both are valid answers**.

vs 3. LLM Testing vs LLM Evaluation

Aspect	LLM Testing	LLM Evaluation
Goal	Check if LLM performs a task correctly	Measure the quality and usefulness of responses
Type	Functional testing (like automation)	Scoring & grading the answer quality
Example	"Can it generate a ticket when asked?"	"How relevant was the answer to the question?"

🧠 Think of **Testing** like checking if the brain works, and **Evaluation** as checking **how smart or useful** the brain is.

📊 4. LLM-Based Metrics vs Non-LLM-Based Metrics

✅ Non-LLM-Based Metrics (Traditional / Rule-Based)

- **BLEU**: Checks similarity based on word overlap.
- **ROUGE**: Compares overlap of words or phrases.
- **Cosine Similarity**: Measures how close two sentence vectors are.

👉 Easy to compute, but:

- Can fail to understand meaning
- Penalizes correct answers with different wording

🤖 LLM-Based Metrics (Smarter, Meaning-Based)

- Uses another LLM to **grade** answers.
- Understands **semantics** and **intent**.
- Can give scores like “this is 80% relevant”.

Example:

- Human says: “What is the capital of France?”
- LLM answers: “Paris is the capital city of France.”
- Even if not exact match to expected: score could be 1.0

🔍 5. How to Test RAG vs Chatbot

📘 RAG (Retrieval-Augmented Generation)

- Combines documents + LLM.
- Fetches **relevant context** before answering.
- Testing:

- Is it retrieving the right document?
- Is it using the context correctly?

Chatbot

- Directly responds to user.
- May or may not have retrieval.
- Testing:
 - Intent identification: What does the user want?
 - Follow-ups and field filling (e.g. ticket form)

6. RAGAS vs DeepEval (Evaluation Frameworks)



Feature	RAGAS	DeepEval
Focus	RAG system evaluation	General LLM output evaluation
Metrics	Faithfulness, Context Recall, Answer Relevance	GEval, Embedding similarity, Custom scoring
Input	RAG pipeline output	Any chatbot / LLM response
Uses LLM?	✅ Yes	✅ Yes
Example Use	Did the RAG system fetch and use correct context?	Was the answer smart, helpful, and correct?

Example Test + Evaluation (Putting it All Together)

Scenario:

User says: “Create a ticket for my WiFi not working.”

✅ Testing (Functional):

- Does it detect **intent**: ticket creation? → 
- Does it ask for **follow-up**: "Is it for self or someone else?" → 

Evaluation (Scoring):

- Is the reply relevant? (Answer Relevance)
 - Did it use the right context? (Faithfulness)
 - Does it sound like a human? (Fluency)
-

Conclusion: Why Do We Need Both?

- **Testing** makes sure the LLM **works**.
- **Evaluation** tells us how **well** it works.

Together, they help you:

- Identify functional bugs
 - Measure quality
 - Improve user experience
-

Summary in 1 Line:

Traditional testing checks if something works. LLM testing checks if it works AND makes sense.