# RAG Evaluation: The Good, The Bad, and The Tricky

leboncoin
tech

**Reka HALMAI**

**Machine Learning
Engineer
@ leboncoin**

**leboncoin**

**Anis ZAKARI**

**Machine Learning
Engineer
@Hymaïa x @leboncoin**

**hymaïa**

# Index

# 01
## Introduction

# Large Language Models (LLMs)
## What are they ?

**Key Characteristics:**

- Massive deep learning models using transformer architecture
- Billions of parameters
- Require multiple GPUs for high computational power
- ChatGPT popularized LLMs

**Training:**

- Pre-trained on vast amounts of text data
- Training corpus often unknown (IP concerns)
- Learn language patterns, facts, and reasoning

**Functionality:**

- Predict next word/token based on previous context
- Generate text sequentially until stop condition

# How can we ensure they work?

## By evaluating!

Evaluation is about **trust** to ensure the **accuracy** and **relevance** of the generated outputs and avoid **regression.**

**Validates** the model's performance in high-stakes applications (e.g., healthcare, customer support).
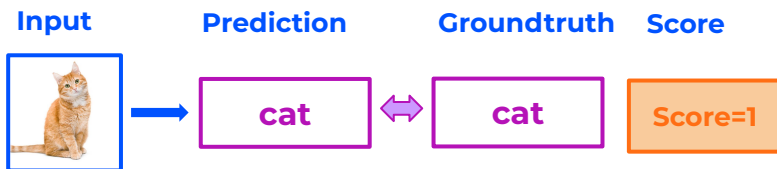
What we need is **good metrics. But what is a good metric?**

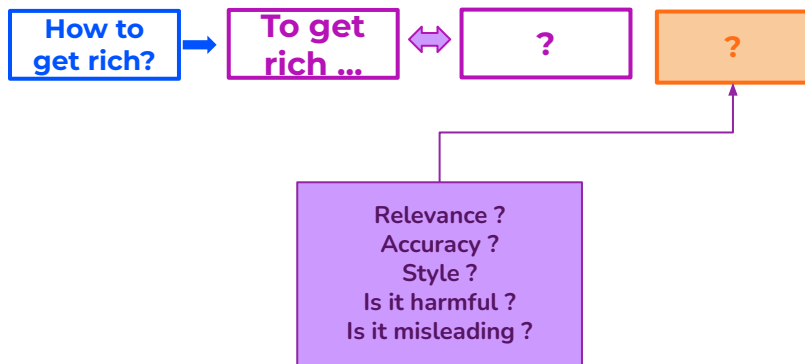# 02

# Why LLM evaluation is tricky?

# Evaluate something that is generated
## How would you do that ?

**Classification model**

| Input | Prediction | Groundtruth | Score |
|-------|-----------|-------------|-------|
| 🐱 | cat | ⟺ cat | Score=1 |

**Generative model**

| How to get rich? | To get rich ... | ⟺ ? | ? |

Relevance ?
Accuracy ?
Style ?
Is it harmful ?
Is it misleading ?

**Other Difficulties with LLM to consider:**

- One input, many possible answers
- Subjectivity
- Sensitivity: slight changes might affect the output
- Non-determinism
- Continuous Evaluation

# Evaluation
## Types of evaluation

### Text similarity metrics
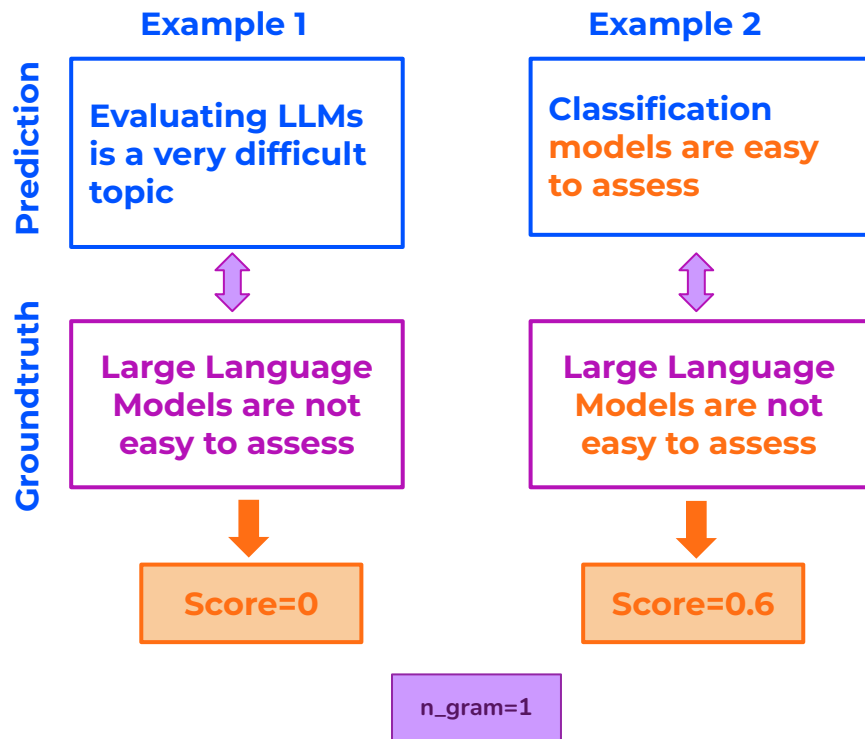**(ex ROUGE, BertScore...)**



### Human Evaluation



### Model evaluation/ LLM as a judge

# Evaluation with text similarity scores
## ROUGE: Naive Text Similarity Metric with Inherent Limitations

**Example 1**

**Prediction**

Evaluating LLMs is a very difficult topic

**Groundtruth**

Large Language Models are not easy to assess

Score=0

**Example 2**

Classification models are easy to assess

Large Language Models are not easy to assess

Score=0.6

n_gram=1

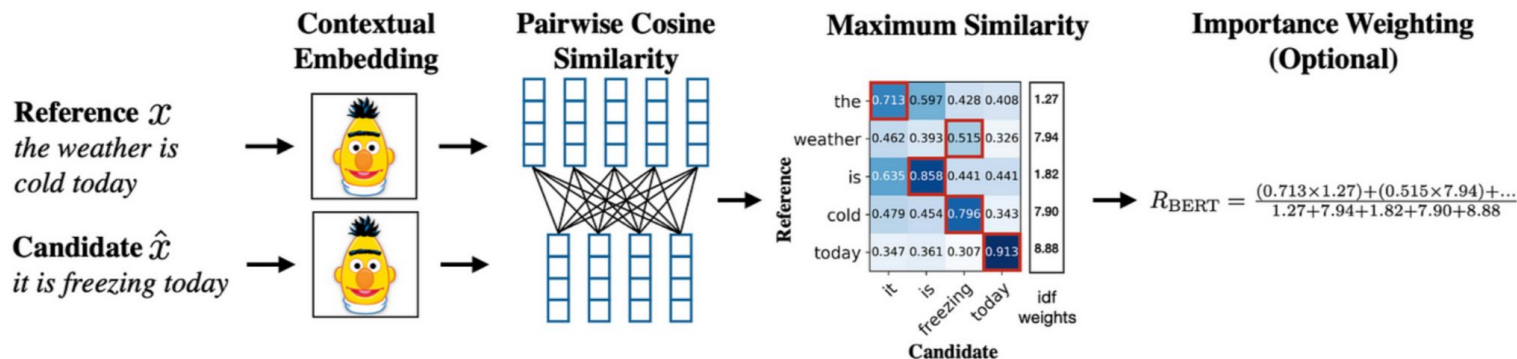**Text similarity scores are:**

- **Limited** to **surface-level text similarity**
    - ignores context and semantics
    - low correlation with human judgments
- Unable to measure **factual accuracy** or **user satisfaction**
- **Biased** toward rigid matching of reference texts

\* Source: Blagec, Dorffner, Moradi: A global analysis of metrics used for me
\* ROUGE: Recall-Oriented Understudy for Gisting Evaluation

# Evaluation with semantic similarity scores
## BERTScore - Capturing Meaning Beyond Surface Matching



**Contextual Embedding**

**Reference** $x$
*the weather is cold today*

**Candidate** $\hat{x}$
*it is freezing today*

**Pairwise Cosine Similarity**

**Maximum Similarity**

|  | it | is | freezing | today | idf weights |
|---|---|---|---|---|---|
| the | 0.713 | 0.597 | 0.428 | 0.408 | 1.27 |
| weather | 0.462 | 0.393 | 0.515 | 0.326 | 7.94 |
| is | 0.635 | 0.858 | 0.441 | 0.441 | 1.82 |
| cold | 0.479 | 0.454 | 0.796 | 0.343 | 7.90 |
| today | 0.347 | 0.361 | 0.307 | 0.913 | 8.88 |

Reference / Candidate

**Importance Weighting (Optional)**

$$R_{\text{BERT}} = \frac{(0.713 \times 1.27) + (0.515 \times 7.94) + \dots}{1.27 + 7.94 + 1.82 + 7.90 + 8.88}$$

**Advantages:**

- Captures semantics beyond surface matching
- Handles paraphrases and linguistic variations
- Correlates better with human judgments

**Limitations:**

- Biased toward reference text
- Semantic similarity =/= relevance
- Computationally expensive

# Human Evaluation Challenges
## The Human Factor

### Advantages

- Captures **nuanced aspects**
- Identifies **safety risks** and biases
- **Qualitative insights** for improvement

### Disadvantages

- **Cost and Time - difficult to scale**
- **Need for Expertise:** Requires domain experts/ specific knowledge
- **Subjectivity and Bias:**
  - Evaluators' differing opinions and biases can skew results.
  - Inconsistent judgments reduce reliability and reproducibility

# LLM as a judge
## Challenges of AI-Based Evaluation

### Advantages

· **Cost-effective** alternative to human evaluation

· **Scalable** for large-scale assessments

· **Consistent scoring (more or less)** across multiple samples

· Empirically **good results**

### Disadvantages

· Potential **biases**

- ➢ Inherited from training data
- ➢ Prefer a single score
- ➢ Prefer own model prediction
- ➢ Order counts

· Requires **careful prompt engineering** and fine-tuning

· **Lack of determinism** in iterations

**Best used in conjunction with human oversight**

# 03

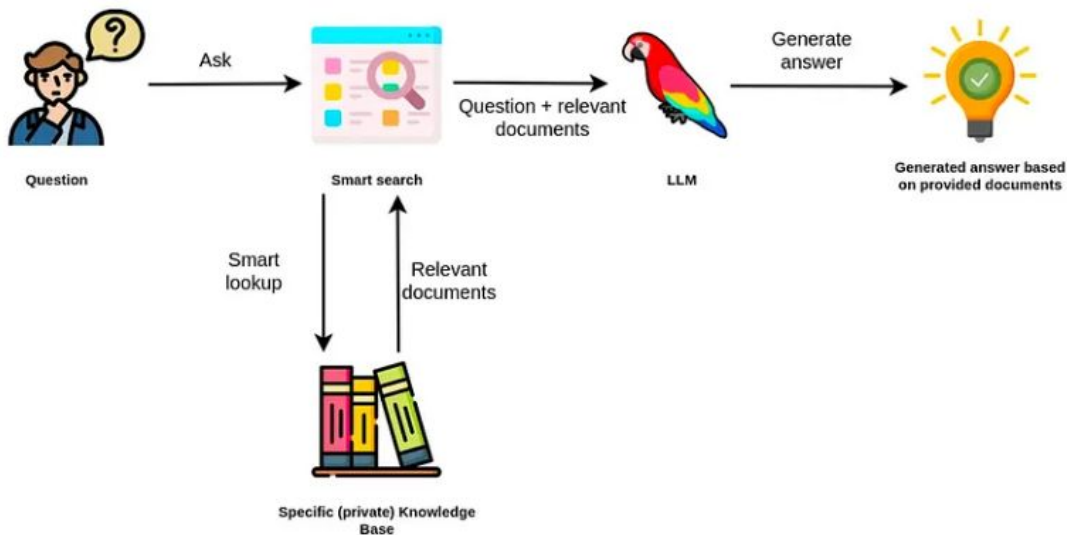# Why RAG evaluation is even trickier?

# What is RAG?
## Understanding Retrieval Augmented Generation (RAG)

RAG combines **retrieval** (locating relevant data from a knowledge base) with text **generation** (leveraging the retrieved information as context for generation).

**Key benefits**:

- Adapts the LLM to **rare or unseen data**

- **Enhances factual accuracy and reduces hallucinations.**

# RAG Evaluation
## The metrics should answer the questions...



**Did we find the right information?**

- Reference context via LLM as a judge
- Information retrieval / Ranking metrics (Precision, recall, F1 score, MRR, nDCG)

**Given the retrieved information, did we generate the correct answer?**

- Rule based metrics
- Groundedness via LLM as a judge

**Is the prediction the correct answer?**

- Standard scores (ROUGE, METEOR)
- Semantic based scores (BertScore)
- Human evaluation
- Answer relevance via LLM as a judge

# RAG-Specific Challenges 1.
## Duality

Combination of **retrieval** and **generation**, each requiring different evaluation metrics.
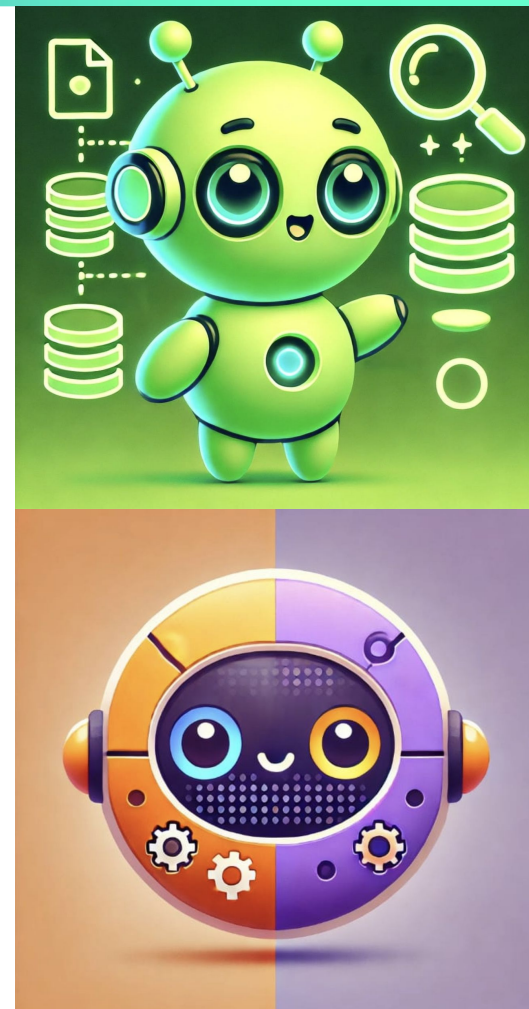
⇒ **Challenges related to the retriever**

- Irrelevant information "creates noise"
- Lack of relevant information makes it impossible to reply well

⇒ **Challenges related to the generation**

- Preserving relevant information
- Mitigating hallucinations/incorrect generations
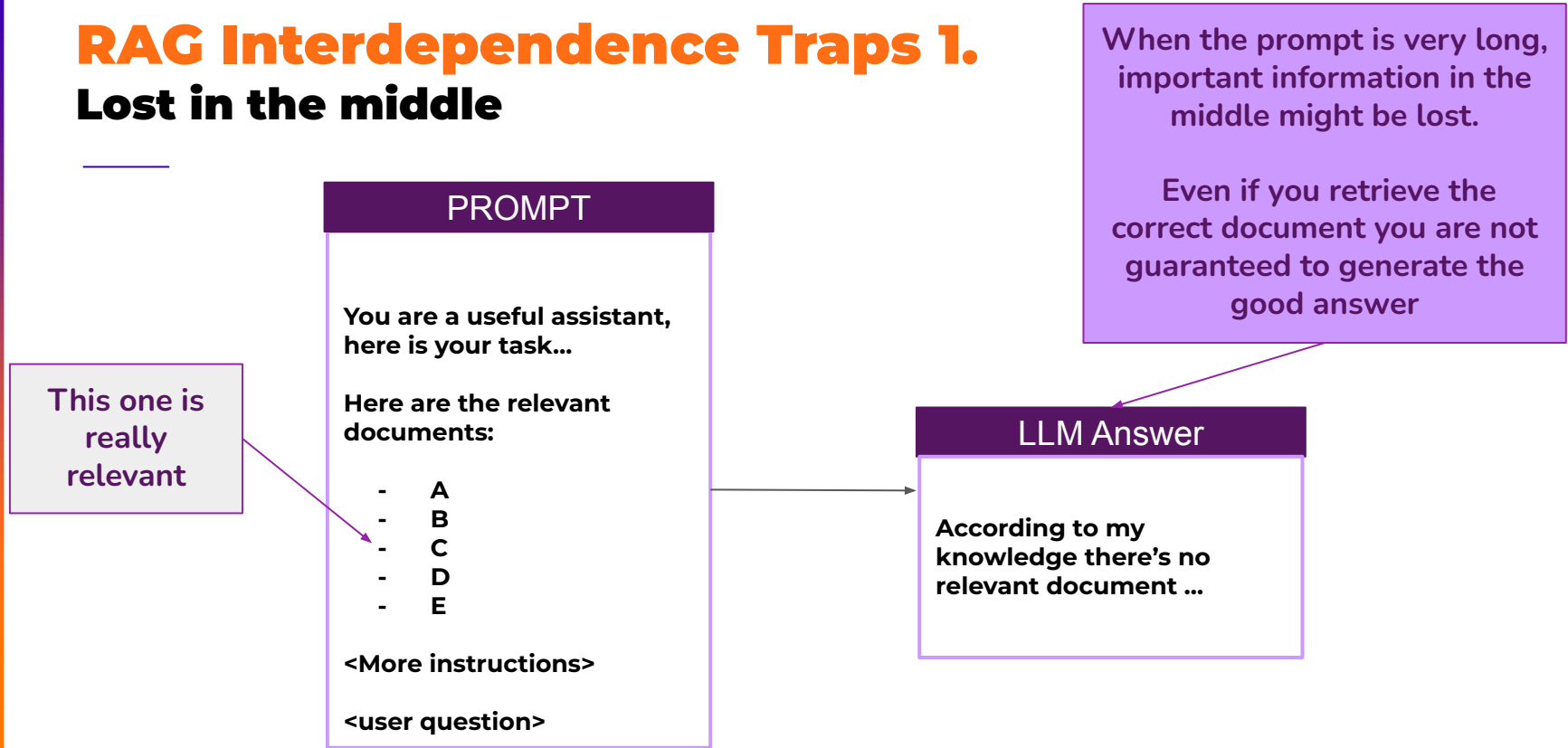- Handling query diversity effectively

Interdependence between them makes it hard to isolate their impact on performance
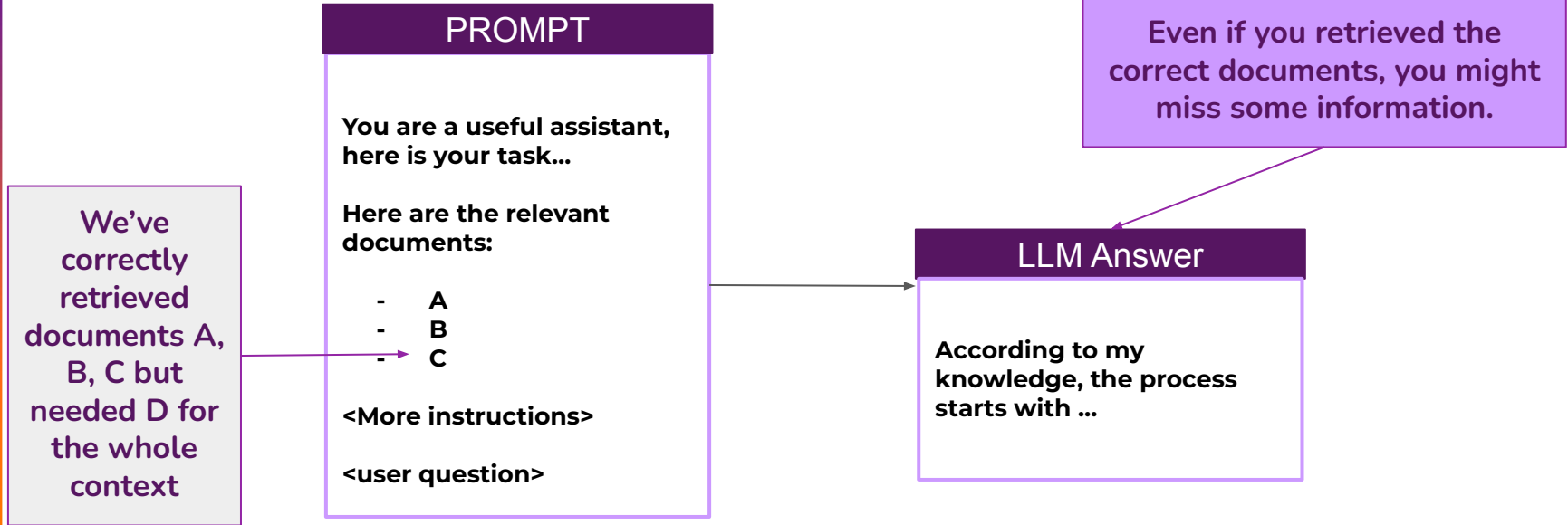
⇒ **Interdependence Traps**

# RAG Interdependence Traps 1.
## Lost in the middle

**PROMPT**

You are a useful assistant, here is your task...

Here are the relevant documents:

- A
- B
- C
- D
- E

\<More instructions\>

\<user question\>

**This one is really relevant**

When the prompt is very long, important information in the middle might be lost.

Even if you retrieve the correct document you are not guaranteed to generate the good answer

**LLM Answer**

According to my knowledge there's no relevant document ...

# RAG Interdependence Traps 2.
## Context Truncation

**We've correctly retrieved documents A, B, C but needed D for the whole context**

### PROMPT

**You are a useful assistant, here is your task...**

**Here are the relevant documents:**

- **A**
- **B**
- **C**

**<More instructions>**

**<user question>**

The generation component fails to add the proper project ending.

Even if you retrieved the correct documents, you might miss some information.

### LLM Answer

**According to my knowledge, the process starts with ...**

# RAG Interdependence Traps 3.
## Hallucination Due to Weak Retrieval

The generation component hallucinates because of the irrelevant context retrieved.

Even if you retrieve the correct document, you might need to control for not retrieving the incorrect ones.

This doc is relevant

These docs are not relevant

**PROMPT**

**You are a useful assistant, here is your task...**

**Here are the relevant documents:**

- **A**
- **B**
- **C**

**<More instructions>**

**<user question>**

**LLM Answer**

**According to my knowledge, ...**

# RAG-Specific Challenges 2.
## Complexity related to the knowledge base

**Knowledge base quality:**

- Incompleteness, biases, and inaccuracies
- Outdated or stale information
- Duplicates and redundant information
- Inconsistencies and contradictions within the knowledge base

**Dynamic External Knowledge:**

- Frequently updated knowledge bases
- Difficulty in consistent evaluation
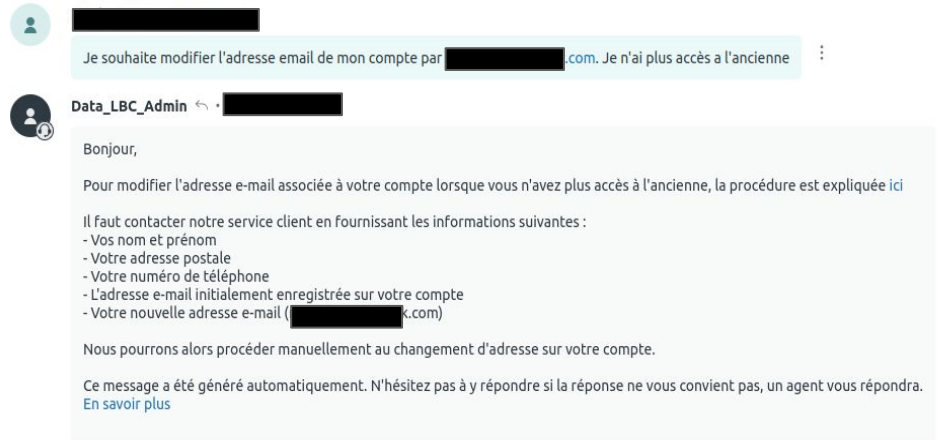
**Multilingual bias in LLM/Embedding models**

# 04
## Our use cases

# Use case 1: Markus
## leboncoin's Customer care Assistant

**What:**

- Reply to users tickets automatically

**Why? :**

- ease the burden on customer service.
- ~200k tickets /month

🦜🔗 **LangChain**

♾ **FAISS**

Amazon Bedrock



Je souhaite modifier l'adresse email de mon compte par ████████.com. Je n'ai plus accès à l'ancienne

**Data_LBC_Admin** ↩ · ████████

Bonjour,

Pour modifier l'adresse e-mail associée à votre compte lorsque vous n'avez plus accès à l'ancienne, la procédure est expliquée ici

Il faut contacter notre service client en fournissant les informations suivantes :
- Vos nom et prénom
- Votre adresse postale
- Votre numéro de téléphone
- L'adresse e-mail initialement enregistrée sur votre compte
- Votre nouvelle adresse e-mail (████████k.com)

Nous pourrons alors procéder manuellement au changement d'adresse sur votre compte.

Ce message a été généré automatiquement. N'hésitez pas à y répondre si la réponse ne vous convient pas, un agent vous répondra.
En savoir plus

# Evaluating Markus
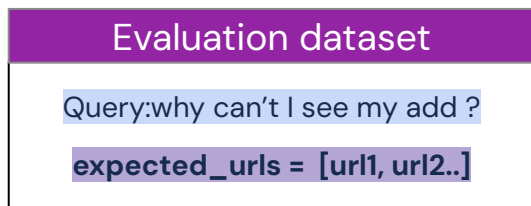## Different evaluation steps

**Pre-production evaluation:**

- Retriever evaluation
- LLM's ability to choose the right documents
- LLM as a judge to assess generated answers on specific scenarios
- Customer care support to assess if answers are acceptable

**Post-production evaluation:**
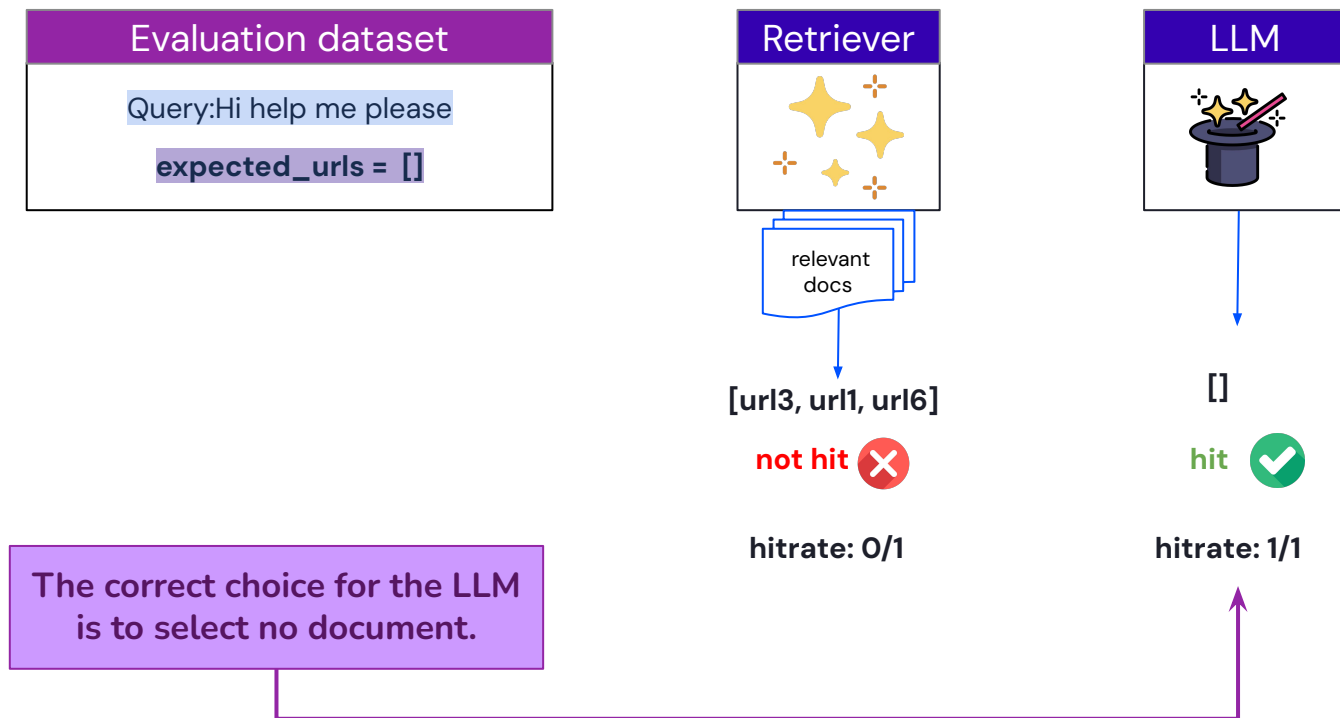
- Guardrails
- Feedback from customer care

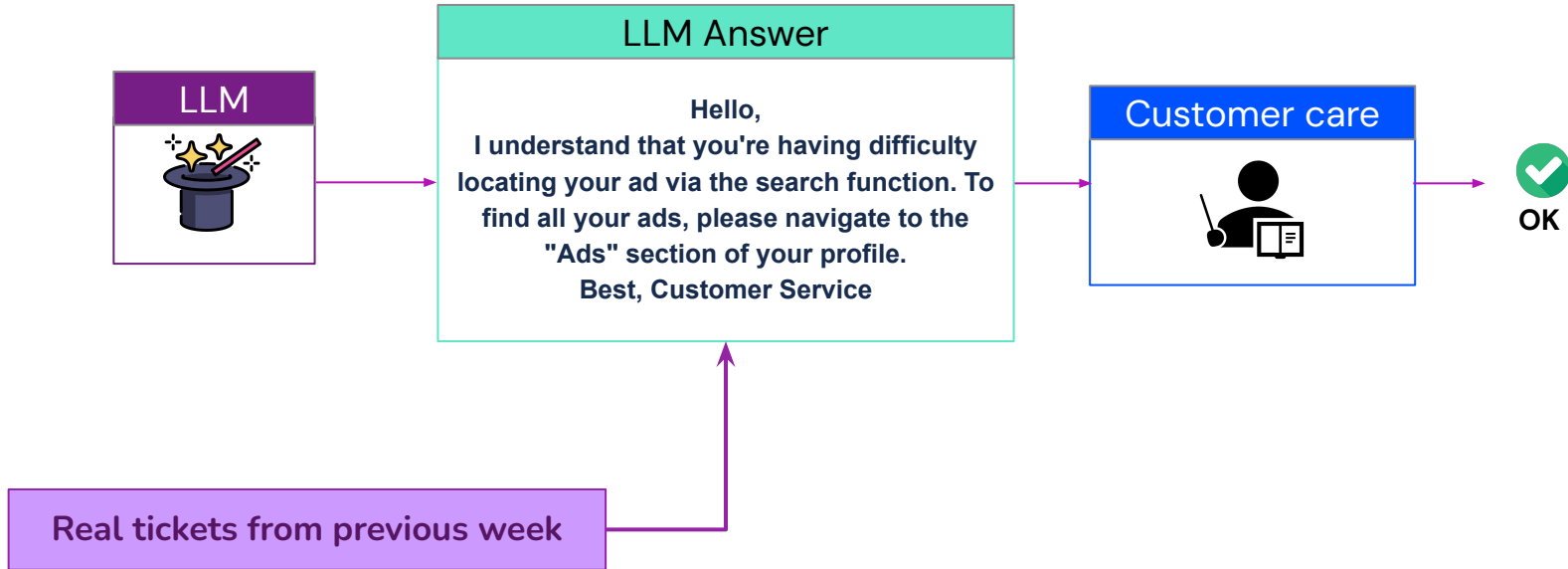# Small dataset on specific scenarios to be covered
# Positive examples



**Evaluation dataset**

Query:why can't I see my add ?

expected_urls = [url1, url2..]

**Retriever**

relevant docs

**LLM**

[url3, **url1**, url6]

url6

**hit** ✅

**not hit** ❌

**hitrate: 1/1**

**hitrate: 0/1**

# Small dataset on specific scenarios to be covered
# Negative examples

**Evaluation dataset**

Query:Hi help me please

expected_urls = []

**Retriever**

relevant docs

[url3, url1, url6]

**not hit** ❌

hitrate: 0/1

**LLM**

[]

**hit** ✅

hitrate: 1/1

The correct choice for the LLM is to select no document.

# Customer care support
# Validate behaviour and accuracy



**LLM**

**LLM Answer**

Hello,
I understand that you're having difficulty locating your ad via the search function. To find all your ads, please navigate to the "Ads" section of your profile.
Best, Customer Service

**Customer care**

OK

**Real tickets from previous week**

# LLM as a judge
# Validate behaviour and accuracy

**User**

Query:why can't I see my add ?

**LLM**

**LLM Answer**

Hello,
I understand that you're having difficulty locating your ad via the search function. To find all your ads, please navigate to the "Ads" section of your profile.
Best, Customer Service

**LLM as a judge**

**4/5**

**User**

Query:why can't I see my add ?

**Example**

Possible answer

# Feedback from Customer Service
# Human in the loop

**LLM**

**LLM Answer**

I understand that you want to reach us by telephone, here is our direct number: 01 xx xx xx xx

**Customer Care**

**Markus Team**

Feedback

# Post Production
# Guardrails

**LLM**

**LLM Answer**

I understand that you want to reach us by telephone, here is our direct number: 06 xx xx xx xx

**Guardrails**

STOP

# Use case 2: Ada - AI Assistant
## leboncoin's internal chatbot

leboncoin's internal chatbot developed as a result of privacy concerns.

**Key Services:**

- General chatbot (UI and Slack)
- Slack thread and channel summarization
- Document upload and processing

**RAG Integration for Internal Knowledge:**

- Company policies (expenses, procurement, etc...)
- Projects, onboarding, and team information (Confluence)

# Evaluating Ada RAG?
## Process

| Off-line evaluation | Release in dev, ask for beta user feedback | Deploy |
|---|---|---|
| **Evaluate** iterations of the RAG system with **metrics** and store them in **Langfuse** | 1. **Duration:** typically a week | Presently, **no online evaluation** because we don't monitor what users ask about. |
| We evaluate each major iteration on an evaluation set: | 2. **Purpose:** | We plan to add **likes** with option to reuse the example to improve the performance. |
| • FAQs for policies<br>• Confluence: Encourage users to contribute **Q&As**.<br>• **User-generated data**: collect questions during user testing, feedback or support tickets. | • Allow users to explore new features/behaviors<br>• Gather feedback<br>• Collect evaluation data | |

# Evaluating Ada

## Scores cont'- LLM as a judge

**LLM as a judge scores**

- Context relevancy
- Answer relevancy
- Groundedness/ Faithfulness

These scores are normalised between 0 and 1.



The RAG Triad

**Answer Relevance:** Is the response relevant to the query?

**Context Relevance:** Is the retrieved context relevant to the query?

**Groundedness:** Is the response supported by the context?

# Evaluating Ada continued
## LLM as Judge: How the Metrics Work

These metrics are based on **prompting** and an **LLM call.**

These metrics are **sensitive to**:

- small changes in the prompt
- the order of the inputted information
- the length of inputted information
- the model used

Some python packages (Trulens, ragas) have them implemented.

```
context_relevance_system_prompt = """You are a RELEVANCE grader; providing the
relevance of the given CONTEXT to the given QUESTION.
        Respond only as a number from 0 to 10 where 0 is the least relevant and 10 is
the most relevant.
        A few additional scoring guidelines:
        - Long CONTEXTS should score equally well as short CONTEXTS.
        - RELEVANCE score should increase as the CONTEXTS provides more RELEVANT
context to the QUESTION.
        - RELEVANCE score should increase as the CONTEXTS provides RELEVANT context to
more parts of the QUESTION.
        - CONTEXT that is RELEVANT to some of the QUESTION should score of 2, 3 or 4.
Higher score indicates more RELEVANCE.
        - CONTEXT that is RELEVANT to most of the QUESTION should get a score of 5, 6,
7 or 8. Higher score indicates more RELEVANCE.
        - CONTEXT that is RELEVANT to the entire QUESTION should get a score of 9 or
10. Higher score indicates more RELEVANCE.
        - CONTEXT must be relevant and helpful for answering the entire QUESTION to get
a score of 10.
        - Never elaborate."""

context_relevance_user_prompt = """QUESTION: {question}
        CONTEXT: {context}
        Please answer with the template below for all statement sentences:
        Supporting Evidence: <Identify the sentences from the CONTEXT where the
information is most relevant to the QUESTION.>
        Score: <Output a number between 0-10 where 0 is no information overlap and 10 is
all information is overlapping>
    """
```

# RAG evaluation UI
## Langfuse

**Open-source observability platform for LLM applications. Helps track, analyze, and optimize AI interactions**

**Features:**

- Tracing and logging datasets and metrics via experiments
- Prompt management
- Performance monitoring
- Cost tracking

**Benefits:**

- Easy integration to the code.
- Centralized location for tracking performance metrics.
- Enhances transparency and reproducibility in AI evaluation workflows.
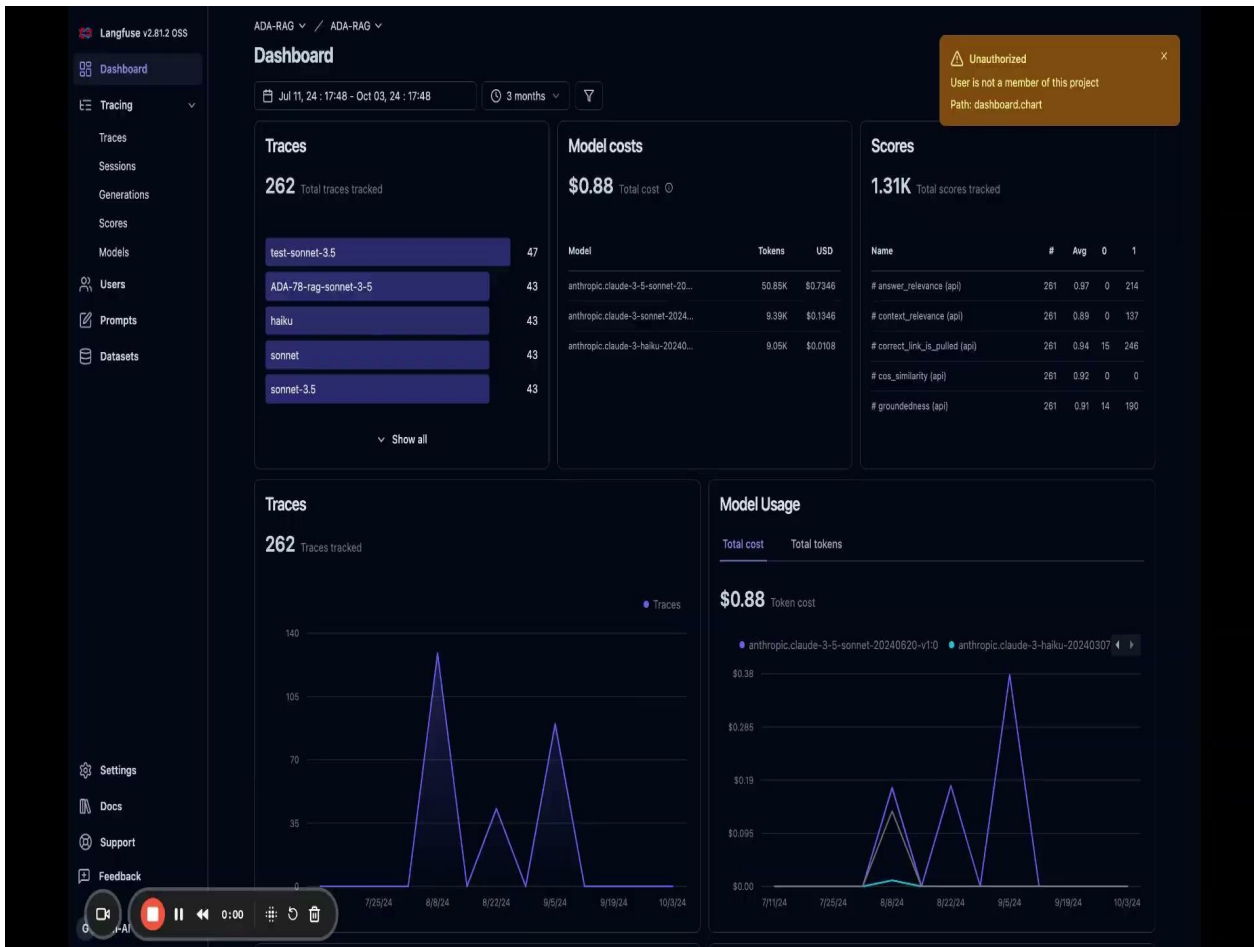
# RAG evaluation UI
## Langfuse

**Disadvantages of Langfuse:**

- Limited Customization
- No native filtering capacity
- Limited analytical features
- Not easy to analyse the results

**Alternatives to Langfuse:**

mlflow™

CLEAR|ML

comet

LangChain

wandb/weave

Weave is a toolkit for developing AI-powered
applications, built by Weights & Biases.

neptune.ai

# 04
## Wrap up

**Evaluation ensures RAG:**

· Improves response accuracy, style and behaviour.

· Handles domain-specific queries effectively

· Reduces hallucinations and factual errors.

**and identifies where we can improve the system**



**But, evaluation is difficult because of:**

·Growing knowledge bases complicate assessment

· Time-consuming manual review is often necessary

· Challenge to evaluate constantly



**And it's even more tricky because of:**

· Defining good metrics and aggregating them is challenging

Non-reproducibility and biases

· Difficulties to pinpoint where an error comes from due to interdependence

# Key takeaways
## What we have learned

- Text-similarity metrics are often inadequate for comprehensive evaluation.
- distinguishing between information retrieved and hallucinated by the model can be unclear.

- Usually, best is to **mix a variety of metrics**:
  - Rule based
  - Semantic based
  - LLM as a judge

- Effective strategies need to consider real-world contexts and use cases but getting good evaluation data can be challenging.
- Guardrails are essential

# Thank you!
# Questions?

leboncoin
tech

# RAG evaluation UI
## Langfuse - Dashboard



You can see the traces/ model usage for an interval…

But these are evaluation tests, filtering and interacting with the results could be more useful!

# RAG evaluation UI
## Langfuse

- You can access **Traces**/ **Sessions**/ **Generations**/ **Scores**/ **Prompts** and **Datasets**

- I find **Datasets** the most useful - here I can run many different versions of the same QA pairs and compare different solutions.

- But when you want to compare the results, it's not that easy…

| | Langfuse v2.78.0 OSS | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| ADA-RAG ∨ / ADA-RAG ∨ / Datasets / wave-question-answers | | | | | | | | | | |

**wave-question-answers**

| | □ (11/11) ∨ | | | Runs | Items |

| Name | Description | Run Items | Latency (avg) | Total Cost (avg) | # answer_relevance ... | # correct_link_is_pu... | # cos_similarity (api) | # groundedness |
|---|---|---|---|---|---|---|---|---|
| ADA-78-rag-sonnet... | | 43 | 14.82s | $0.0042 | Ø 0.9605 | Ø 0.9302 | Ø 0.9147 | Ø 0.8907 |
| haiku | | 43 | 12.42s | $0.0003 | Ø 0.9930 | Ø 0.9767 | Ø 0.9180 | Ø 0.9349 |
| sonnet3.5 | | 43 | 15.38s | $0.0041 | Ø 0.9571 | Ø 0.9286 | Ø 0.9153 | Ø 0.9095 |
| sonnet | | 43 | 15.46s | $0.0031 | Ø 0.9674 | Ø 0.9535 | Ø 0.9209 | Ø 0.9233 |

Tracing
Traces
Sessions
Generations
Scores
Models
Users

Dashboard

# RAG evaluation UI
## Langfuse - comparing experiments

- For an individual example it's easy to compare different experiments

- But you can't group experiments together in an easy way, can't order them by increasing scores…

- It stays a very human exercise to translate the findings. To understand why a score was low you'll need to access each generation individually.

| Run At | Dataset Item | Trace | Latency | Cost | # answer_relevance ... | # correct_link_is_pu... | # cos_similarity (api) | # groundedness (api) | # context_relevance... |
|---|---|---|---|---|---|---|---|---|---|
| 22/08/2024, 17:49:45 | 21c9e4a1-7d... | | 16.73s | $0.0046 | 1.0000 | 1.0000 | 0.9491 | 1.0000 | 1.0000 |
| 08/08/2024, 12:13:48 | 21c9e4a1-7d... | | 10.15s | $0.0002 | 1.0000 | 1.0000 | 0.9295 | 1.0000 | 1.0000 |

# RAG evaluation UI
## Langfuse - Individual trace

A trace contains the **question**, the **expected answer**, the **llm generated response** and anything else that you want to add to it:

- the retrieved documents
- the explanation for LLM generated scores
- Prompt

page_content: "with integrity in everything we do. Our reputation is important to us and is essential for our s... and ability to operate today and in the future.
We serve our colleagues, customers and users best when we communicate openly, speak freely and share relevant information. If you observe or suspect any conduct that concerns you or that seems to violate our Code of Ethical Conduct and/or our policies and/or procedures, you are encouraged to speak up.

groundedness_score_explanation: "Supporting Evidence: The source document "Speak Up and Investigations Policy V2 - Final version (EN).txt" provides extensive details that match the key points mentioned in the hypothesis. Specifically:

- Page 5 lists the types of concerns that can be raised through Speak Up channels, including fraud, bribery, conflicts of interest, data misuse, discrimination, etc., aligning with the hypothesis.
- Page 5 also mentions that Speak Up channels can be used to report any violation of laws or dangers

GENERATION sonnet3.5
08/08/2024, 11:59:15

Prompt: RAG_PROMPT_NO_UPLOADED_DOC_BUT_RETRIEVED_DOC - v1    Latency: 17.90s    6 prompt → 349 completion (Σ 355)

anthropic.claude-3-5-sonnet-20240620-v1:0    $0.0053

💬    ✏ Annotate    ➕ Add to dataset

Pretty ✨  JSON

**Input**

What is Speak Up?

**Output**

Speak Up is Adevinta's internal reporting system that allows employees to raise concerns about suspected misconduct or violations of the company's Code of Ethical Conduct, policies, or procedures. It's designed to encourage open communication and promote integrity within the organization.

TRACE  sonnet3.5    ⊕ ⊖
17.90s

GENERATION  sonnet3.5
17.90s  6 → 349 (Σ 355)

answer_relevance: 1.00    context_relevance: 0.90

correct_link_is_pulled: 1.00    cos_similarity: 0.90

groundedness: 0.90

# Evaluation with scores
## BertScore

ROUGE score measures how similar are the prediction and a possible ground truth but in a naive setup it can easily go wrong:

| Evaluating LLMs is a very difficult topic | Large Language Models are not easy to assess | Score: 0 |
|---|---|---|
| prediction | groundtruth | |

wut?

| Classification models are easy to assess | Large Language Models are not easy to assess | Score: 0.6 |
|---|---|---|
| prediction | groundtruth | n_gram=1 |

# Issues with standard metrics
## Why Traditional Metrics Fall Short

**Limited Scope, focus on surface-level text similarity:**

- Metrics like BLEU and ROUGE **ignore context and meaning.**
- Metrics like BertScore takes into account the semantics but is biased toward a "groundtruth".
- They tend to be **biased**, correlation between human judgement and some benchmark scores are low.*

**Lack of Factual Accuracy Measurement:**

- Unable to assess if the generated content is **factually correct or misleading**.
- No Consideration for **User Satisfaction.**
- Doesn't capture qualitative aspects like **coherence, fluency, or helpfulness**.

**Bias Toward Rigid Matching**

# RAG Evaluation

## The metrics should answer...

### Does the prediction answer the question correctly?

**Evaluation set with reference answers**

**Human evaluation**

**LLM as a judge**
(Answer relevance)

**Standard scores**
(ROUGE, METEOR)

**Semantic based scores**
(BertScore)

**Models/ LLM as a judge**

# RAG Evaluation
## The metrics should answer...

**Did we find the right information? (Retrieval metrics)**

No "reference context": LLM as a judge

Pairwise comparison of "reference context" and retrieved context (Precision, recall, F1 score, MRR, nDCG)

**Given the retrieved information, did we generate the correct answer? (Generation metrics)**

Generation based metrics - LLM metrics

Rule based metrics

LLM as a judge

# RAG-Specific Challenges 4.
## Hallucinations due to the LLM and or The Knowledge base

**The question isn't whether it will come, but rather when.**

Therefore a guardrail bloc has to be included in the architecture to avoid harmful or misleading answers

**<add schema here>**

# Use case 2: Markus - ..

•••

_____

- Retrieval part is important
- Prompt is very sensitive –
- *If possible fix the problems without changing the prompt.*
- Hallucinations are still a pain
- Non deterministic answers are still a challenge for evaluation and reproducibility
- Documents that are close semantically to the input are not always necessarily the most relevant.
- Constant "surprise" <-> continuous debugging
- You must have guardrails to mitigate those surprises (add scripts to detect hallucinations for example).
- Evaluation frameworks are new and not yet matures

## Takeaways and learnings

# Evaluating CCA RAG
## Process

| Off-line evaluation | User testing process | Deploy |

**Blabla** CCA

**Blabla** CCA

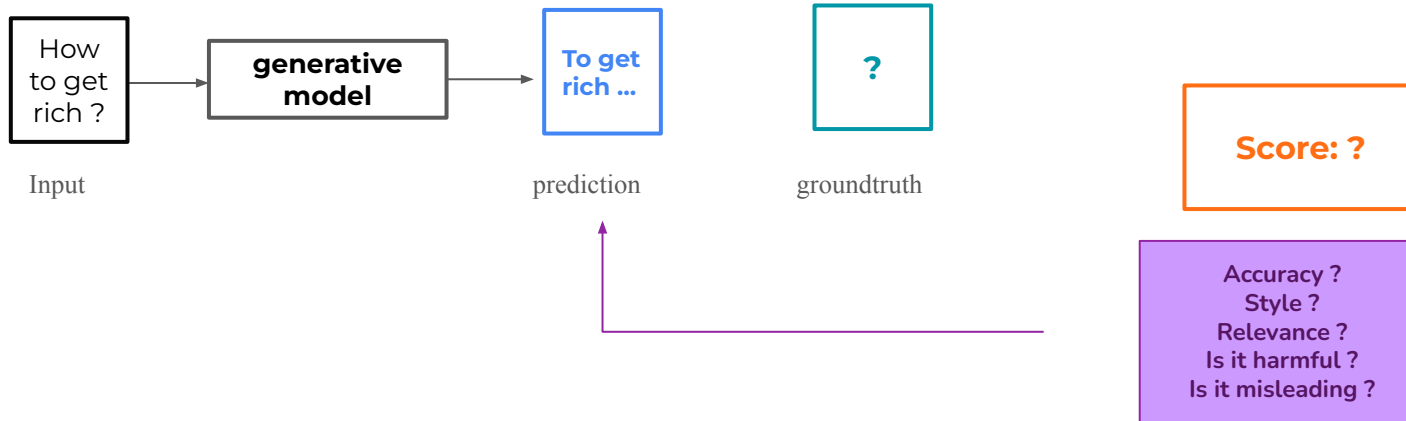**Blabla** CCA

# Challenges of LLM Evaluation
## Complexities in Evaluating LLMs

- **Output Variability:** Multiple valid responses possible for one input

- **Subjectivity:** Human and model judgments can be inconsistent and biased

- **Context Sensitivity:** Subtle input changes may significantly alter outputs

- **Non-determinism:** Outputs may vary across runs for the same input

- **Continuous Evaluation:** Needed due to constant model improvements

# Evaluate something that is generated
## How would you do that ?

# Use case 1: Markus - Customer Care Assistant

**leboncoin customer service total tickets : ~200k/month**

Many of these tickets could be automated using LLMs and leboncoin's FAQ to ease the burden on customer service.

Je souhaite modifier l'adresse email de mon compte par ▮▮▮▮▮.com. Je n'ai plus accès a l'ancienne

**Data_LBC_Admin**

Bonjour,

Pour modifier l'adresse e-mail associée à votre compte lorsque vous n'avez plus accès à l'ancienne, la procédure est expliquée ici

Il faut contacter notre service client en fournissant les informations suivantes :
- Vos nom et prénom
- Votre adresse postale
- Votre numéro de téléphone
- L'adresse e-mail initialement enregistrée sur votre compte
- Votre nouvelle adresse e-mail (▮▮▮▮▮k.com)

Nous pourrons alors procéder manuellement au changement d'adresse sur votre compte.

Ce message a été généré automatiquement. N'hésitez pas à y répondre si la réponse ne vous convient pas, un agent vous répondra.
En savoir plus

LangChain

FAISS

Amazon Bedrock

# Evaluating Ada
## Scores used

**Rule-based scores: make sure that the answer is in a certain way (dependant of the use case)**

The end message should always finish with a caution *"Be aware that I can make mistake. I found the source here..."*

Check if source is in response

**BertScore: use reference responses, embed prediction and reference response, compare the similarity**

**Compare the stability of a response in several runs by comparing their embeddings**

**Retrieval metrics: check whether correct page is pulled:**
**precision/ recall/ F1 scores**

**LLM as a judge scores**