

# GPTx und RAG in der Praxis

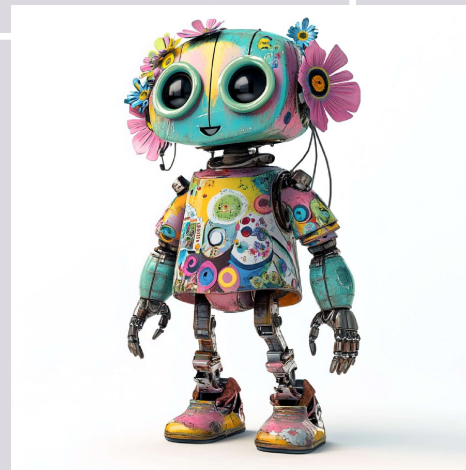
**Schluss mit Prototyp**

Christian Hidber  
Oliver Zeigermann

data2day, Heidelberg, September 2024

Chef:

mein Enkel kann das auch...

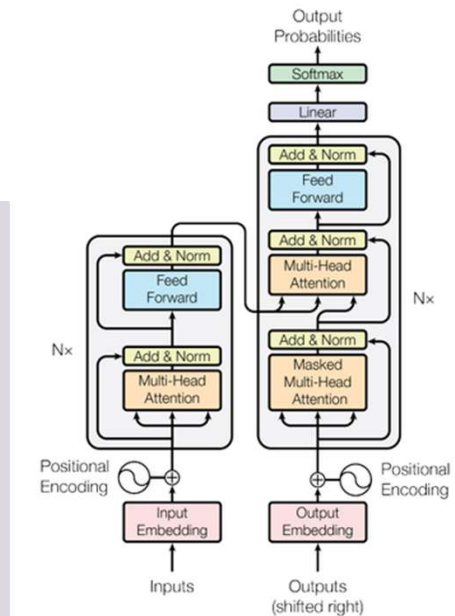


# LLM Intro

---

## Transformers, LLMs, Encoder, Decoder: WTF?

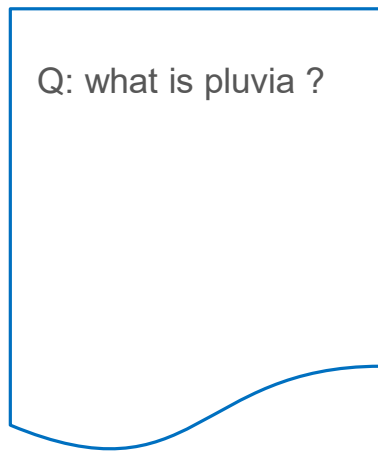
- **Transformers:** A flexible architecture that uses self-attention to process sequential data efficiently.
- **LLMs:** Large-scale Transformer models trained on extensive text datasets to perform various language tasks.
  - **Encoder Models:**
    - Part of the Transformer architecture focused on understanding and interpreting input data (e.g. *BERT*)
    - Instrumental for Embedding Models
  - **Decoder Models:**
    - Part of the Transformer architecture focused on generating sequential output based on the interpreted inputs or prior outputs
    - Instrumental for GPT-style Models like **Llama, Mistral or OpenAI GPT**



# Decoder Models

---

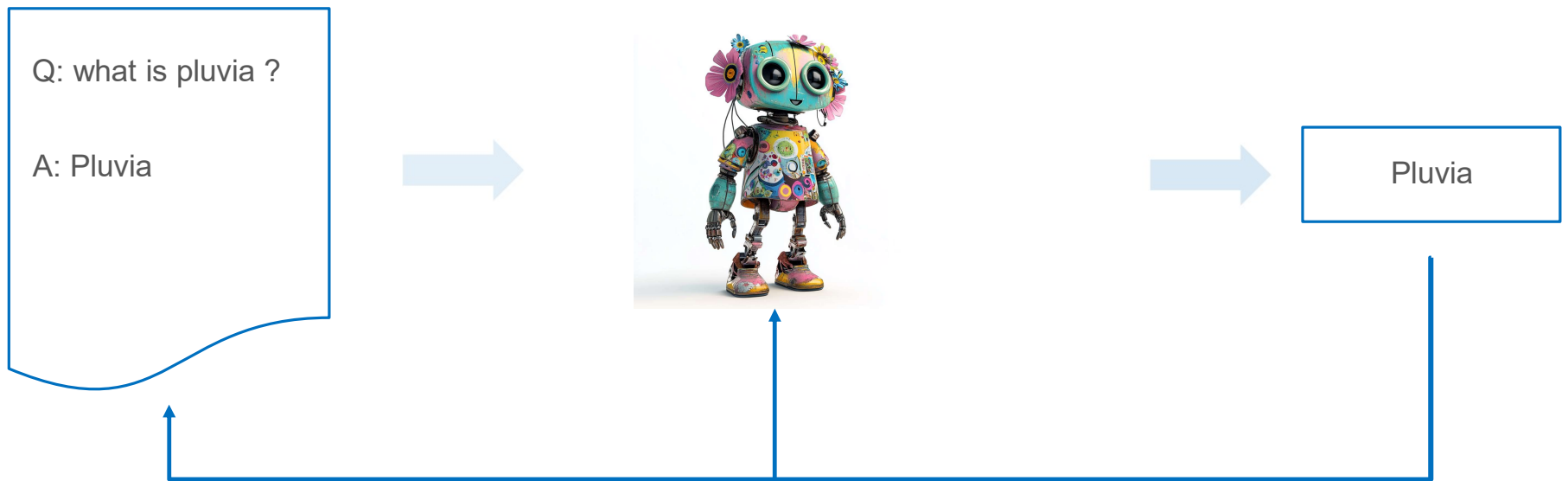
## How does a Decoder Model work ?



- Depends on users goal
- Unique for each chat & user
- Contains the chat history
- «**the context**»

- Trained on huge datasets
- Does not change
- Same for all users
- «**the model**»

## How does a Decoder Model work ?

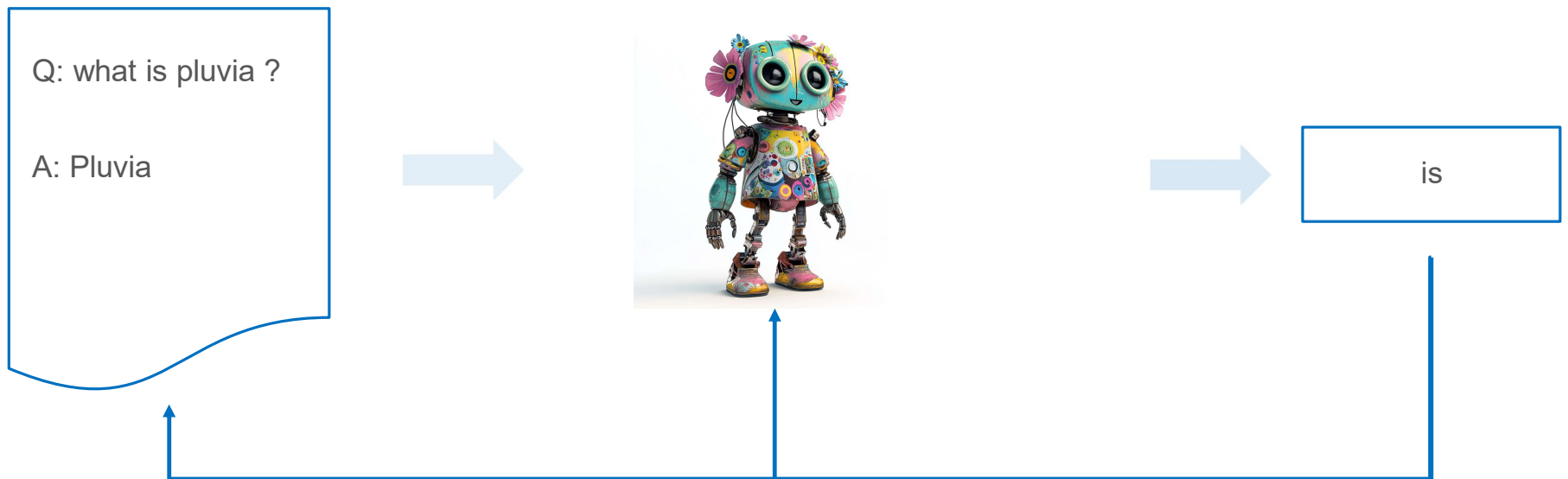


- Depends on users goal
- Unique for each chat & user
- Contains the chat history
- **«the context»**

- Trained on huge datasets
- Does not change
- Same for all users
- **«the model»**

- Single «word»
- Depends on context and model
- **«the token»**

## How does a Decoder Model work ?

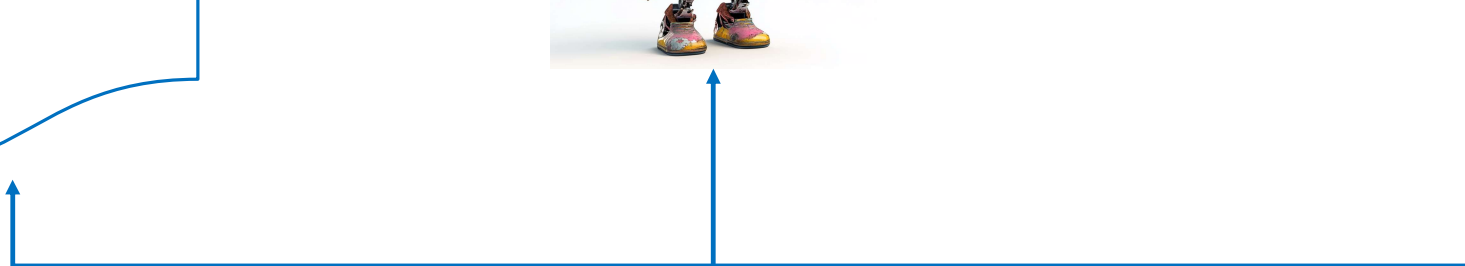
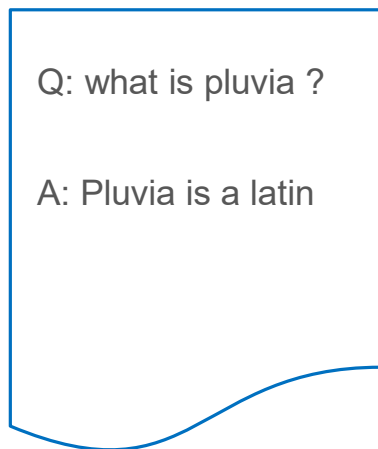


- Depends on users goal
- Unique for each chat & user
- Contains the chat history
- **«the context»**

- Trained on huge datasets
- Does not change
- Same for all users
- **«the model»**

- Single «word»
- Depends on context and model
- **«the token»**





- Depends on users goal
- Unique for each chat & user
- Contains the chat history
- «the context»

- Trained on huge datasets
- Does not change
- Same for all users
- «the model»

- Single «word»
- Depends on context and model
- «the token»

## How does a Decoder Model work ?



- Depends on users goal
- Unique for each chat & user
- Contains the chat history
- «the context»

- Trained on huge datasets
- Does not change
- Same for all users
- «the model»

- Single «word»
- Depends on context and model
- «the token»

## Naïve Approach

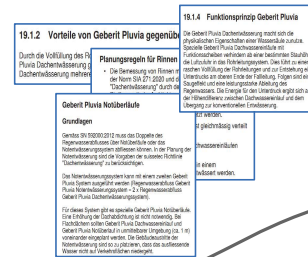


Idea: just a few pages

You are an expert in .....

## What is Pluvia ?

Use the following facts:



User:  
Asking a Question



answer



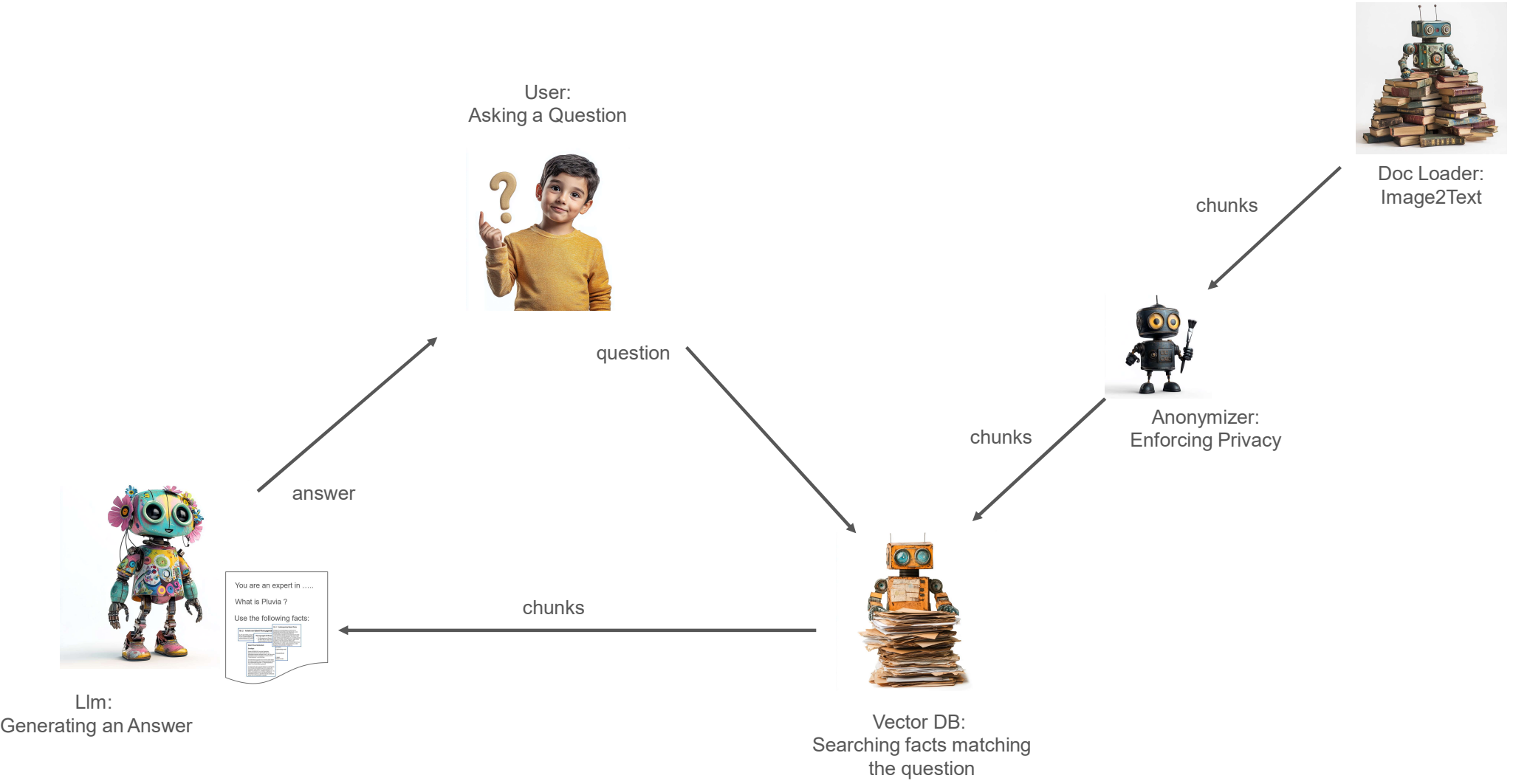
Llm:  
Generating an Answer

# **RAG**

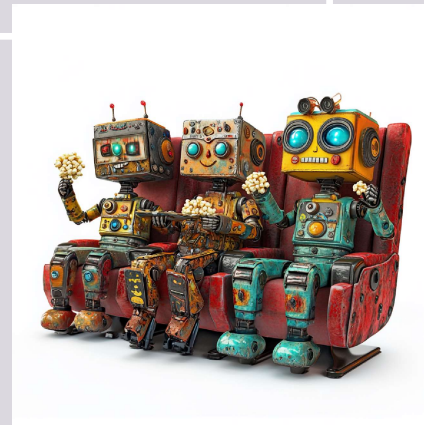
## **Retrieval Augmented Generation**

---

# RAG System Architecture



Demo:  
Low Risk RAG Applications



# Choosing an application

---



Low Risk, but nice benefit

### **Low Risk**

- What is the worst thing that could happen and how to mitigate that?
- Low profile
- Failures should be ok
- Human in the loop

### **Nice Benefit**

- Impossible to do by humans or
- Humans don't like to do
- Let the whole organization learn
- Management likes it, but is afraid
- Can it be used for (internal) marketing?

# **From Prompt Hacking to Production**

---

## The Small Handyman vs. Engineer Task

**Ad-hoc prompting is something very different from writing a prompt for a service**

### **With ad-hoc prompting**

- you can immediately see if it works.
- there's a high level of human oversight.
- it only needs to work for a specific example

### **With prompting for a system**

- It needs to generalize for all expected use cases
- Has no or less human supervision
- Stability is expected



Talking about  
stability:  
**Evaluation**

# Evaluation

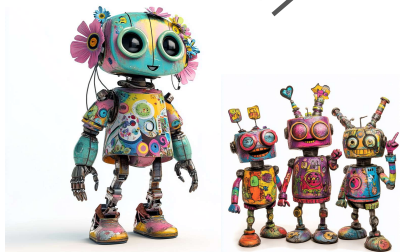
---

## Evaluation on text results

User:  
Asking a Question



answer



Llm:  
Generating an Answer

Human Eval

Question

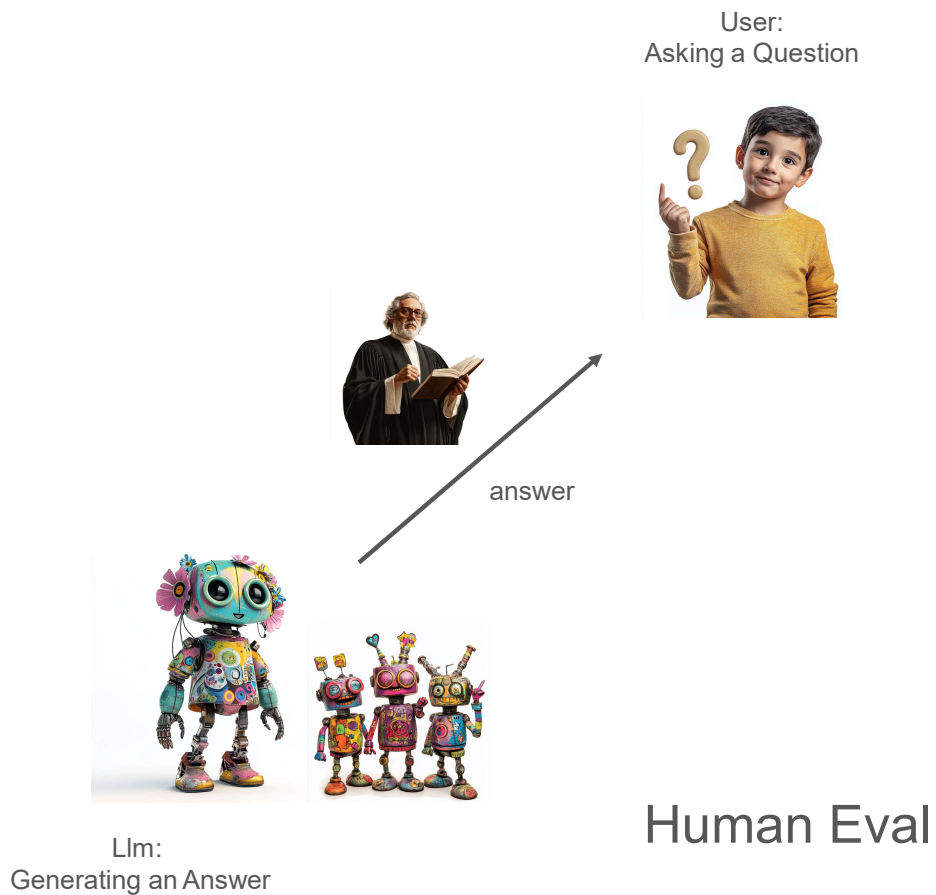
- What is Pluvia ?

Answer

- Pluvia is a latin word meaning rainfall.
- The latin word for rainfall.
- ....

=> equality not an option

## Evaluation on text results

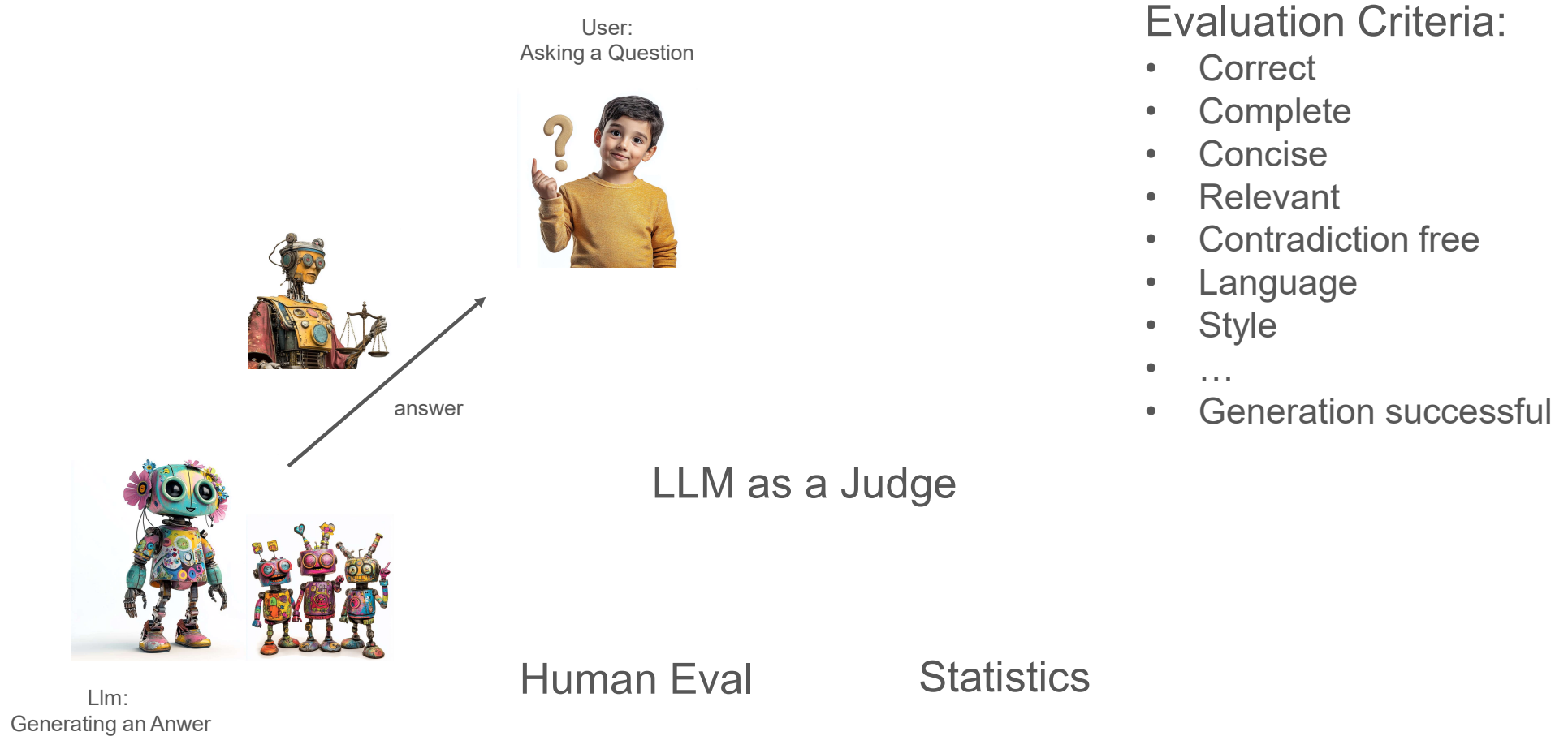


### Evaluation Criteria:

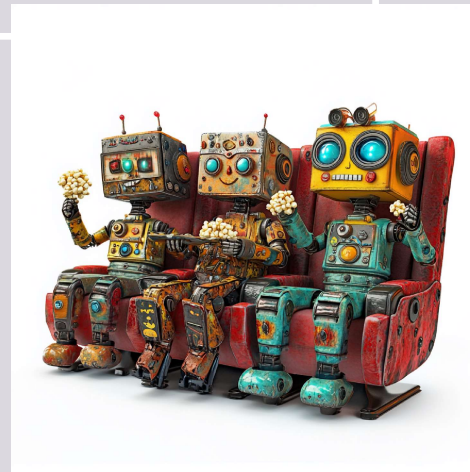
- Correct
- Complete
- Concise
- Relevant
- Contradiction free
- Language
- Style
- ...
- Generation successful

Statistics

## Evaluation on text results



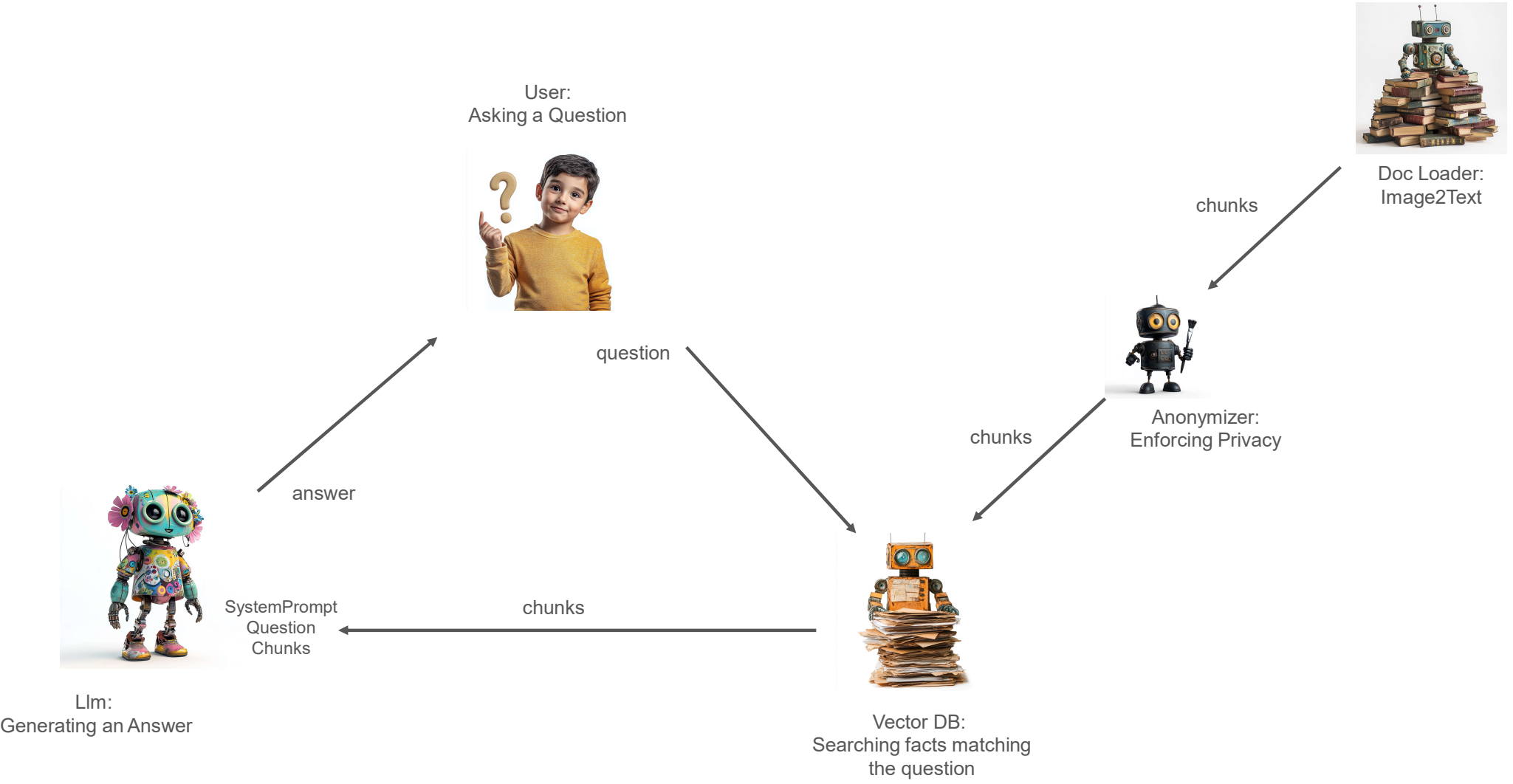
## Demo: Evaluation Notebook



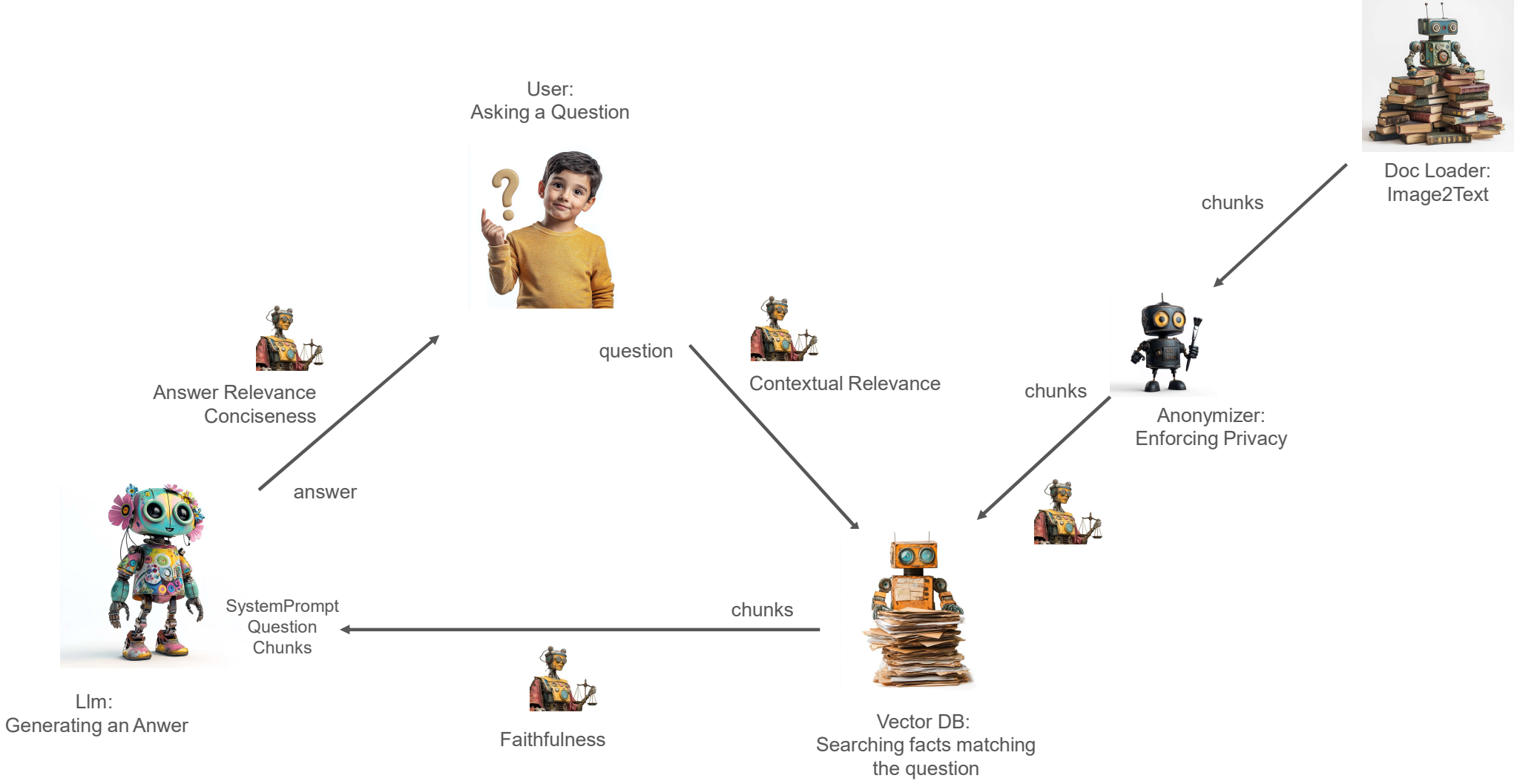
<https://colab.research.google.com/github/DJCordhose/llm-from-prototype-to-production/blob/main/Eval4pptx.ipynb>



# RAG System Architecture



# RAG System Architecture: Online Evaluation



## Online Eval: Example

```
[07:56:37 INF] POST "https://[REDACTED].azurewebsites.net/"eval succeeded in 4.91s with response={
  "Metadata": {
    "Answer": "Der Artikel Sigma20 BetÄmtigungsplatte dient zur Steuerung der 2-Mengen-SpÄhlung bei Geberit UP-SpÄlkÄsten. Sie ermÄtzt",
    "CreateDate": "2024-09-10T07:56:37.744166Z",
    "DeepEval": {
      "Answer_Relevancy": {
        "reason": "The score is 1.00 because the response directly addresses the purpose of the Sigma20 BetPl. article without any irre",
        "score": 1.0
      },
      "Conciseness_(GEval)": {
        "reason": "The output is somewhat concise but includes unnecessary details about materials and suitability that could be omitted",
        "score": 0.6
      },
      "Contextual_Relevancy": {
        "reason": "The score is 0.33 because the context discusses various models and specifications of flushing systems but does not p",
        "score": 0.3333333333333333
      },
      "Faithfulness": {
        "reason": "The score is 0.80 because the actual output inaccurately generalizes the material of the BetÄmtigungsplatte, stating",
        "score": 0.8
      }
    },
    "ElapsedSeconds": 4.85,
    "EvalType": "deep_eval",
    "EvalVersion": "240903",
    "Input": "Wozu dient der Artikel Sigma20 BetPl., fÄr 2-Mengen-SpÄhlung weiÄ / weiÄ matt ?"
  },
  "Metrics": {
    "Answer_Relevancy": 1.0,
    "Conciseness_(GEval)": 0.6,
    "Contextual_Relevancy": 0.3333333333333333,
    "Faithfulness": 0.8
  },
  "Score": 0.6833333333333333
}
```

## Online Eval: Example

```
LlmDesc germany W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 242,242]
LlmDesc germany W.240820_C.240625_E.240903: Answer_Relevancy=0.962 Conciseness_(GEval)=0.628 Contextual_Relevancy=0.341
LlmDesc switzerland W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 140,140]
LlmDesc switzerland W.240820_C.240625_E.240903: Answer_Relevancy=0.943 Conciseness_(GEval)=0.613 Contextual_Relevancy=0.534
LlmFp germany W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 4,4]
LlmFp germany W.240820_C.240625_E.240903: Answer_Relevancy=0.964 Conciseness_(GEval)=0.595 Contextual_Relevancy=0.739
LlmFp switzerland W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 8,8]
LlmFp switzerland W.240820_C.240625_E.240903: Answer_Relevancy=0.984 Conciseness_(GEval)=0.624 Contextual_Relevancy=0.621
LlmHelp german W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 66,66]
LlmHelp german W.240820_C.240625_E.240903: Answer_Relevancy=0.992 Conciseness_(GEval)=0.692 Contextual_Relevancy=0.731
LlmSi germany W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 6,6]
LlmSi germany W.240820_C.240625_E.240903: Answer_Relevancy=0.987 Conciseness_(GEval)=0.615 Contextual_Relevancy=0.827
LlmSi switzerland W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 8,8]
LlmSi switzerland W.240820_C.240625_E.240903: Answer_Relevancy=0.985 Conciseness_(GEval)=0.623 Contextual_Relevancy=0.904
LlmTtp english W.240820_C.240625_E.240903: Answer_Relevancy=1.000 Conciseness_(GEval)=0.500 Contextual_Relevancy=0.000
LlmTtp french W.240820_C.240625_E.240903: Answer_Relevancy=1.000 Conciseness_(GEval)=0.600 Contextual_Relevancy=0.000
LlmTtp german W.240820_C.240625 : Score=0.000 TextGenerated=0.000 [counts 74,74]
LlmTtp german W.240820_C.240625_E.240903: Answer_Relevancy=0.925 Conciseness_(GEval)=0.611 Contextual_Relevancy=0.442
[09:21:55 INF proplanner] HTTP POST /api/descriptions responded 200 in 25.1855 ms
```

## Evaluation Issues

- Online Performance impact on LLM
  - Eval may call 10x more often, but have less output tokens
- Which LLM do you use ? Same ? Faster ? Most Powerful ?
- What Dimensions do you eval ?
  - Toxicity, Conciseness, Answer Relevance ?
  - Ground Truth available ?
- Human Feedback from your users ?
- Interpretation of the Scores ?

## Eval Frameworks

- **DeepEval** <https://docs.confident-ai.com/>
- Ragas <https://ragas.io/>
- TruLens <https://www.trulens.org/>
- Evidently <https://www.evidentlyai.com/>
- Ares <https://ares-ai.vercel.app/>
- ...

# **Wrap Up**

---

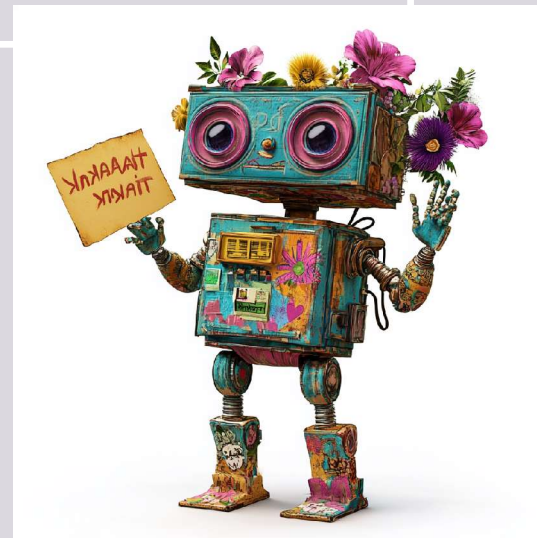
## Key takeaways

- GPT-style models and RAG are the key to a new era of applications
- Choose low-risk, nice-benefit applications (first)
- Ad-hoc prompting is different from prompting for a system
- Human Eval is a great starting point
- LLM-as-a-Judge works, but take the scores with a grain of salt
- Use a strong LLM for evaluation
- Getting the Documents & keeping them up-to-date can be painful

**Vorsicht vor dem Enkel des Chefs...**



Thank you



## Llm-as-a-judge: Idea

Actual Output:

```
Witing texts is painful,  
caus im making mitakes.
```

Prompt:

```
You are an expert on  
english language. Grade  
a students text...
```

```
Answer with a Json  
containing scores &  
reason..
```

```
Students Text:  
Witing texts is...
```

```
{  
  "score": 2,  
  "reason": "Multiple grammatical errors  
            such as 'witing' and ..."  
}
```

## Your Experience ?

- Anyone doing RAG ? In Production ?
- Do you do evaluation ? By humans ?
- What else do you use for evaluation ?

# Llm-as-a-judge: G-Eval

