

# How hard can it be to build a customer support chatbot? - Lessons learned

AWS Deep Dive Days –  
Generative AI

Leonard Zucht (HDI), Elias Stiegelmeier  
(HDI), Akarsha Sehwal (AWS)  
Munich | 20. & 21.02.2024

# The vision – „Alle Infos auf einen Blick“



**Knowledge-based customer service**



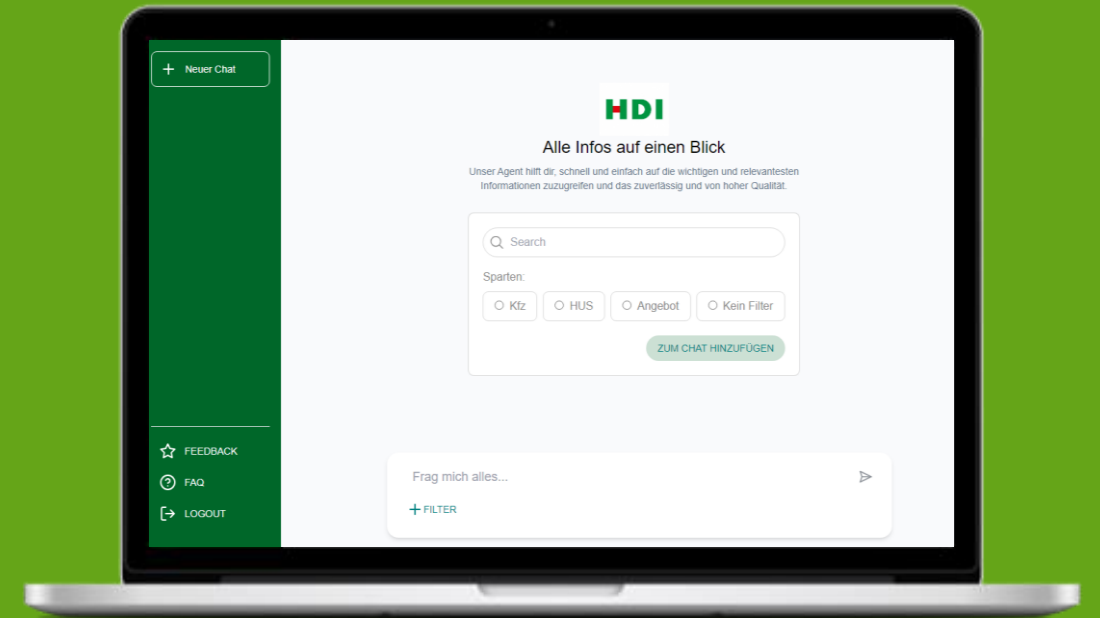
**Fast and flexible access to information**



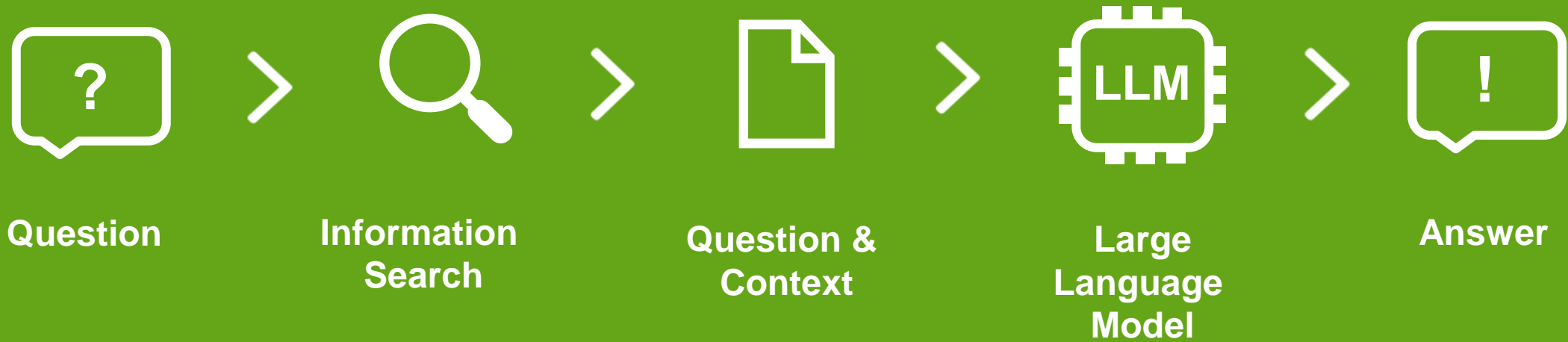
**Chatbot-interface**



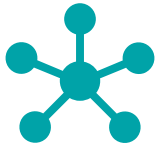
**Start with internal-only scope and scale**



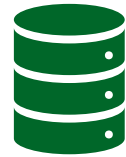
# Retrieval Augmented Generation (RAG)



# Setting up a first prototype is easy



Langchain



FAISS



Streamlit

<https://www.langchain.com/>  
<https://ai.meta.com/tools/faiss/>  
<https://streamlit.io/>



# Setting up a first prototype is easy

## The next steps



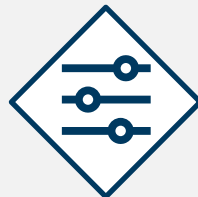
Scalability



Professional frontend



Improved document preprocessing  
for tables and layout



Modularize solution for optimization

# We need a product team






Build a scalable solution which solves the limitations and can be adapted in a modular manner


# The journey



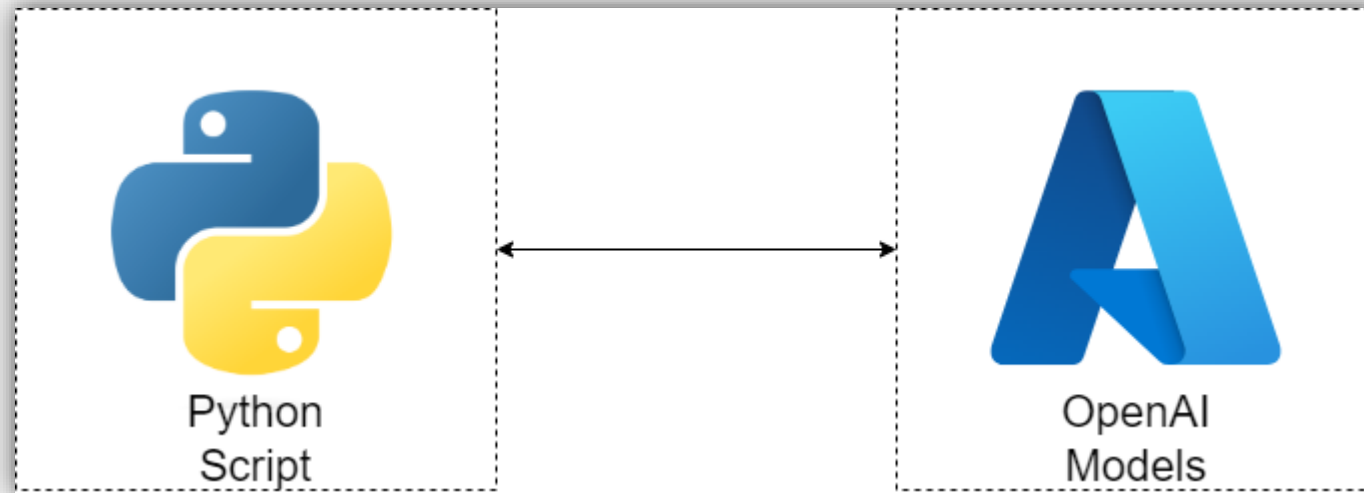
## Backend

	Simple ingest	Simple ingest	OCR-based ingest	Adapted OCR-based ingest
	In-memory retrieval	OpenSearch-based retrieval	OpenSearch-based retrieval	Adapted OpenSearch-based retrieval
	Local	Public-facing	Public-facing	HDI-internal VPC

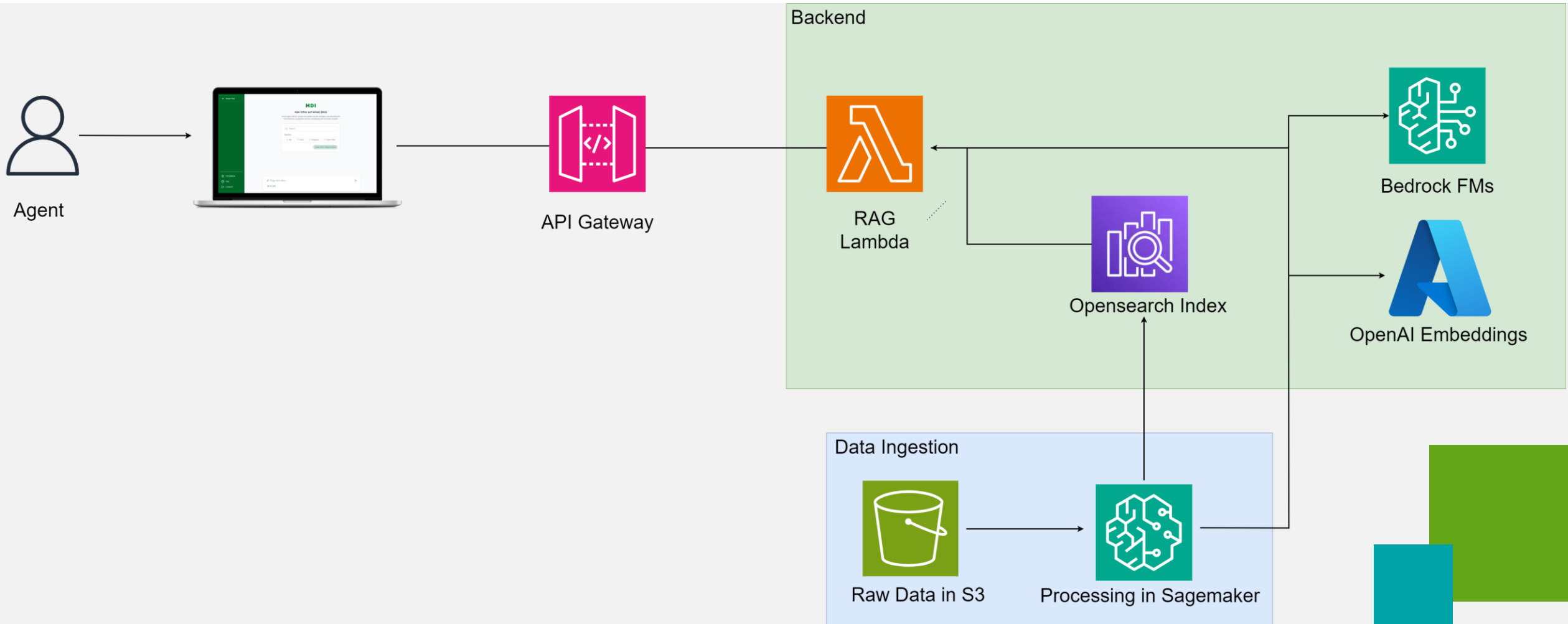
## Frontend

	Simple design (Streamlit)	Simple design (Gradio)	Simple design (Streamlit)	Professional design (React) Streaming enabled
---	---------------------------	------------------------	---------------------------	---

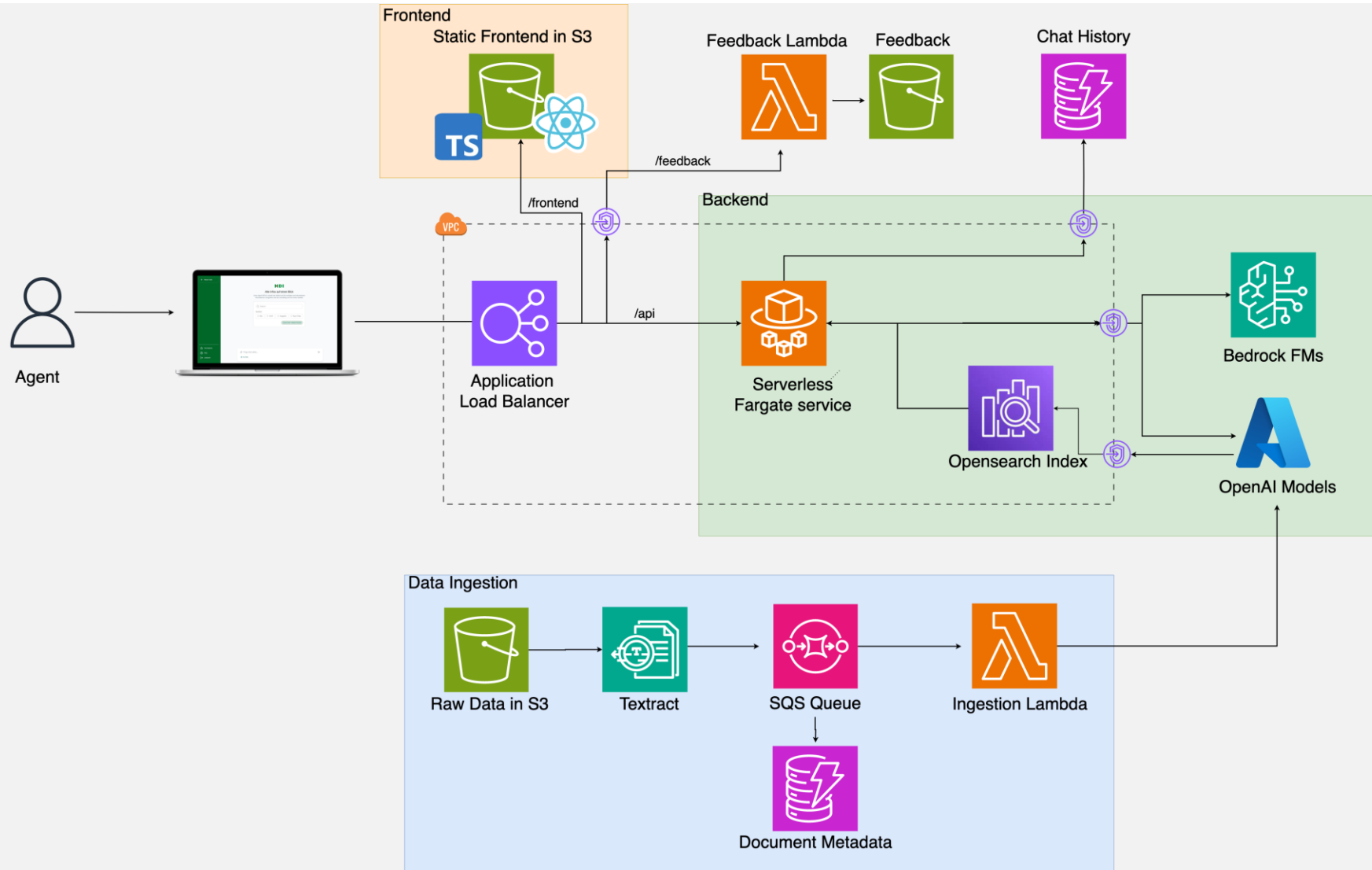
# The journey - Streamlit



# The journey - blogpost



# The journey – own application



# There are many moving parts

Challenges	Contents	Approaches
Complex document parsing	Tables, layout, pictures	OCR, ML-Parser
Chunking	Columns, paragraphs, lists, tables	Layout
Domain knowledge	Domain language, abbreviations, conventions	Prompt engineering, look-ups, finetuning
Variable parameters	Top-k, chunk-size, relevance, overlap, temperature, ...	Performance test
Document similarity	Retrieval, embedding, sources	Filter, hybrid-search, reranking
Performance	Recall, answer quality	Performance test, advanced prompting
User intent	Information search, summarization, follow-up questions	Prompt engineering, routing, agents / tools





# Spotlight – third party dependencies

## Challenges

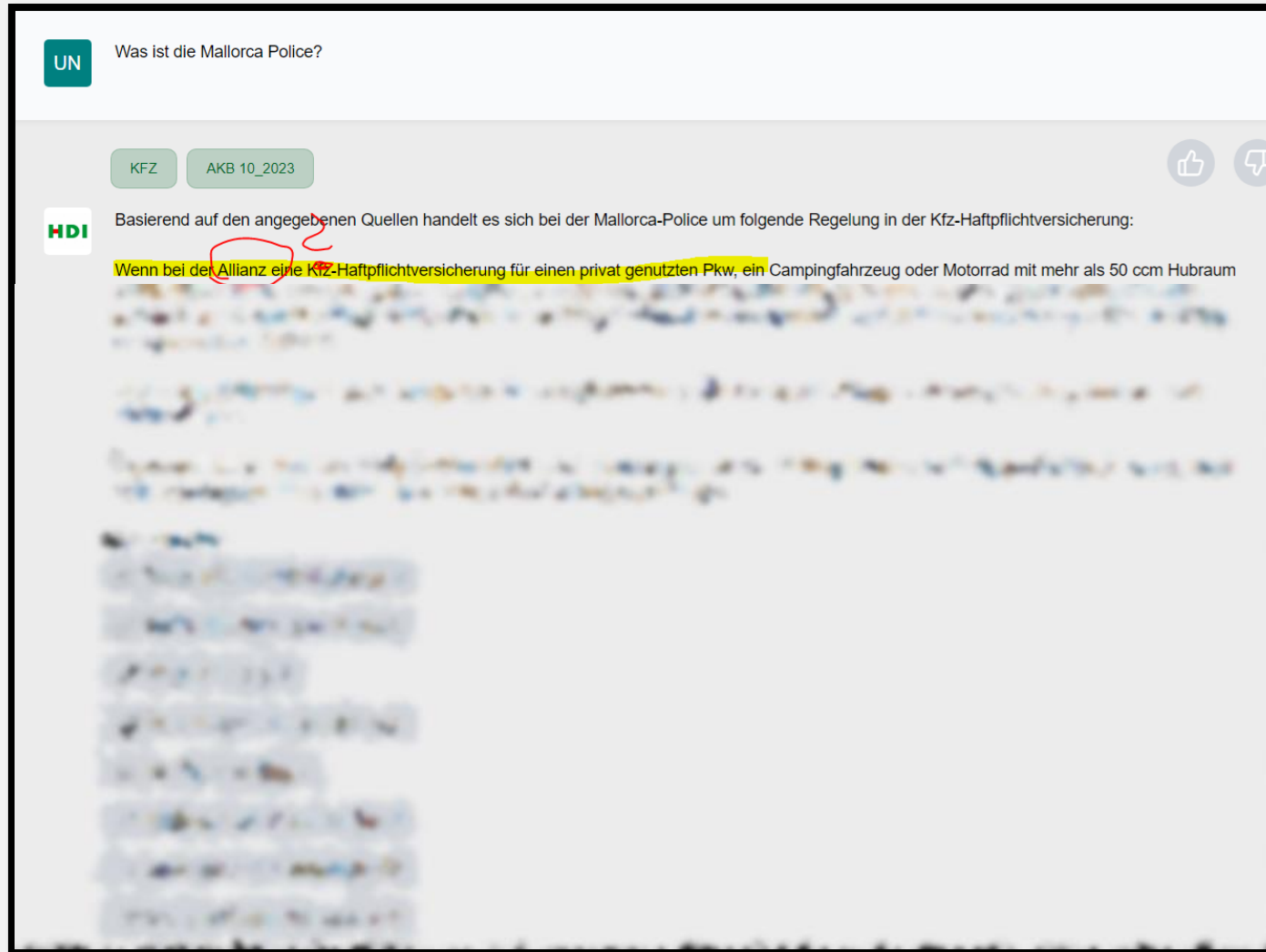
- Inconsistent
- Incomplete documentation
- Little transparency e.g. for timeouts



## Learnings

- Not everything popular is good.
- Rather prefer low level code (customizations and flexibility as needed)

# Spotlight - one prompt to rule them all?



# Spotlight - one prompt to rule them all?

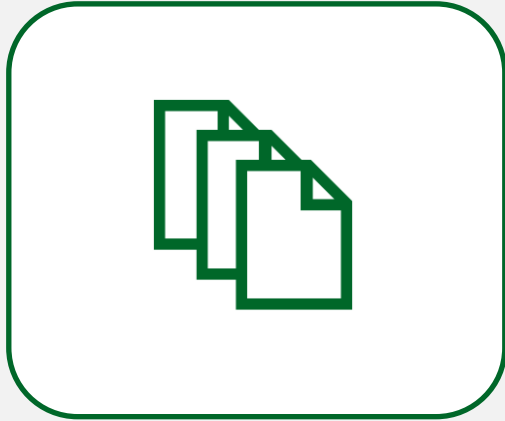
- Different system prompt can be useful for questions with differing complexity
- System prompt sometimes does not work as assumed
- How do you want your answers to look like?
- What should be included in the system prompt?
- Do you need to prevent prompt injection?



# Spotlight - one prompt to rule them all?

- Different system prompt can be useful for questions with differing complexity
- System prompt sometimes does not work as assumed
- How do you want your **Compromise**
- What should be included in the system prompt?
- Do you need to prevent prompt injection?

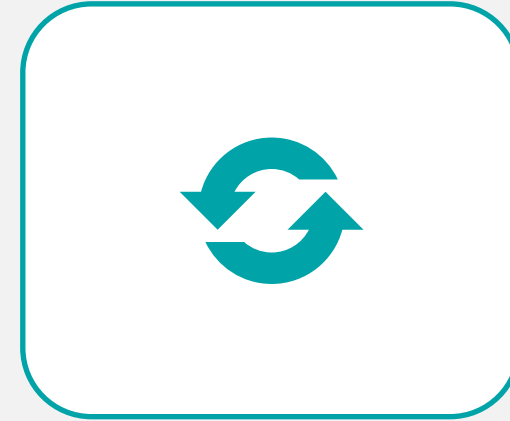
# How to test our solution? - The setup



**Ground truth**  
(Questions, Sources,  
Answers)

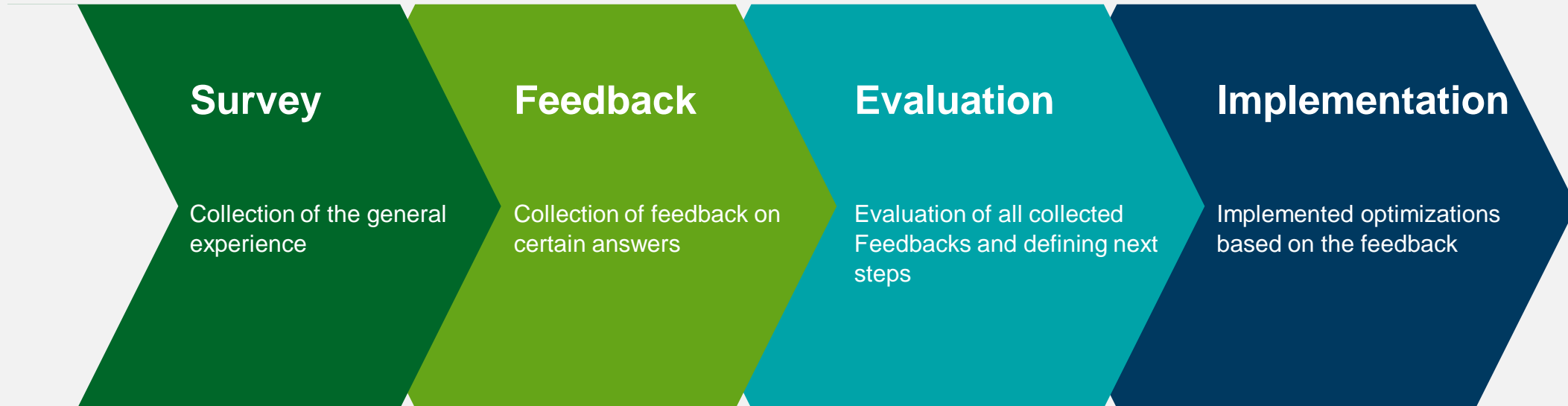


**Business Testers**

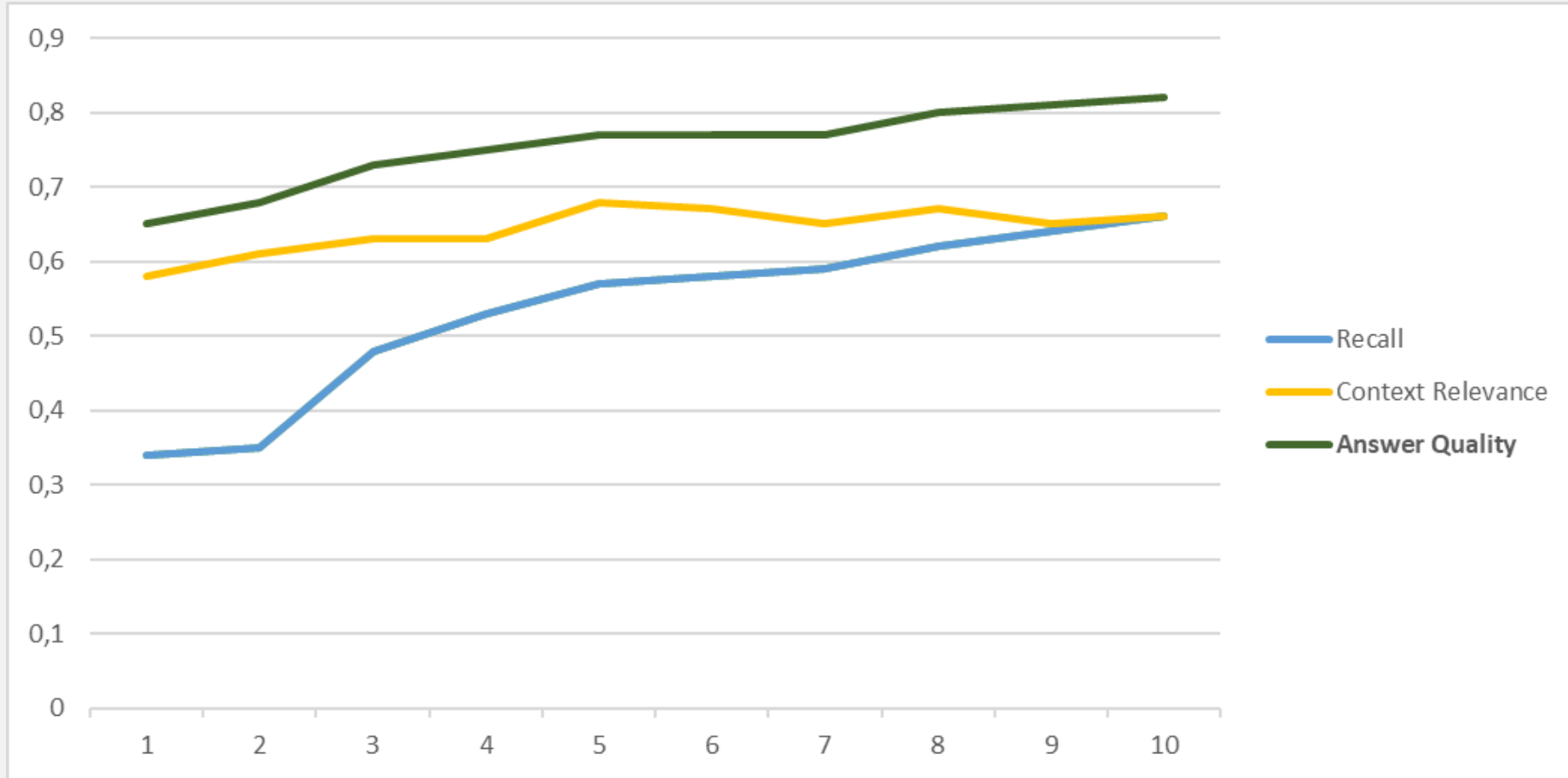


**Feedback Loop  
&  
Incremental  
Improvements**

# How to test our solution? - Feedback loop



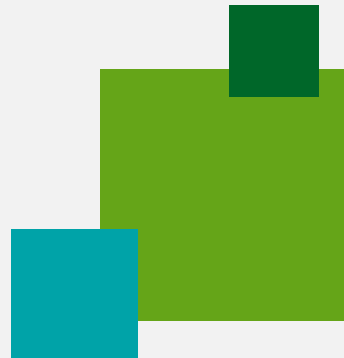
# Our results (so far)



# Key take-aways

## **It is not easy, it is a journey!**

- User feedback and evaluation matters
- Testing pipeline and experiments are helpful
- Start small and think big
- Modular approach for a fast-changing landscape
- Hallucinations need to be monitored to some extent.



# Outlook

## Improve data quality



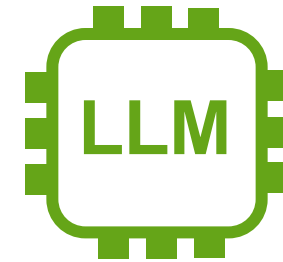
Database is refined

## Connection to further data



Relevant data

## Finetuning and hosting of own models



Models tailored to use-case

# Thank you for your attention!

## Questions?

