~~SiloX-GPT~~

# Content Search

## Datensilos verbinden mittels Multi-Agenten-Systemen

Data2Day 2024

# Moritz Glauner

Head of Data Science @ Bertelsmann

in moritzglauner

# Nico Kreiling

Data Scientist @ scieneers

Host of techtiefen.de     in nico-kreiling

# Why classical RAG was not an option…

**€20.2 Billion**
Group revenues

**€3,119 Million**
Operating EBITDA

**€1,326 Million**
Group profit

**80,418**
Employees

**Gütersloh**
Headquarters

**No central**
Data Storage

# At our starting-point, is was not yet clear, that LLMs can work with Knowledge-Graphs (very well)
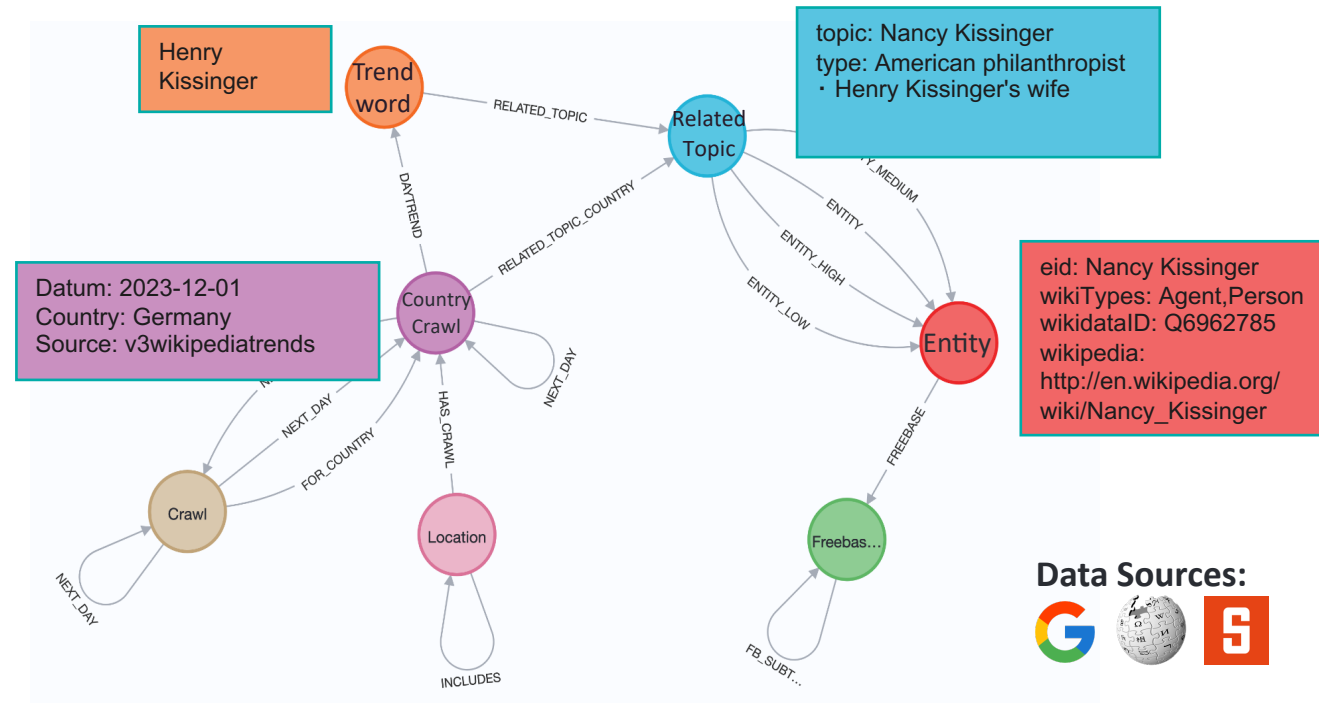
**Initial Situation:**

➢ Existing knowledge graph from predecessor project in Neo4J

**Research-Question:**

➢ Can natural questions be answered by an LLM and a Knowledge-Graph?

**Learnings:**

- GPT4 can write decent cypher querries when prompted with schema and some examples
- Simple-intuitive schemas work better then complex-formel
- We want to have a natural-language answer



**Data Sources:**

# Agents enable LLMs to do real work

**Agents give LLMs the ability to interact with the world by executing programs**

- Minimal version: Simple function-call with one parameter
- Complex version: Browsing or writing and testing own source code



| User |
|------|
| Welche Trends stehen in Bezug zu Angela Merkel? |

| LLM |
|------|
| Trigger function: `execute_cypher` With `cypher="MATCH (c:CountryCrawl)-[d:DAYTREND]->(t:Trendword)` `WHERE t.trend= 'Angela Merkel'` |

| Tool |
|------|
| [{year: 2015, topic: Flüchtlingskriese}, {year: 2020, topic: Corona Pandemie}] |

| LLM |
|------|
| Relevante Trends in Bezug auf Angela Merkel und ihre Kanzlerschaft waren:<br>* Flüchtlingskrise: 2015 traf Merkel die kontroverse Entscheidung, Deutschlands Grenzen für Flüchtlinge offen zu halten.<br> * Corona-Pandemie: In ihrer letzten Amtszeit war die Bewältigung der COVID-19-Pandemie(2020) |

Diagram labels: User, LLM, Tools, Outside world

Trends related to Angela Merkel

Angela Merkel related to the trends...
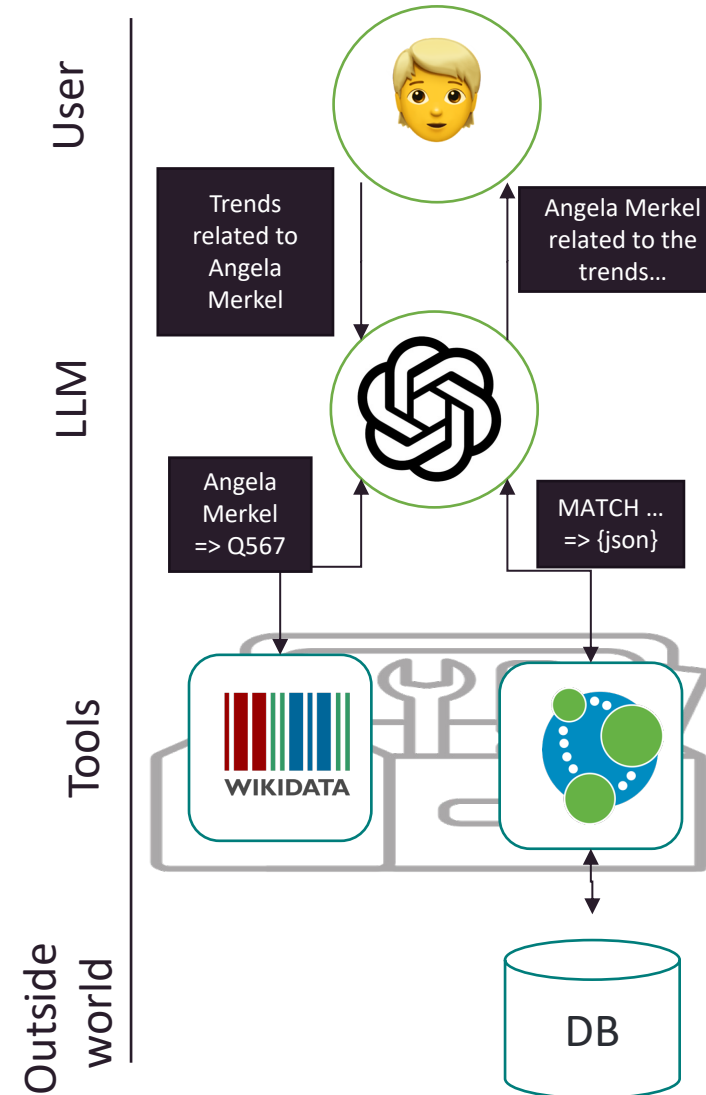
MATCH .. => {}

DB

# Agents enable LLMs to do real work

**Agents give LLMs the ability to interact with the world by executing programs**

- Minimal version: Simple function-call with one parameter
- Complex version: Browsing or writing and testing own source code



| User | |
|---|---|
| Welche Trends waren besonders relevant für die Kanzlerschaft von Merkel? | |

| LLM | |
|---|---|
| Trigger function: `execute_cypher` With `cypher="MATCH … WHERE t.trend= 'Kanzlerschaft Merkel'` | |

| Tool | |
|---|---|
| [] | |

| LLM | |
|---|---|
| Trigger function: `execute_cypher` With `cypher="MATCH … WHERE t.trend= 'Kanzler'` | |

| Tool | |
|---|---|
| [{topic: Angela Merkel}, {topic: Olaf Scholz}] | |

| LLM | |
|---|---|
| Angela Merkel war Kanzlerin der Bundesrepublik Deutschland. Olaf Scholz ist ihr Nachfolger. | |

# Agents enable LLMs to do real work

**Agents give LLMs the ability to interact with the world by executing programs**

- Minimal version: Simple function-call with one parameter
- Complex version: Browsing or writing and testing own source code

**Learnings:**

- Agents are much more powerfull then simple LLMs
- Important to keep the focus of an agent



User
Welche Trends waren besonders relevant für die Kanzlerschaft von Merkel?

LLM
Trigger function: `lookup_wikidata_id`
With `name="Angela Merkel"`

Tool
Q567

LLM
Trigger function: `execute_cypher`
With `cypher="MATCH (c:CountryCrawl)-[d:DAYTREND]->(t:Trendword) WHERE t.wikidata_id = 'Q567'`

Tool
[{year: 2015, topic: Flüchtlingskriese}, {year: 2020, topic: Corona Pandemie}]

LLM
Relevante Trends in Bezug auf Angela Merkel und ihre Kanzlerschaft waren:
* Flüchtlingskrise: 2015 traf Merkel die kontroverse Entscheidung, Deutschlands Grenzen für Flüchtlinge offen zu halten.
 * Corona-Pandemie: In ihrer letzten Amtszeit war die Bewältigung der COVID-19-Pandemie(2020)

# Agents are similiar to humans: Specialication and Collaboration increases the scope of possibilities – and the complexity

**Using multiple agents increases flexibility of the system**

- Use different LLMs for different tasks
- Have clear and focused system-prompts
- Control the length of the chat-history
- Gain more

**But they also increase the complexity**

- Requires communication accross agents
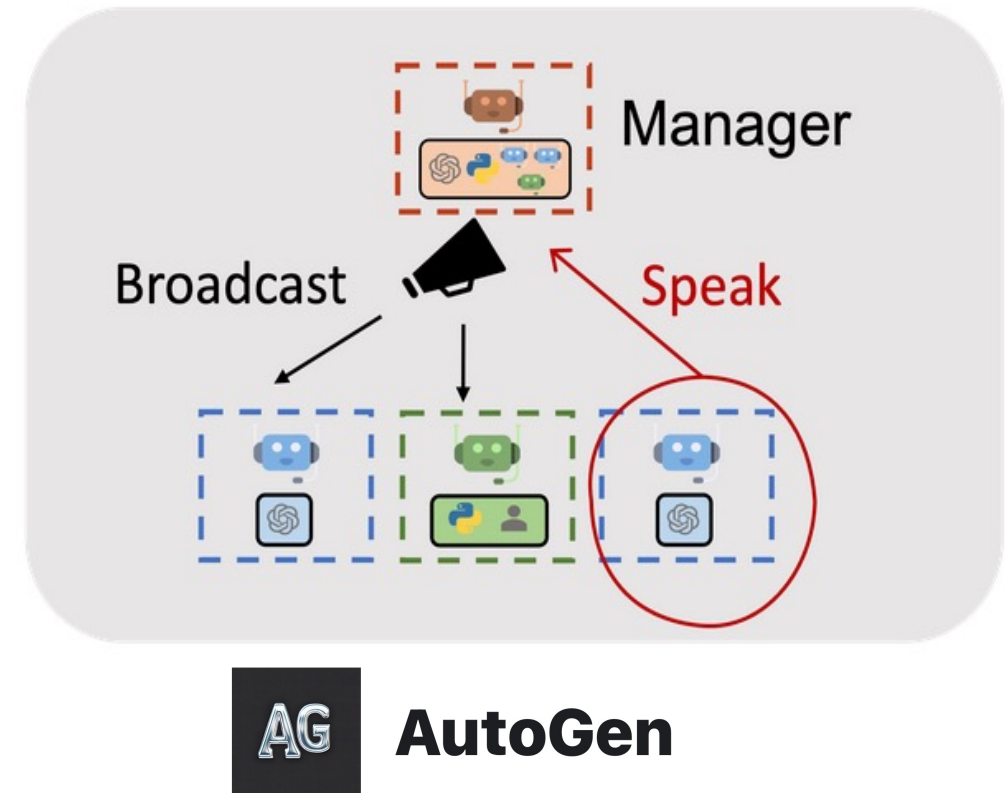- Generel technical complexity increases

# AutoGen uses an open-discussion-style communication

**AutoGen in a nutshell:**
- AutoGen was the first project that specificly targeted cooperation among multiple agents to solve tasks
- Message flow looks like an open discussion between agents
- Has dedicated agents for user, that is also responsible for tool execution

**Learnings:**
- Open message flow limits system relieabilty
- No support for parallel execution of agents
- Customization requires deep interventions in code



Image taken from: https://github.com/langchain-ai/langgraph/blob/main/examples/visualization.ipynb

# LangGraph uses a graph to design the message flow

**LangGraph in a nutshell:**
- Released with LangChain v0.1
- Graph design inspired by networkx
- Provides a nice way to realize very complex message flows in a deterministic way
- LangChain ecosystem brings great tool support
- Much better then early LangChain versions, but still not always intuitive

**Learnings:**
- Graph-Design suits much better to our use-case
- LangChain ecosystem seams more production-ready

Image taken from: https://github.com/langchain-ai/langgraph/blob/main/examples/visualization.ipynb

# Each division has so much individual knowledge!
# But it's sometimes hard to make sense of it

What music do we have the rights to?

What video material do we have about the elections?

**Content Intelligence**

What books do we have about football?

Who is watching dating shows on RTL+?

**Audience Intelligence**

What age is the average viewer of a show?

**Trend Intelligence**

BeTrend

Which international shows are receiving the most attention?

What topics dominated the news recently?

# How can we query this data?

What music do we have the rights to?

What video material do we have about the elections?

What books do we have about football?

Which international shows are receiving the most attention?

Who is watching dating shows on RTL+?

What products do people buy?

What topics dominated the news recently?

Just use a public LLM

Fine-Tune your own LLM

RAG System

At this, Winky howled even harder, her squashed-tomato of a nose dribbling all down her front, though she made no effort to stem the flow.

"Dobby has traveled the country for two whole years, sir, trying to find work!" Dobby squeaked. "But Dobby hasn't found work, sir, because Dobby wants paying now!"

The house-elves all around the kitchen, who had been listening and watching with interest, all looked away at these words, as though Dobby had said something rude and embarrassing. Hermione, however, said, "Good for you, Dobby!"

"Thank you, miss!" said Dobby, grinning toothily at her. "But most wizards doesn't want a house-elf who wants paying, miss. 'That's not the point of a house-elf,' they says, and they slammed the door in Dobby's face! Dobby likes work, but he wants to wear clothes and he wants to be paid, Harry Potter. . . . Dobby likes being free!"

The Hogwarts house-elves had now started edging away from Dobby, as though he were carrying something contagious. Winky, however, remained where she was, though there was a definite increase in the volume of her crying.

"And then, Harry Potter, Dobby goes to visit Winky, and finds out Winky has been freed too, sir!" said Dobby delightedly.

At this, Winky flung herself forward off her stool and lay face-down on the flagged stone floor, beating her tiny fists upon it and positively screaming with misery. Hermione hastily dropped down to her knees beside her and tried to comfort her, but nothing she said made the slightest difference. Dobby continued with his story, shouting shrilly over Winky's screeches.

"And then Dobby had the idea, Harry Potter, sir! 'Why doesn't

Dobby and Winky find work together?' says Winky. 'Where is there enough work for two house-elves?' says Winky. And Dobby thinks, and it comes to him, sir! Hogwarts! So Dobby and Winky came to see Professor Dumbledore, sir, and Professor Dumbledore took us on!"

Dobby beamed very brightly, and happy tears welled in his eyes again.

"And Professor Dumbledore says he will pay Dobby, sir, if Dobby wants paying! And so Dobby is a free elf, sir, and Dobby gets a Galleon a week and one day off a month!"

"That's not very much!" Hermione shouted indignantly from the floor, over Winky's continued screaming and fist-beating.

"Professor Dumbledore offered Dobby ten Galleons a week, and weekends off," said Dobby, suddenly giving a little shiver, as though the prospect of so much leisure and riches were frightening, "but Dobby beat him down, miss. . . . Dobby likes freedom, miss, but he isn't wanting too much, miss, he likes work better."

"And how much is Professor Dumbledore paying you, Winky?" Hermione asked kindly.

If she had thought this would cheer up Winky, she was wildly mistaken. Winky did stop crying, but when she sat up she was glaring at Hermione through her massive brown eyes, her whole face sopping wet and suddenly furious.

"Winky is a disgraced elf, but Winky is not yet getting paid!" she squeaked. "Winky is not sunk so low as that! Winky is properly ashamed of being freed!"

"Ashamed?" said Hermione blankly. "But — Winky, come on! It's Mr. Crouch who should be ashamed, not you! You didn't do anything wrong, he was really horrible to you —"

# Directions

1st Floor: **BMG**
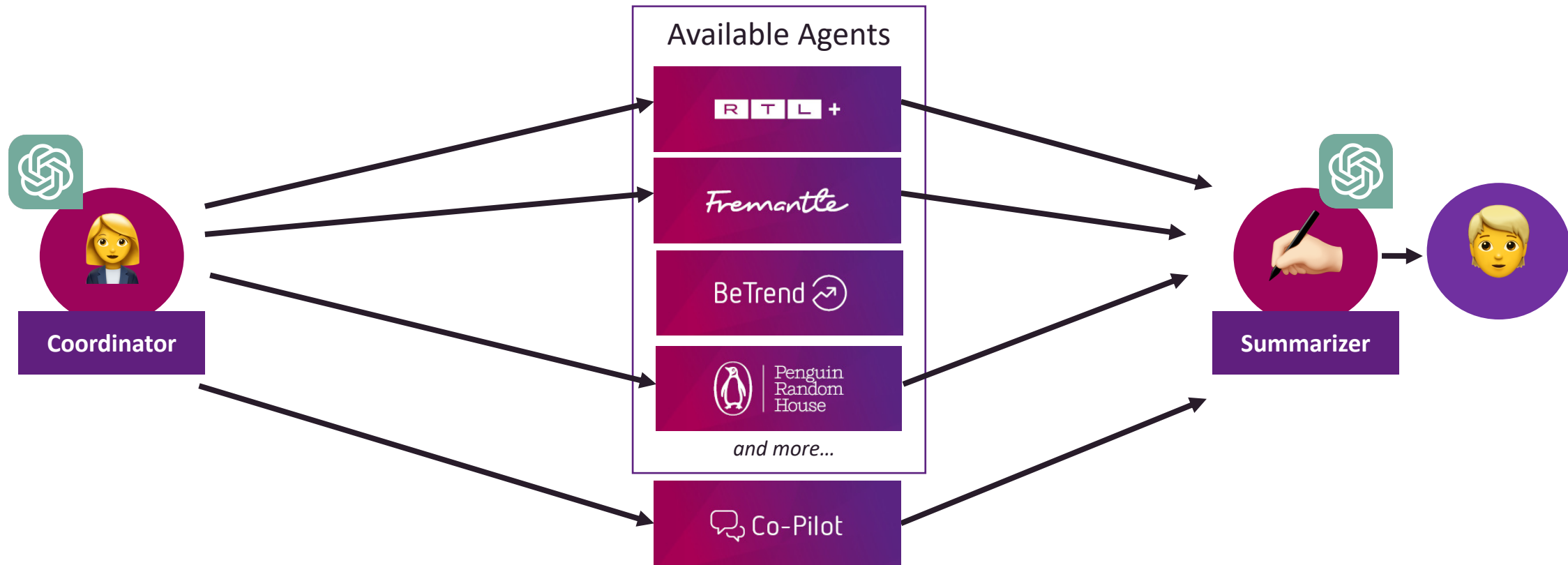
3rd Floor: Penguin Random House

2nd Floor: RTL

4th Floor: BeTrend

# Different agent-categories are used to access different types of data

**Agents using vector search:**
Similar to classical RAG, but with the possibility to add filters and to use tools multiple times

Tools: Qdrant

**API-based Agents**
Using tools that wrap the relevant parts of the specific APIs

*Only custom Tools*

**Graph-Based Agents:**
Executes dynamically filled Cypher-Template-Querries

Tools: neo4j WIKIPEDIA

**Other Agents:**
Perform web-search and add relevancy information for the selected agents

Tools: Tavily

# Based on the user request the relevant agents get triggered to combine them in a single answer

# It's demo time!

**Content Search** beta

Ask me about content at Bertelsmann...

| Welche PRH Autoren kommen aus Berlin? | What TV Series are there about climate change? |
| Recommend me a good action movie starring Denzel ... | What were the highest trending topics in the US in the l... |

Currently we offer access to the following sources:

IPRNOW    Penguin Random House    BeTrend    RTL+    Penguin Random House Verlagsgruppe    BENET

More to come soon...

# The LangChain Ecosystem Offers Much

**Most Popular High-Level LLM Building Framework**
⭐ 92k Github, Python + JS,

**Getting (a bit) more production ready**
- (Mostly) Backwards compatible with major versions now
- v0.3 now with Pydantic 2 support
- $10M Seed funding

**Chains as Core Concept**
- Chain = Input + Prompt + LLM
- Plus, optional components

LangChain

# LangChain Overview – the core package to create LLM chains

BERTELSMANN

scieneers
DRIVEN BY DATA



| LangServe | Chains as Rest APIs |

| Templates | Reference Applications |

## LangChain
Python · JavaScript

**Chains, Agents, Advanced Retrieval Strategies**

## LangChain-Community
Python · JavaScript

| Models I/O | Retrieval | Agent Tooling |
|---|---|---|
| Model<br>Prompt<br>Example Selector<br>Output Parser | Retriever<br>Document Loader<br>Vector Store<br>Text Splitter<br>Embedding Model | Tools<br>Toolkits |

## LangChain-Core
Python · JavaScript

**LCEL - LangChain Expression Language**
Parallelization · Fallbacks · Tracing · Batching · Streaming · Async · Composition

Image taken from: https://python.langchain.com/v0.1/docs/get_started/introduction

# The LangChain Ecosystem goes beyond the LangChain package

## LangChain

Build individual components
by combining LLM, storage
and prompts

## LangGraph

Orchestrate complex
flows through one or
more  components

## LangServe

Wrap LangChain chains
available as FastAPI
endpoints

## Langfuse

Monitor and
continuously evaluate
applications

# LangFuse is the LLMOps jack-of-all-traits –
# Tracking and Tracing

# LangFuse is the LLMOps jack-of-all-traits – Evaluation (Scoring)



Manual annotation

All Scores (API or LangFuse UI)

# LangFuse is the LLMOps jack-of-all-traits – Evaluation (Annotation & datasets)



User input

Desired system output

# Azure OpenAI Latency – What we expected

- LLM-Calls are send to the Azure Open-AI Endpoint, which routes those to the next free GPU.

- Execution time depends on the Input, the Model and the regional system load

- Region of the Azure OAI Endpoint matters, but not the subscription or model

**LLM-Calls**

**Azure OAI-Endpoint**

**GPUs**

# Azure OpenAI Latency – What we got

- LLM-Calls are send to the Azure Open-AI Endpoint, that routes all to the same GPU

- The GPU does not handle requests in FIFO-order but splits its token generation capacity to all requests equally

- Execution time scales linearly with the number of concurrent requests

# Azure OpenAI Latency – How we dealt with it

- We deployed dedicated Azure Open-AI Endpoints per Agent (all serving the same model). Those endpoints route to different GPUs.

- Different agents are don't affect the execution time of others.

- Multiple requests to the same agent still scale linearly.

# Fast technology exploration and productive software development require different skillsets and workflows

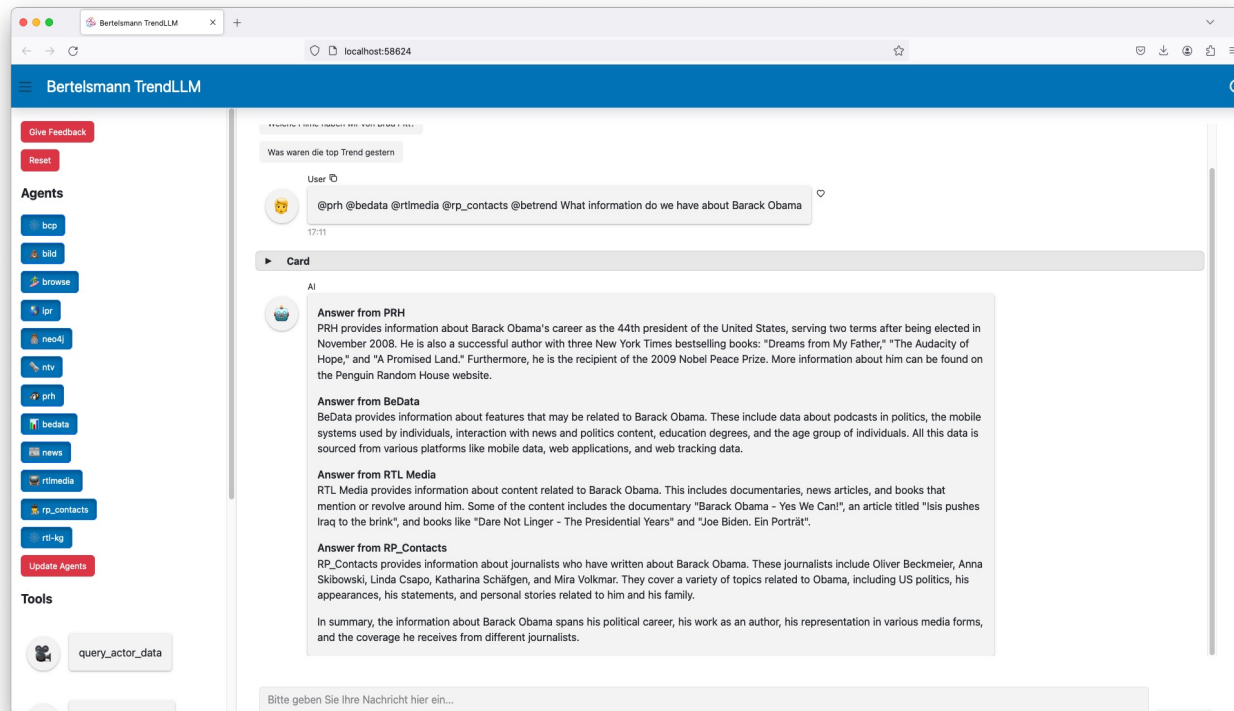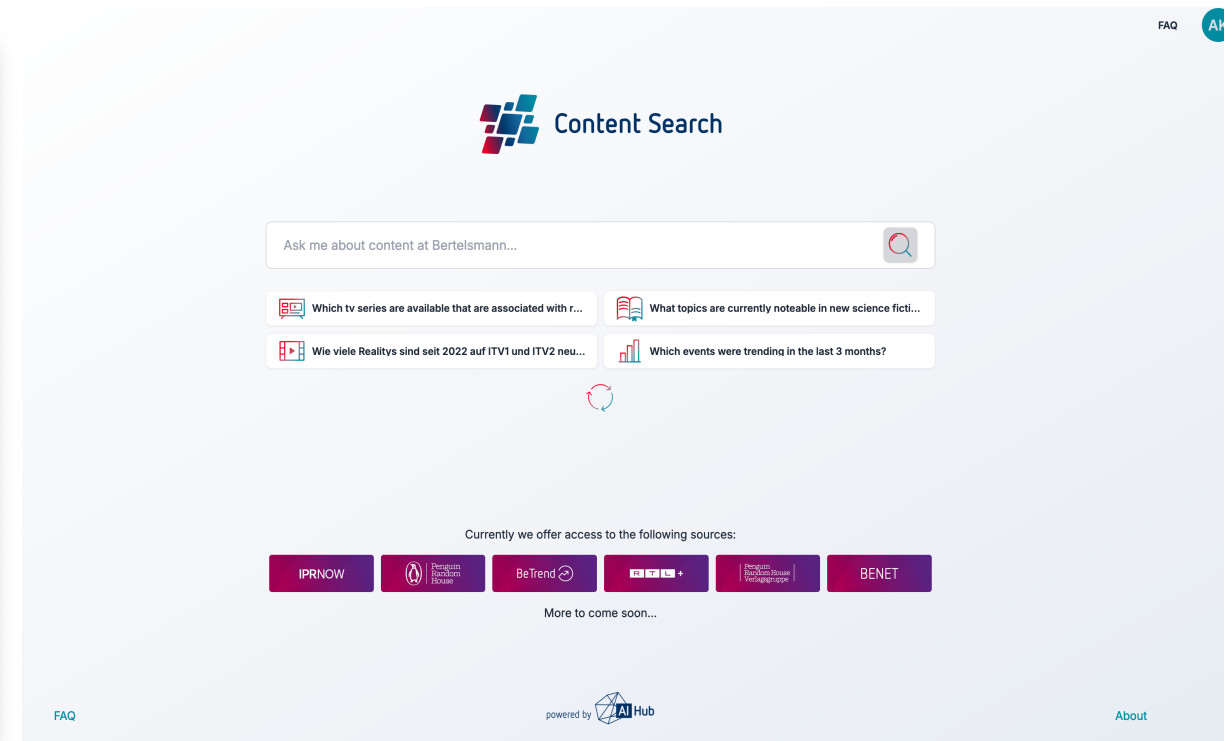|  | **Technology Exploration** | **Software Development** |
|---|---|---|
| Goal | Quick-Insights | Reliable Software |
| Team | A few creative allrounders | Interdiscipilnary team including frontend & infrastructure |
| Collaboration | Quick-Direct Communication, Great personal Freedom to test and play | Kanban/ Scrum with Dailies Clear responsibilities |

# Show, don't tell! Building prototypes is helpful...



**UI Explorative Phase**:
Focus on Usability (Panel / Python)



**Final UI**:
Focus on User-Experience (React)

# Also the technical components required grow in complexity

**Technology Exploration**

A LLM app that handles most cases

Single dump of data

A simple frontend

Deployment on localhost

**Software Development**

A LLM app that handles all cases

Scheduled ETL pipelines incl. update logic

A sophisticated user interface

Deployment in cloud environment

IT security

Load balancing and scaling

Logging, monitoring & alerting
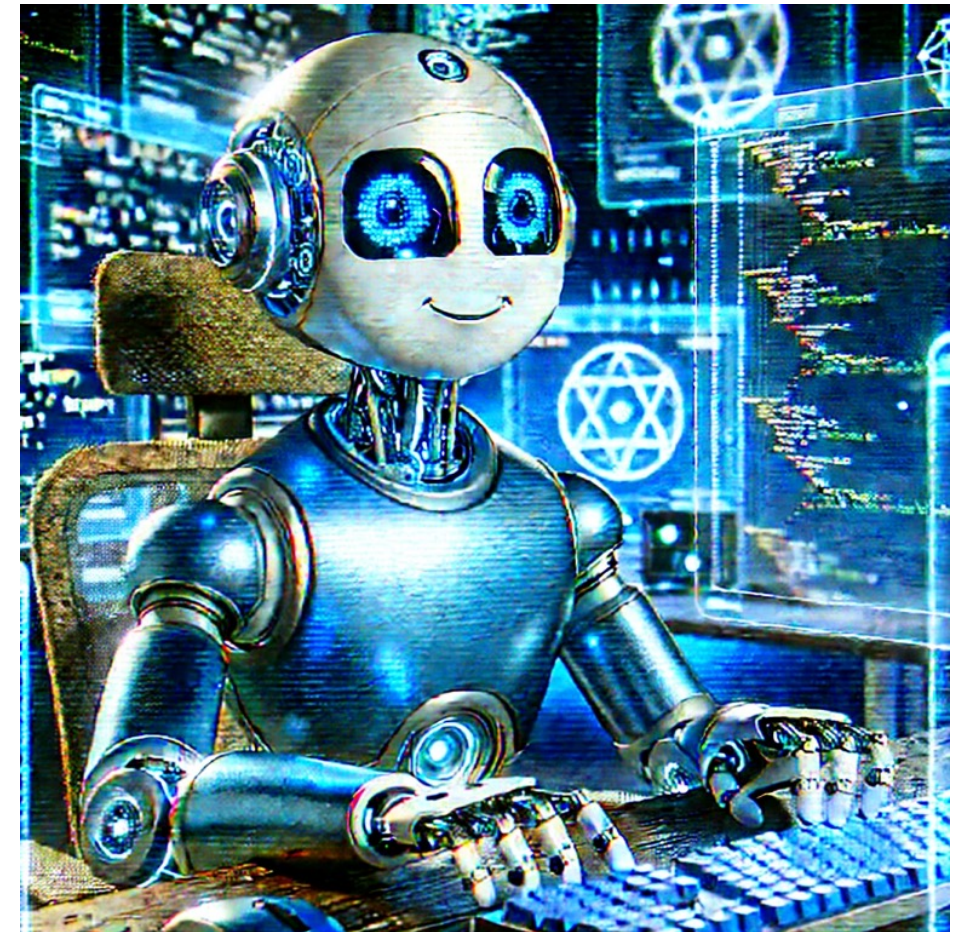
Identity and access management (IAM)

# Best-Practices from classical software development stay relevant in the LLM-Era

Transforming PoCs into **stable** and **reliable software is costly**, due to security, stability and UX/UI claims.

**Tracing** and **Monitoring** are crucial, to understand the system and its limitations

Multi-Agent-Systems are the **Microservice-Design** of LLMs. It increases maintainability, makes side-effects controllable and overall predictability

**Unix-Rules** apply: Make each program do one thing. Tools and LLMs shall have a clear purpose, don't expect magic.

# Feedback und Evaluation

**Evaluation is hard and requires expert knowledge**
- Working with users and data source owners is important
- Confirmation bias stemming from confident LLM responses

**Integrating Feedback buttons early helps**
- With gaining an intuition for the system performance
- Curating test sets for automated test & evaluation runs

**UX of the feedback button matters**
- Position the button centrally, in case of doubt add a second one
- Carefully consider lickart scale ratings vs. binary feedback

**Logging user traces is valuable but can be legally challenging**
- Valuable as it allows to adapt system for user requirements
- Remember it has GDPR implications, pseudonymisation helps

# Next Steps

**Improve quality and range of feedback data**

- Automate scoring of traces to fill feedback gaps
- Gain holistic impression of system performance

**Balance System improvements vs. extensions**

- Onboard new data sources
- Improve existing Agents e.g. through advanced RAG design patterns

# Q&A

**Moritz Glauner**

Head of Data Science @ Bertelsmann

moritzglauner

**Nico Kreiling**

Data Scientist @ scieneers

Host of techtiefen.de      nico-kreiling