SHAPE
A MYTY AGENCY

# AI-Agents

…under the hood

# What's an LLM?

> **G** Pretend you're a medieval scribe accidentally given a modern tech support ticket.

*Adjusts medieval spectacles and unfurls modern support ticket on parchment*

Dearest troubled soul,

I, Brother Bartholomew of the Sacred Help Desk, have received thy electronic plea via messenger pigeon (or what thou callest "email"). Thou speaketh of an accursed "blue screen" that plagues thy computing device?

Verily, I shall guide thee through the sacred rituals of troubleshooting:

1. First, thou must perform the ancient rite of "turning it off and on again"

2. If the demons persist, we shall need to check thy "device manager" - though I confess, I know not if this requires eye of newt

Prithee respond with haste if these holy incantations do not banish thy technical woes.
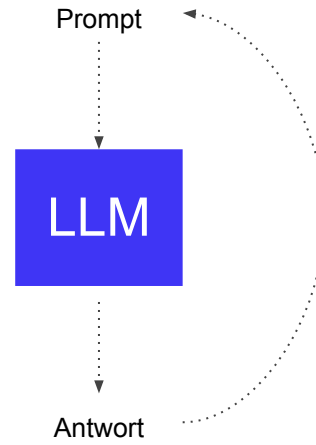
May thy bandwidth be plentiful,
Brother Bartholomew
Tech Support Monastery, Third Cubicle from the Left

*Adds illuminated border of frustrated users and tangled cables to ticket*

Copy  Retry

Prompt

LLM

Antwort

*LLM (e.g. ChatGPT)*

# What if the prompt requires you to…

… create a summary of multiple documents?

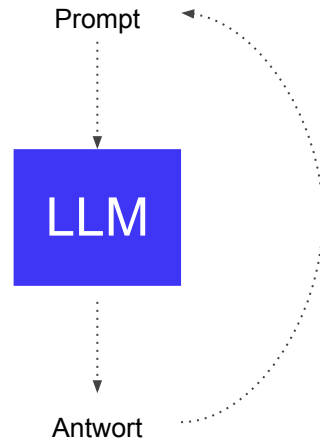*"Provide a summary of our 2023 earnings"*

… query a database?
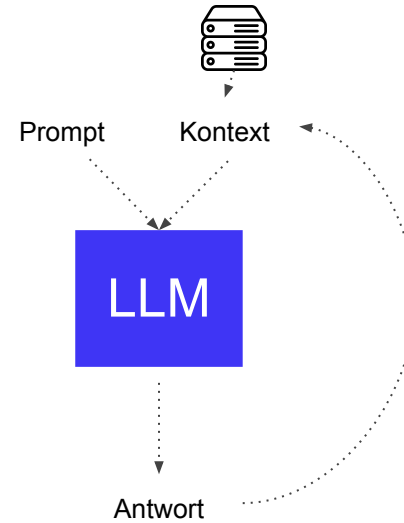
*"What were our top selling product categories in Q2?"*

… break down the questions into multiple subquestions?
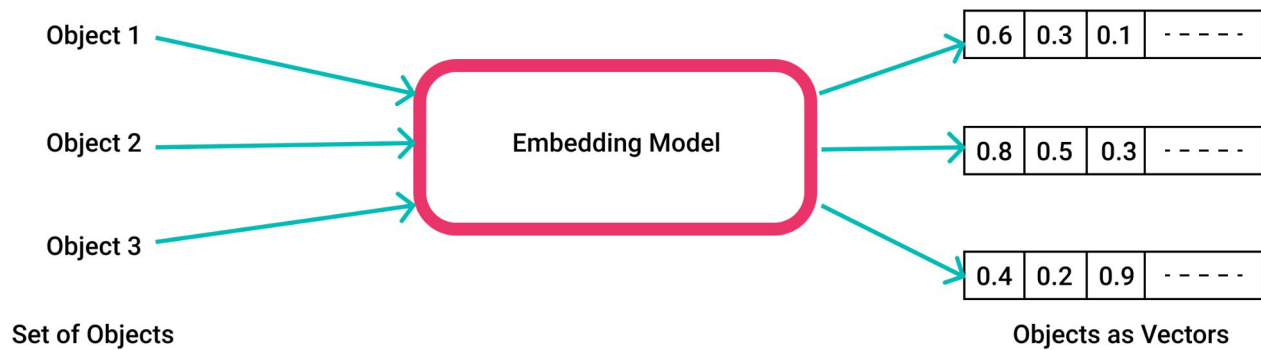
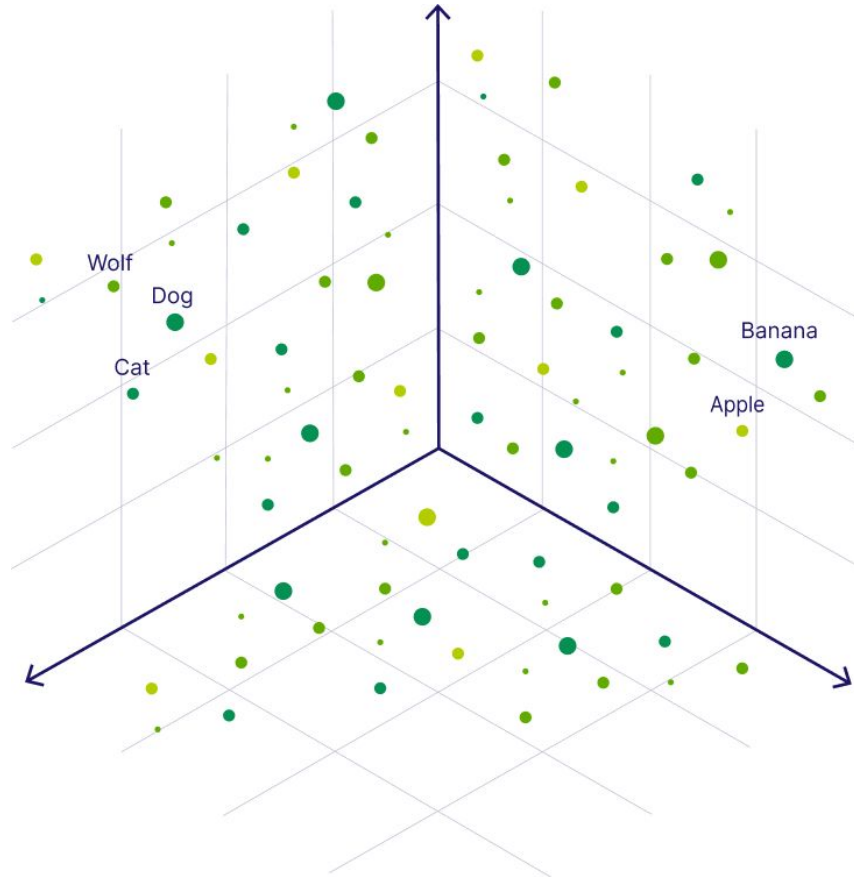*"Compare our revenue growth between Q2 2023 and Q2 2024"*

# RAG



Prompt

Kontext

LLM

Antwort

Prompt

LLM

Antwort

*LLM (z.B. ChatGPT)*

*Retrieval Augmented Generation (z.B. Company Chatbot)*

Excursion: Embedding Models



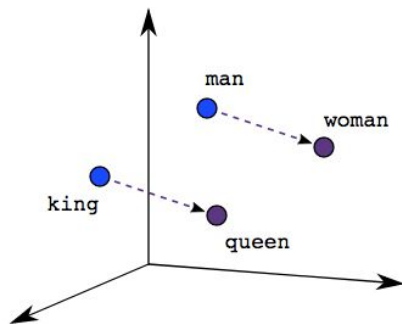| Object 1 | | Embedding Model | | 0.6 | 0.3 | 0.1 | - - - - - |
| Object 2 | | | | 0.8 | 0.5 | 0.3 | - - - - - |
| Object 3 | | | | 0.4 | 0.2 | 0.9 | - - - - - |

**Set of Objects**

**Objects as Vectors**

Embedding models translate text into meaning

Vector-space

Semantically similar terms are close to each other



Male - Female

Verb Tense

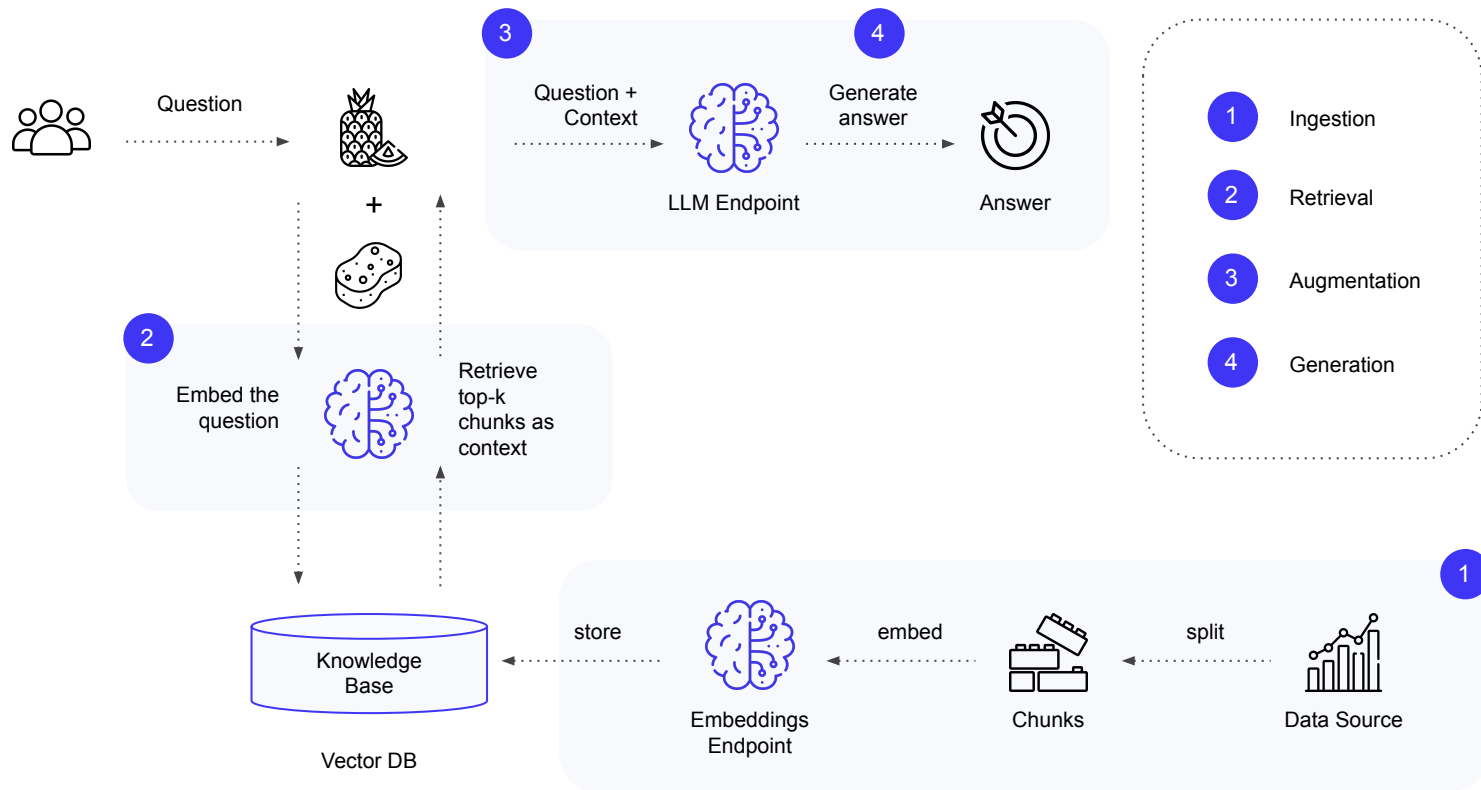Country - Capital

# Retrieval Augmented Generation (RAG) concept

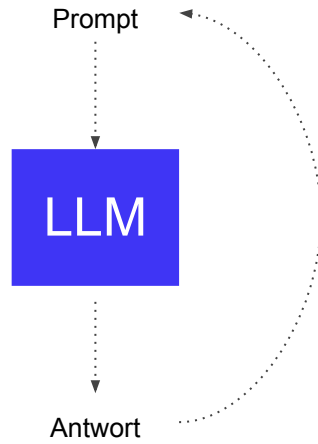# But what if the prompt requires you to…

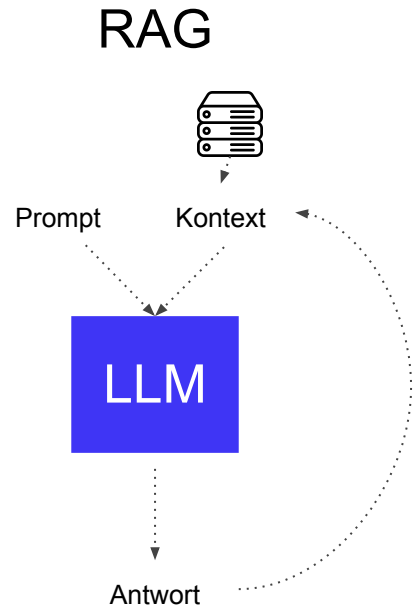… modify external data?

*"Create a JIRA ticket"*

… research the web?

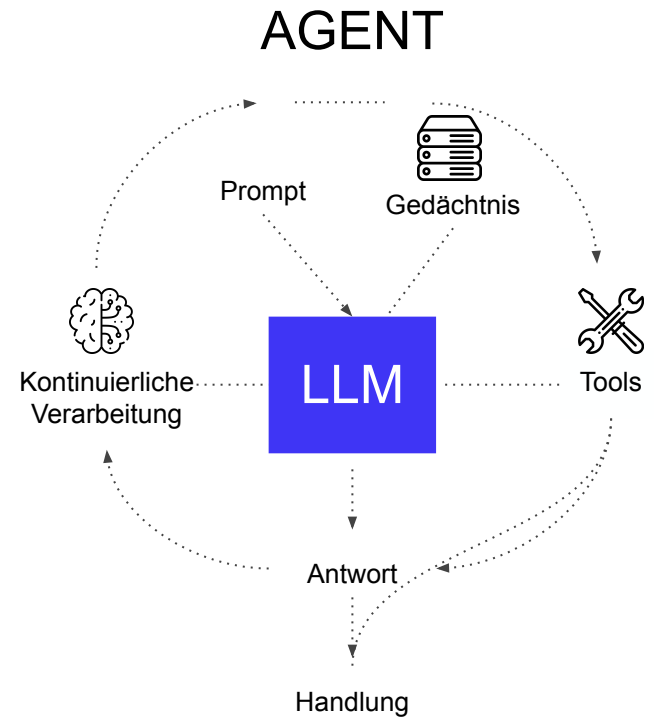*"How's the weather in Zagreb?"*

… reason about complex tasks?

*"Find the best way to implement feature X into our product"*

# RAG

# AGENT

**LLM**

**LLM**

**LLM**

Prompt

Prompt    Kontext

Prompt    Gedächtnis

Kontinuierliche
Verarbeitung

Tools

Antwort

Antwort

Antwort

Handlung

*LLM (z.B. ChatGPT)*

*Retrieval Augmented Generation
(z.B. Company Chatbot)*

*Neue Ära: autonome Systeme*

# Agenticness: A gradual feature



Agenticness

Goal Complexity                                          How difficult and flexible are the goals?

Environmental Complexity                                 How dynamic is the environment?

Adaptability                                             How well does the system adapt to unforeseen situations?

Autonomy                                                 How independent does the system act?
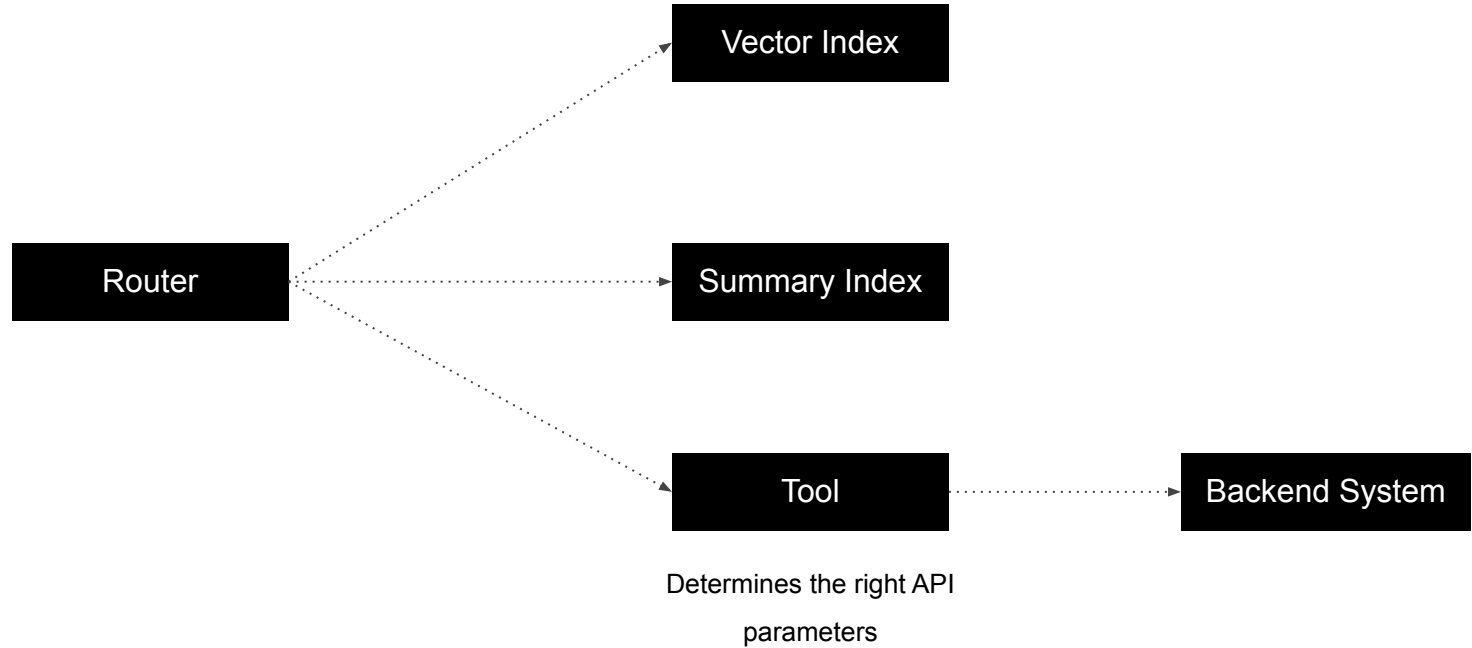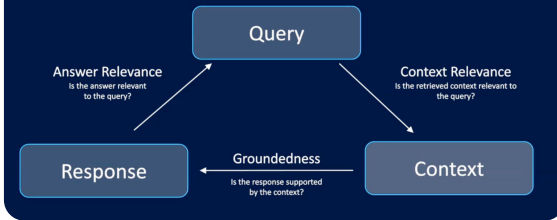
# Adding agentic Layers to RAG

**Routing**

**Tool Use**

**Query Planning**

**ReAct**

Simple

Lower Cost

Lower Latency

Advanced

Higher Cost

Higher Latency

# Stage 1 - Router

# Stage 2 - Query planning

Vector Index

Query Planner

Decides how to break down the

question and picks the right index

Summary Index

# Stage 3 - Tool Use

Vector Index

Router

Summary Index

Tool

Backend System

Determines the right API

parameters

# Stage 3 - Tool Use

"What is your biggest limitation of putting more agents in production?"

Performance                                                                41%

Cost                                                                     18.4%

Safety concerns                                                          18.4%

Latency                                                                  15.1%

Other                                                                      7%

# The more "agentic" the less reliable…



Reliability | Agenticness

| Goal Complexity | How difficult and flexible are the goals? |
|---|---|

| Environmental Complexity | How dynamic is the environment? |
|---|---|

| Adaptability | How well does the system adapt to unforeseen situations? |
|---|---|

| Autonomy | How independent does the system act? |
|---|---|

# Solutions for better reliability



Evaluation tests



Human in the loop



LLM as a Judge

# Use-Case: Healthcare (Falk AI)

# Use-Case: Logistics Provider (HHLA)

# In a Nutshell
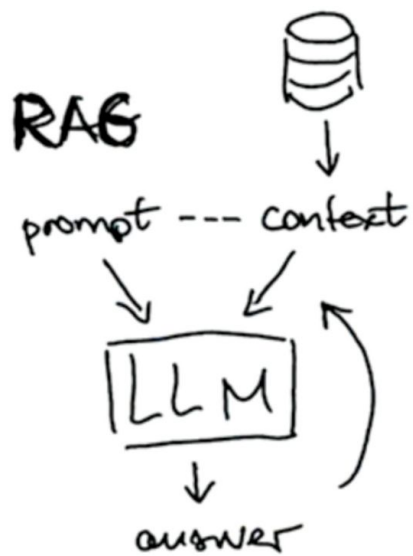
Agentic systems are AI-Systems, that can accomplish…

…complex goals
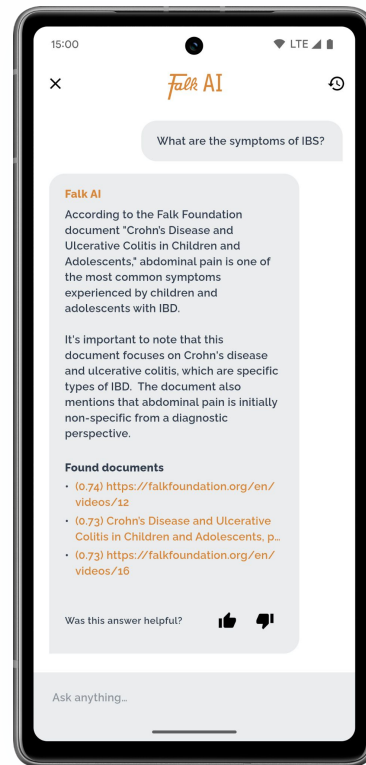
…in complex environments

…with limited direct supervision

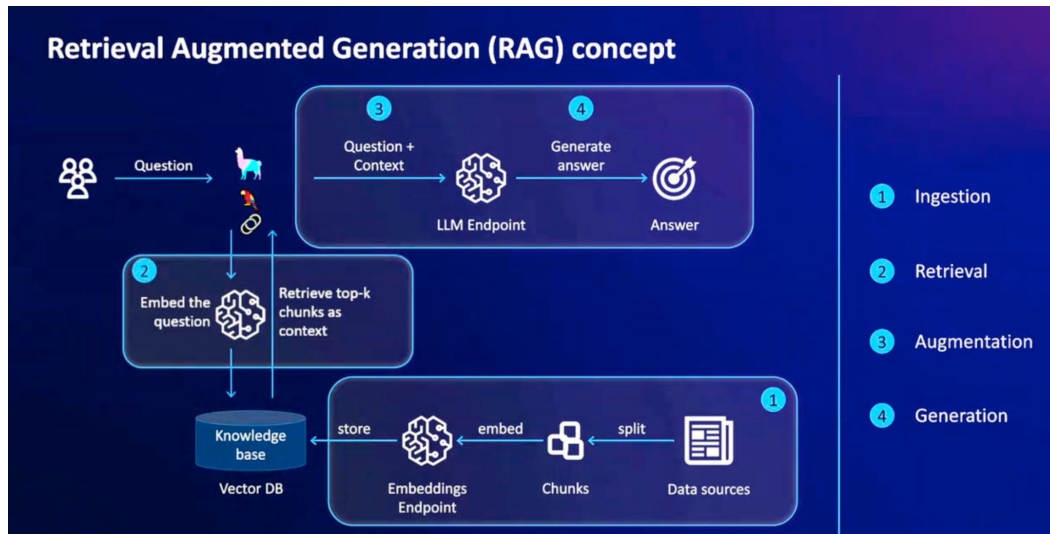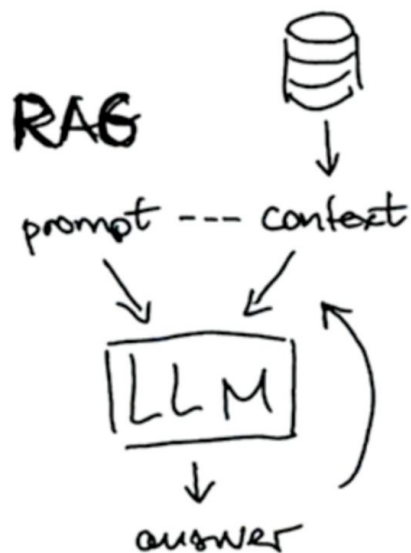# Backup

prompt

↓

LLM

↓

answer

prompt

LLM

answer

RAG

prompt --- context

LLM

answer

# Retrieval Augmented Generation (RAG) concept

RAG

AGENT

prompt

↓

**LLM**

↓

answer

prompt --- context

↓ ↘

LLM

↓

answer

prompt

↓

LLM

↓

answer

↓

action

reasoning

memory

tools

# Adding agentic layers to RAG

Routing

Query Planning

Tool Use

ReAct

Simple
Lower Cost
Lower Latency

Advanced
Higher Cost
Higher Latency

Step 2: Query Planning

# Stage 3 – Tool use



Router

Vector index

Summary index

tool

backend system

Determines the right API parameters

Stage 4: ReAct

"What is your biggest limitation of putting more agents in production?"

Performance quality                                          41%

Cost                                                        18.4%

Safety concerns                                             18.4%

Latency                                                     15.1%

Other                                                        7%

# The more "agentic", the less reliable…



Reliability | Agenticness

🎯 Goal Complexity
How difficult and flexible are the goals?

🌎 Environmental Complexity
How dynamic is the environment?

How well does the system adapt to unforeseen situations?

🔁 Adaptability

How independent does the system act?

👁 Autonomy

# Solutions for better reliability



**Evaluation tests**



**Human in the loop**



**LLM as a Judge**

# Use-Case: Logistics Provider (HHLA)

# In a Nutshell 🥜

Agentic systems are AI-Systems, that can accomplish

- complex goals 🎯
- in complex environments 🌏
- with limited direct supervision 👀

# Improving LLMs 🤖



Prompt Engineering

RAG

Fine-Tuning

Low ← Effort → High

Pre-Training
(Computationally Expensive)

Fine-Tuning
(Cheaper)

LLM → LLM

Large
Unlabeled Corpus

Small
Labeled Corpus