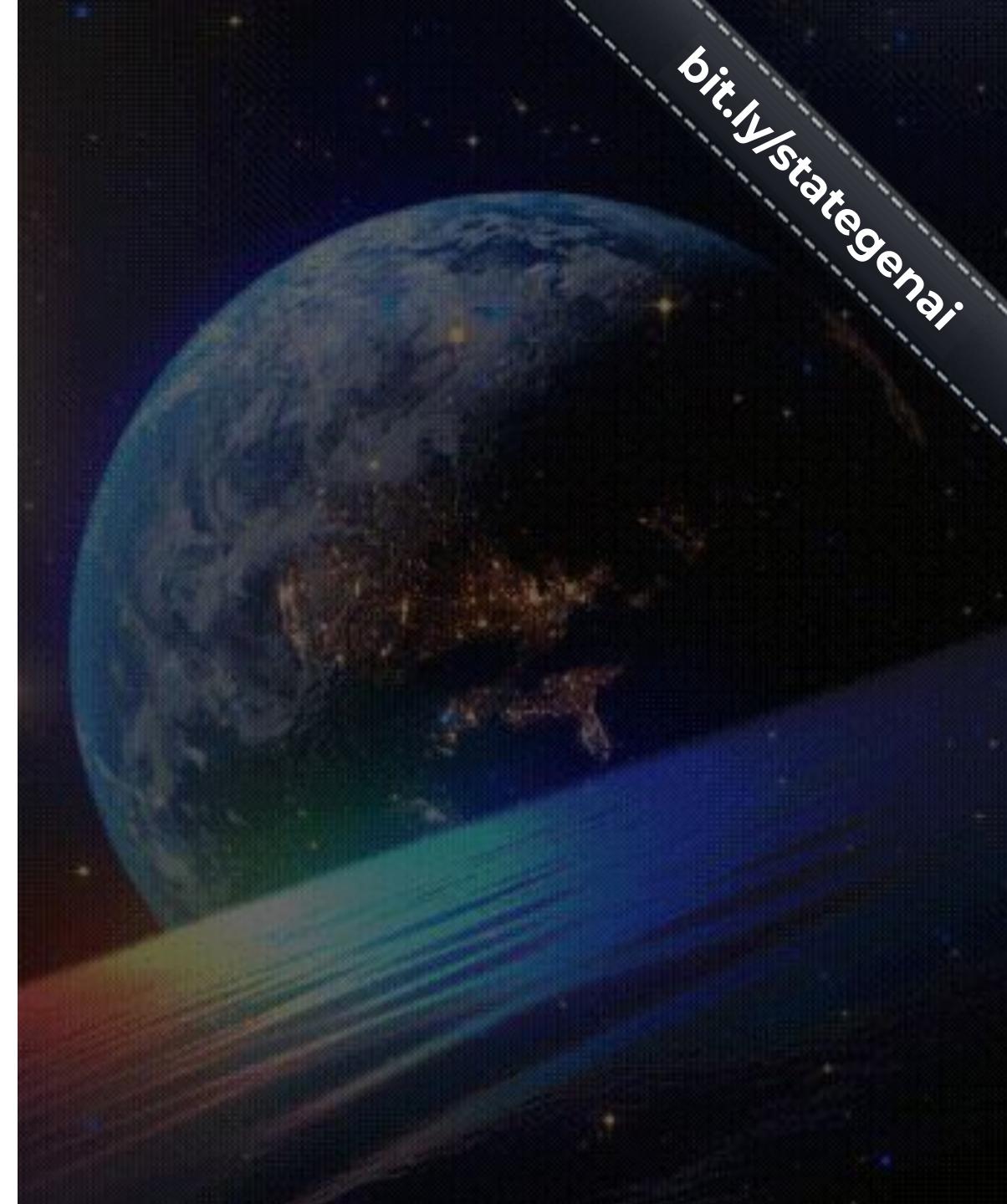
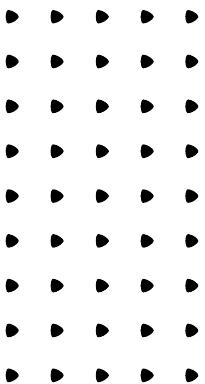


The State of GenAI & ML in the Cloud Native Ecosystem

Alejandro Saucedo
Bartosz Ocytko



zalando



About us

Alejandro Saucedo

Director of Eng, Science & Product



Bartosz Ocytko

Executive Principal Engineer

Zalando | Who We Are

Leading European fashion.

~15bn

Euro GMV

~51m

Active Customers

~2m

Articles





Zalando | GenAI & ML

Across ~80 Applied Science teams



Search & Recommendations	Warehouse Inventory Optimization	Conversational AI GenAI
Replenishment Prediction	Pricing Elasticities	Forecasting Demand, Returns, etc
Size & Fit Recommendations	Machine Translation GenAI	Fraud Detection
Virtual Clothing Try-On GenAI	Marketing Intelligence	Competitive Matching

We're hiring!
Come work with us.

jobs.zalando.com

Find your career

 Search for a role...

All open roles



A photograph of a man and a woman standing in a dense forest. The man is on the left, wearing a dark blue hooded jacket over a light-colored vest and dark pants. The woman is on the right, wearing a light green vest over a white t-shirt and light-colored pants. They are both looking towards the camera. The background is filled with tall, thin trees and dappled sunlight.

The State of GenAI & ML in Cloud Native Ecosystem

Motivations

Market trends,
opportunities and
challenges



GenAI Products

How it started



ChatGPT



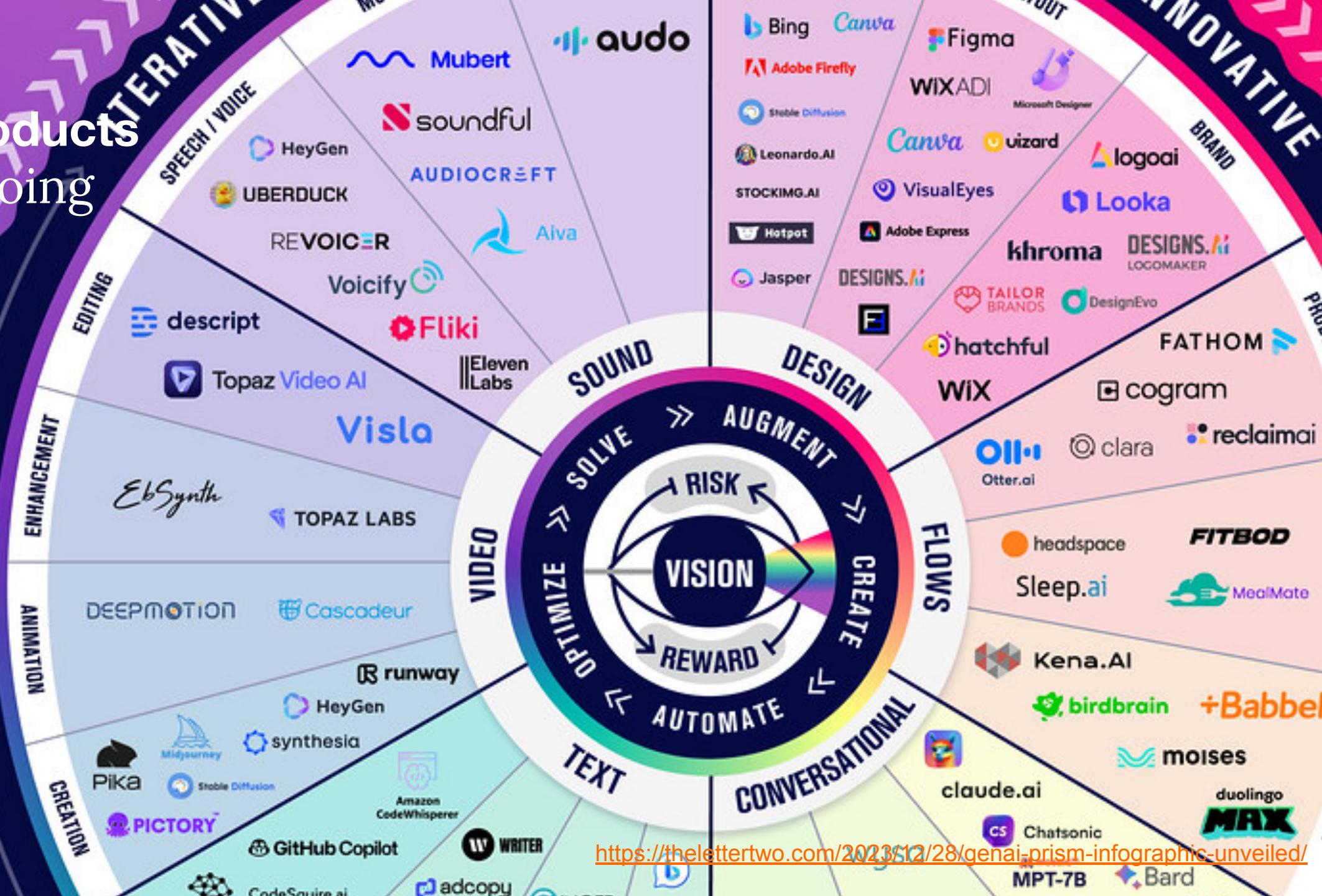
DALL-E



Stable Diffusion

GenAI Products

How it's going



GenAI Products

How it's going



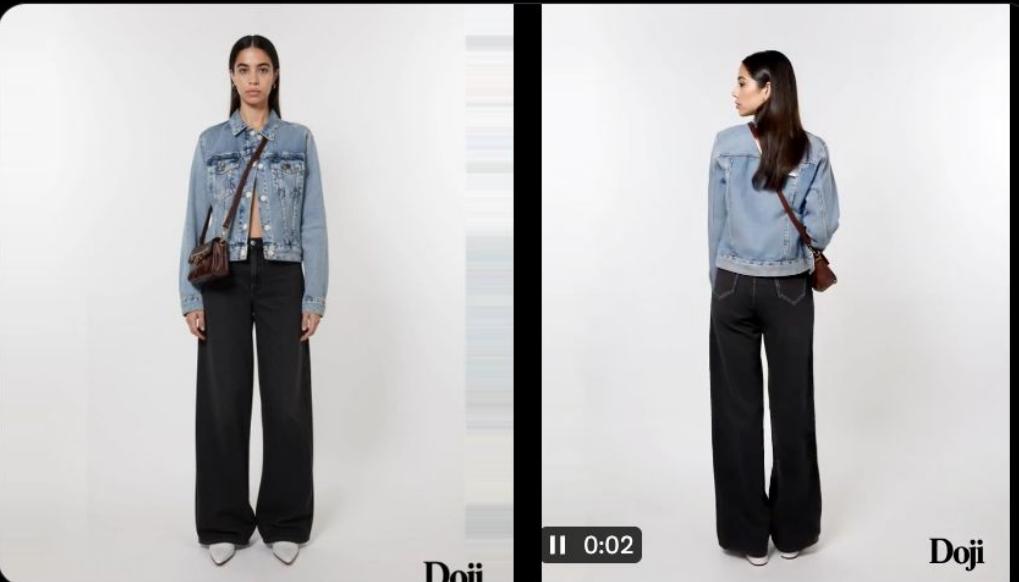
Shweta ✅

@shweta_ai

I turned myself into an AI model! 🤖✨

Everything - the digital me, the outfits, the image-to-video transformation - is 100% AI-generated.

Here's exactly how I did it using Doji, Channel42, Kling AI & fal. 🧵 ↗



5:37 PM · Feb 27, 2025 · 194.2K Views

securityJustInterferesWithVibes

Meme

leo ✅
@leojr94_

my saas was built with Cursor, zero hand written code

AI is no longer just an assistant, it's also the builder

Now, you can continue to whine about it or start building.

P.S. Yes, people pay for it

4:34 am · 15 Mar 2025 · 52.2K Views

leo ✅
@leojr94_

guys, i'm under attack

ever since I started to share how I built my SaaS using Cursor

random thing are happening, maxed out usage on api keys, people bypassing the subscription, creating random shit on db

as you know, I'm not technical so this is taking me longer than usual to figure out

for now, I will stop sharing what I do publicly on X

there are just some weird ppl out there

9:04 am · 17 Mar 2025 · 53.6K Views



Pres Mihaylov ✅

@PreslavMihaylov

A non-engineer vibe-coded a simple book-suggestion app using Cursor & Claude, ended up extremely frustrated past the 85% mark.

"it feels like building an endless Jenga tower that just doesn't get higher"

in 3 to 6 months, AI
ode software
of

Share

Save

The AI Market and it's numbers.

~\$390Bn.

Market valuation [1]

+5x.

Projected 5y increase [1]

+100m

Working in AI space by EOY [2]

[1] <https://explodingtopics.com/blog/ai-statistics>

[2] https://www3.weforum.org/docs/WEF_Future_of_Jobs_2020.pdf





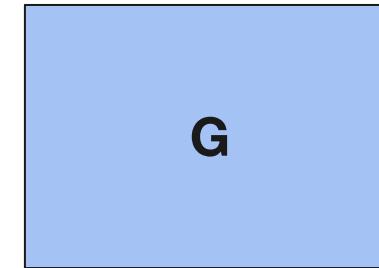
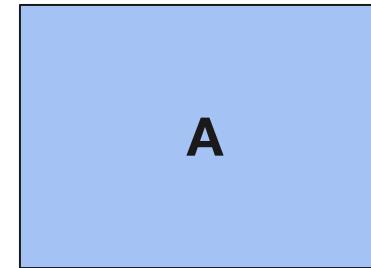
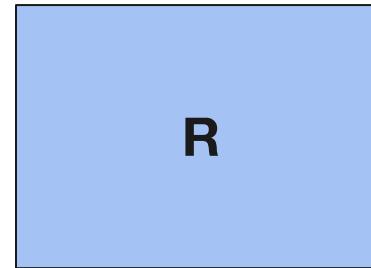
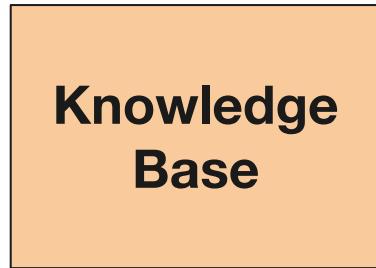
The State of GenAI & ML in Cloud Native Ecosystem

Challenges

GenAI Journey Overview
and Challenges reflected

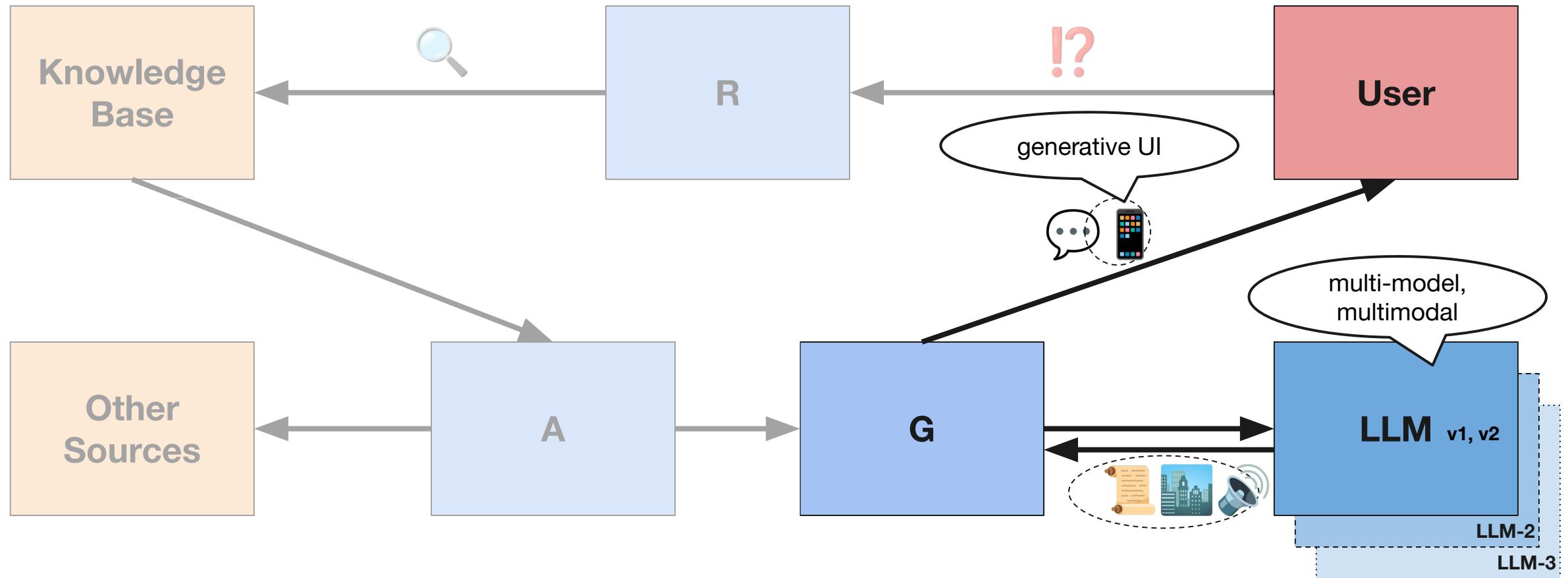
RAG

The “Hello World” for LLMs



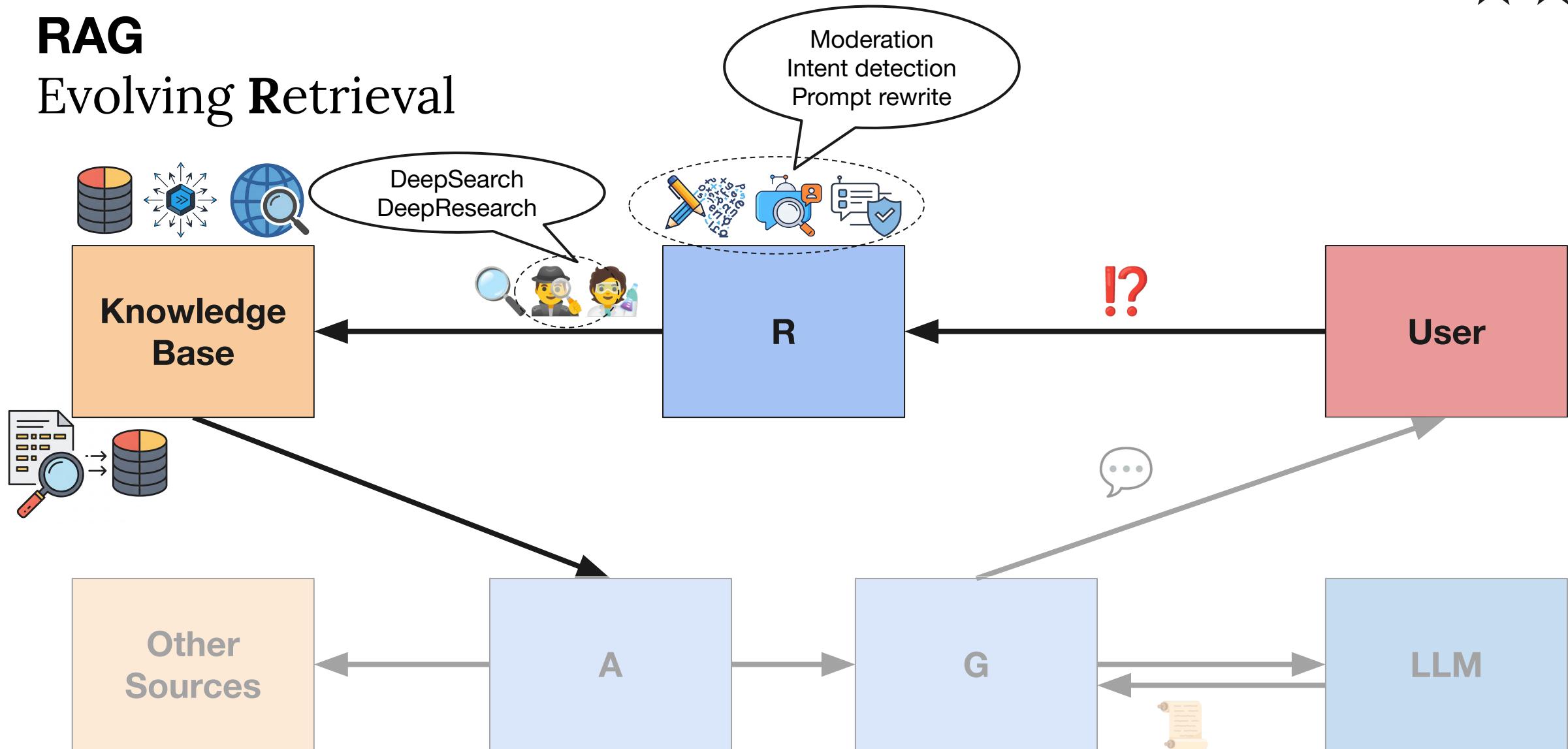
RAG

Evolving Generation



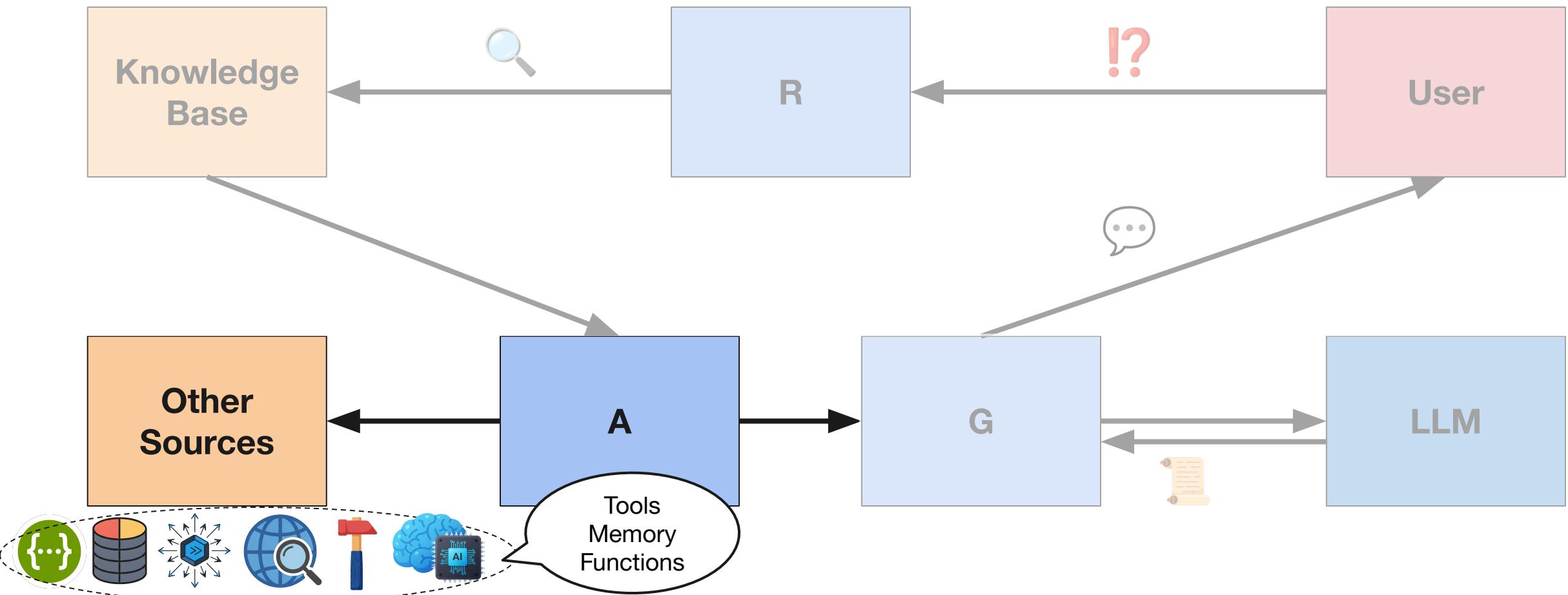
RAG

Evolving Retrieval



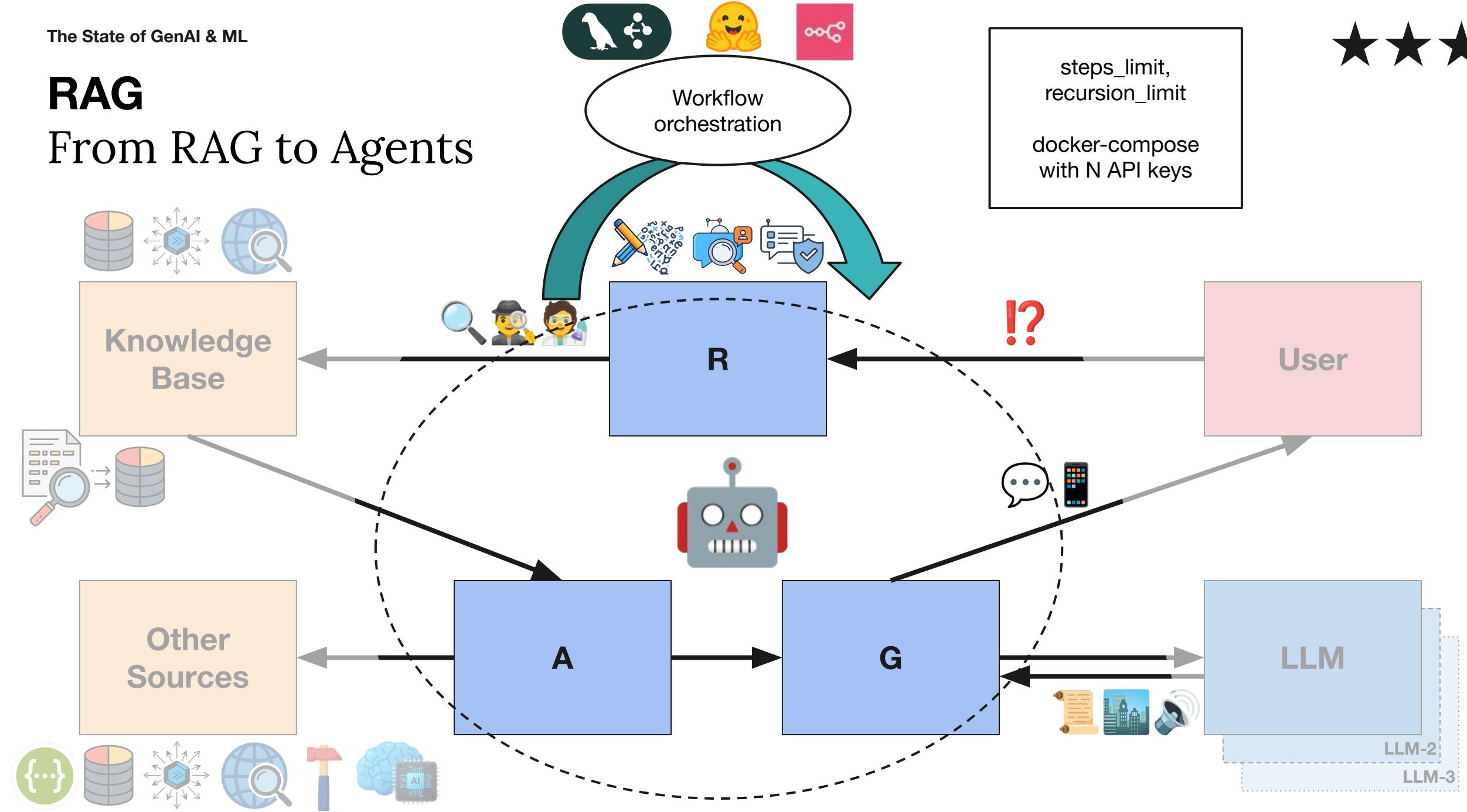
RAG

Evolving Augmentation



RAG

From RAG to Agents

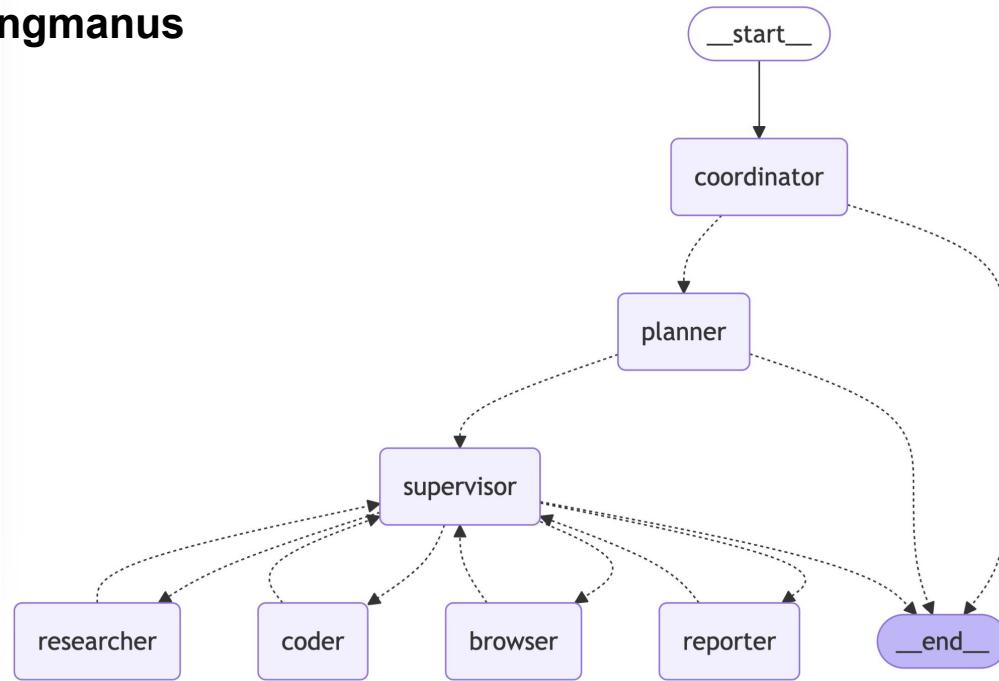




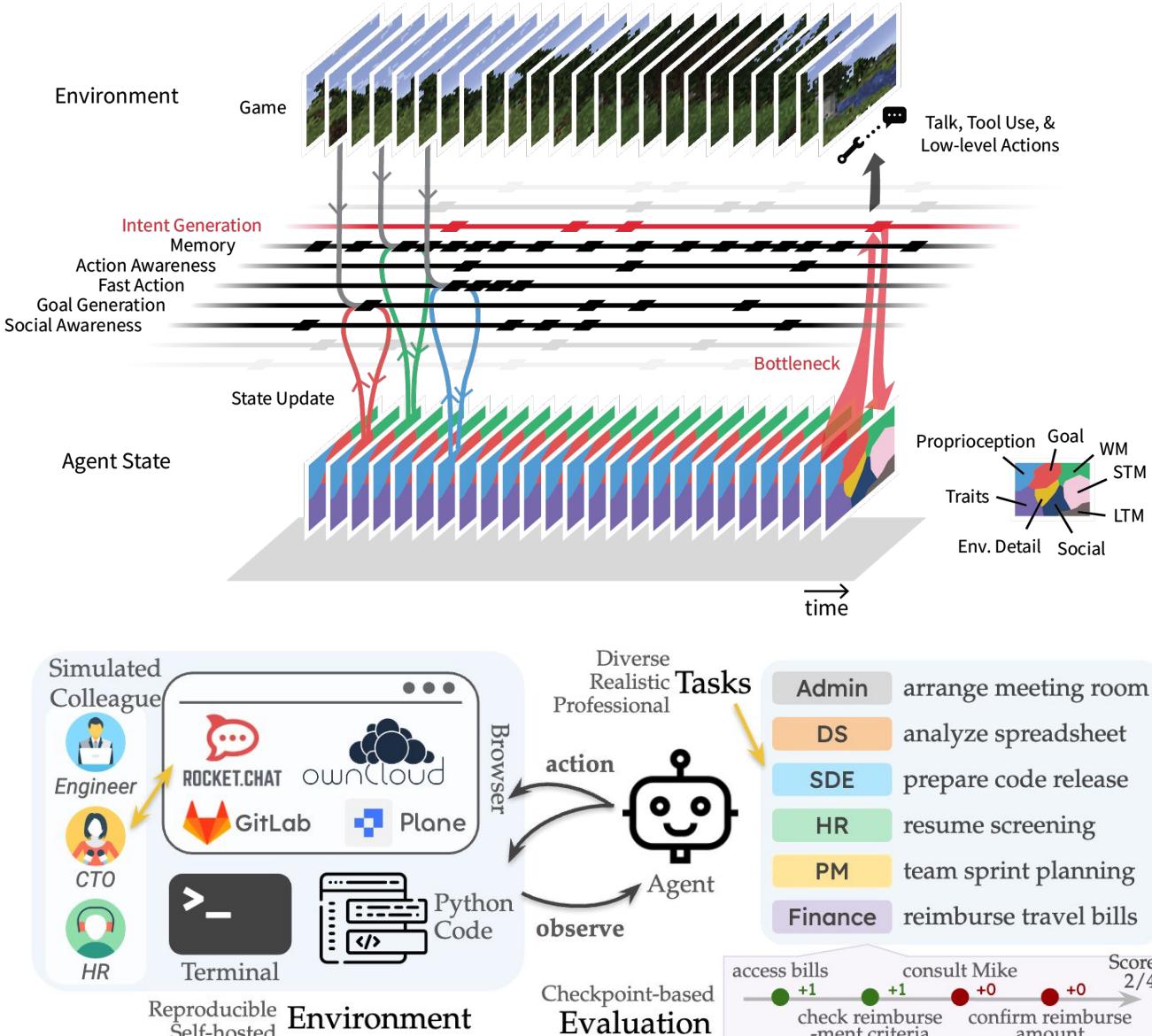
Agents on Steroids

Autonomous Agents

Langmanus



“Let loose in Minecraft, 1,000 autonomous AI agents collaborate to build their own society” (Project Sid from Altera)
<https://digitalhumanity.substack.com/p/project-sid-many-agent-simulations>



Benchmarks:

- **The Agent Company (Dec, 2024)**
- GAIA (Nov, 2023)
- AgentBench (v2, Oct, 2023)
- WebShop (v4, Feb 2023)

<https://the-agent-company.com/>



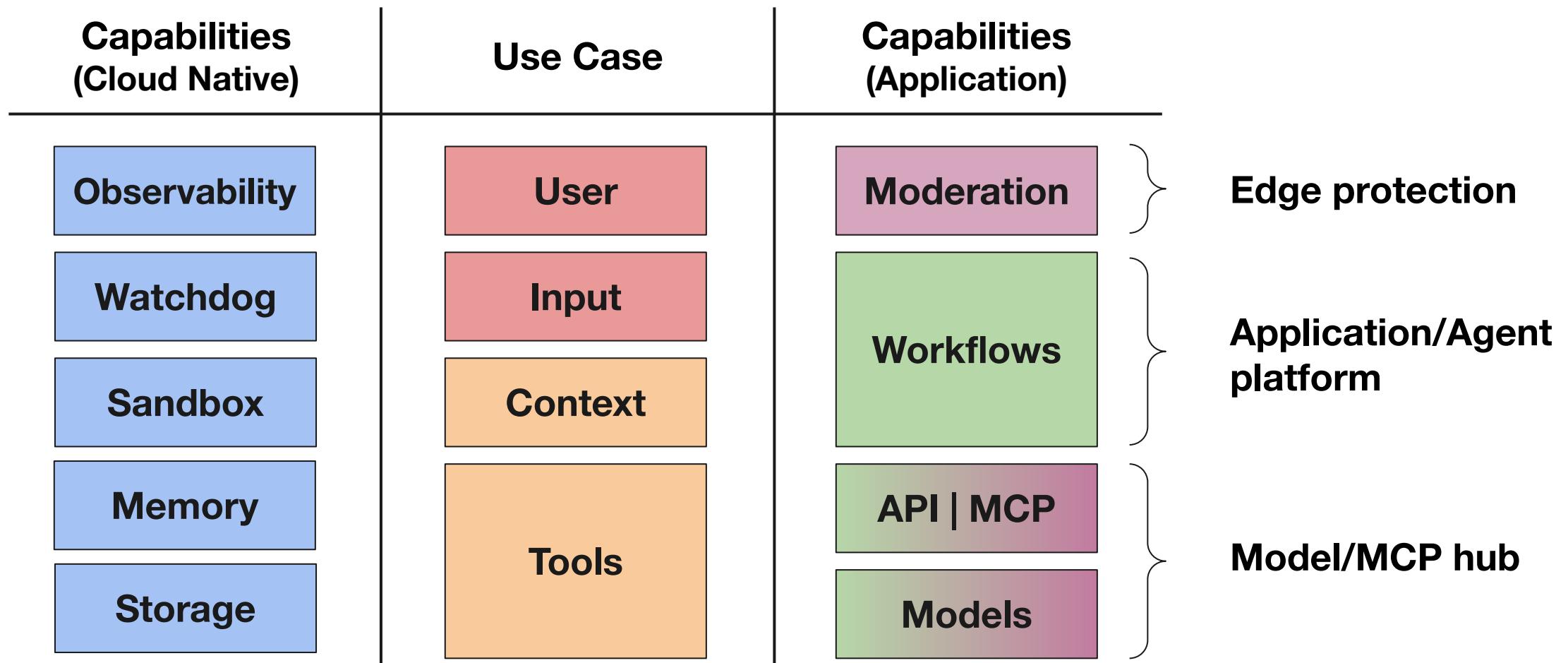
Challenges per agency level

Agency level	Type	Challenges
★★★	Simple processor	hallucinations, latency, model access, prompt quality, prompt injection, data poisoning
★★★	Router	limited capabilities: pre-defined, static workflows
★★★	Tool caller	access to diverse tool libraries, wrong tool , params, invalid output
★★★	Multi-step agents	(non-) predictability, planning failures, communication, large prompts
★★★	Multi-agent	security, rogue agents, endless loops



Cloud Native Agents

Capabilities



A photograph of a woman with short brown hair, wearing a flowing pink dress, dancing joyfully in a bright green grassy field under a clear blue sky.

The State of GenAI & ML in Cloud Native Ecosystem

Best Practices

Overview of Best
Practices in Industry



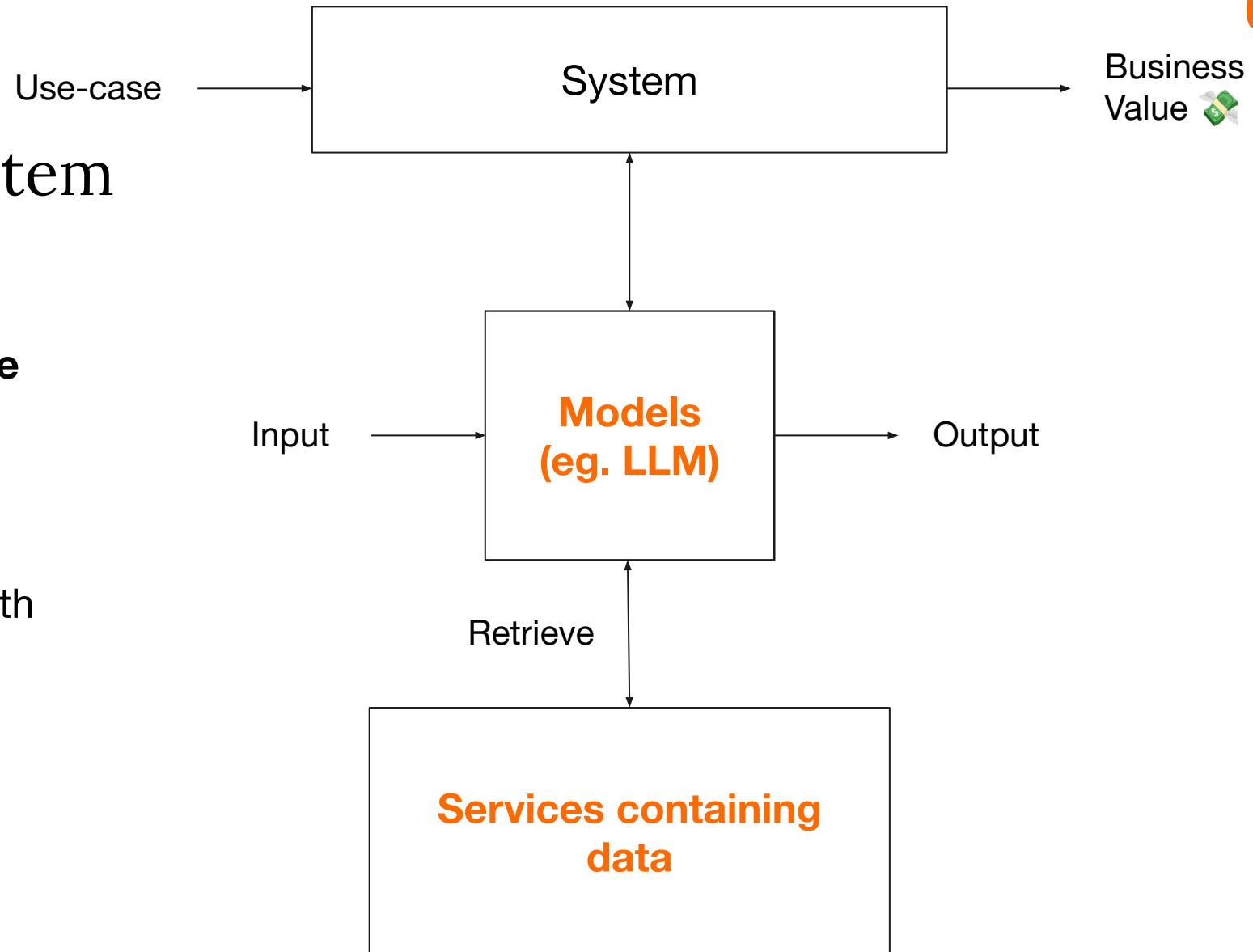
AI Ecosystem

From AI Model to AI System

Evolution from specific **applied science capabilities** to sophisticated **organisational capabilities**.

Encompassing **complex use-cases** with clear tangible **business value**.

Leveraging advancements in **software operations** applied into **AI delivery**.





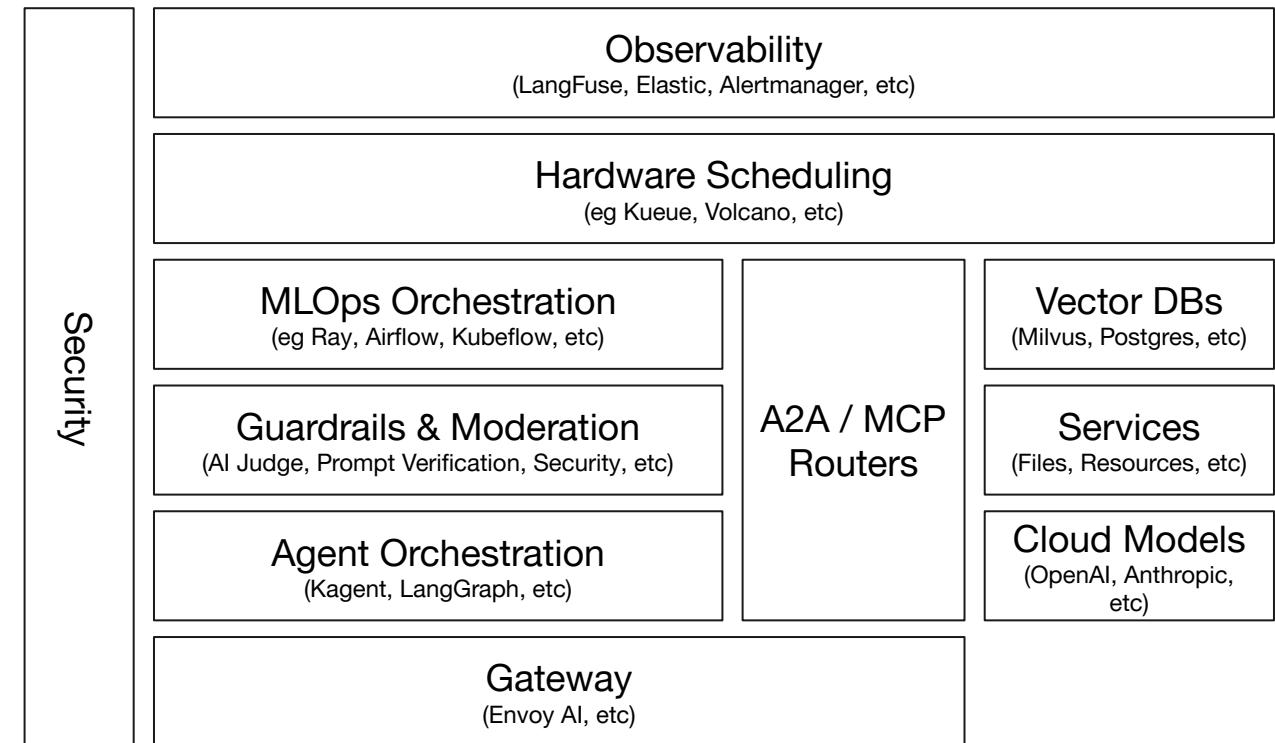
AI Ecosystem

From Model- to Data-Centric

The **core potential** is not on the models themselves, but in the **interconnected data-centric ecosystem**

This includes **complex orchestration** across **applications, infrastructure** and **complex data flow** in agentic systems

This extracts best practices from **traditional MLOps** but brings **further complexities** and **tooling requirements**



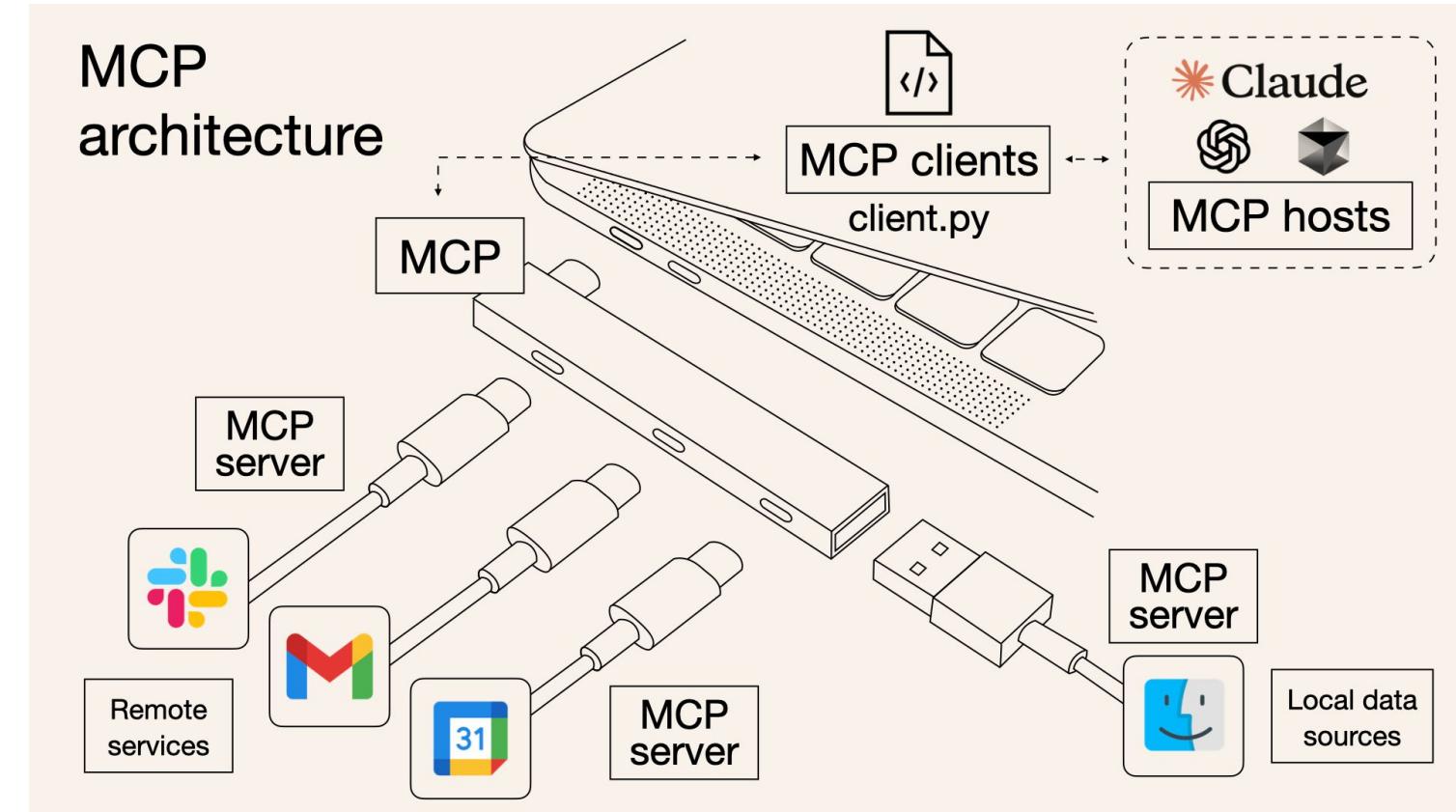
AI Ecosystem

The “Boring” Data

The archetypal **garbage-in = garbage-out** is true now more than ever across our intelligent systems and processes

This requires equally new protocols that drive **standardisation** and best practice for **interoperability** (eg envoy, MCP, etc)

Strategic investment in the “Boring” data is key to be **successful in the long term**, whilst **capturing value in the short term**.





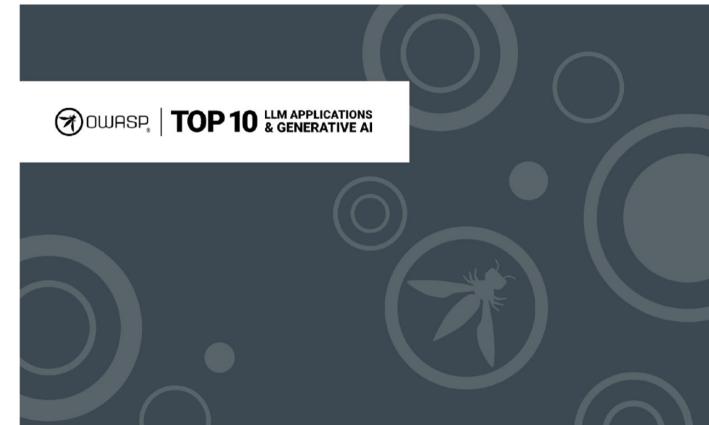
ML Security

Considerations of Security

Taxonomy on Agents and Agentic Systems, both for **single- and multi-agent** contexts

Taxonomy on Agentic System **Vulnerabilities and Exploits**

Examples include: Data poisoning (scrape data) Memory poisoning, remote code execution, identity spoofing, multi-agent chain exploits, cascading hallucination attacks, intent breaking, etc



OWASP | **TOP 10** LLM APPLICATIONS & GENERATIVE AI

Agentic AI – Threats and Mitigations

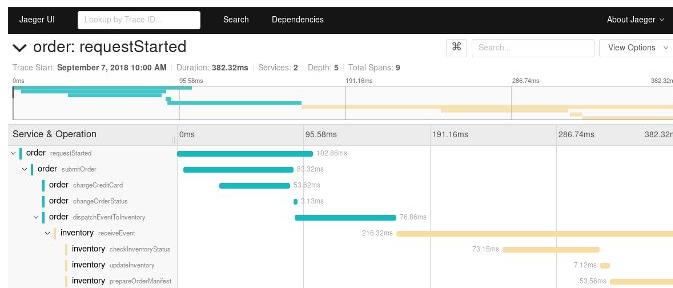
OWASP Top 10 for LLM Apps & Gen AI
Agentic Security Initiative

Version 1.0
February 2025

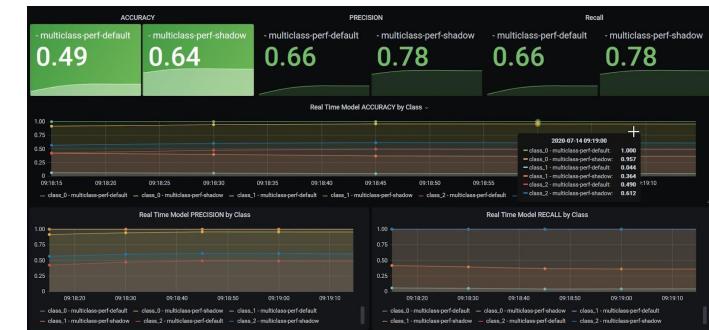
MLOps Monitoring

Paradigms of Observability

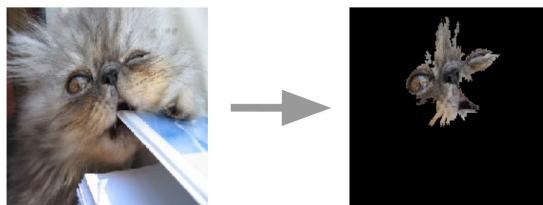
Service Performance



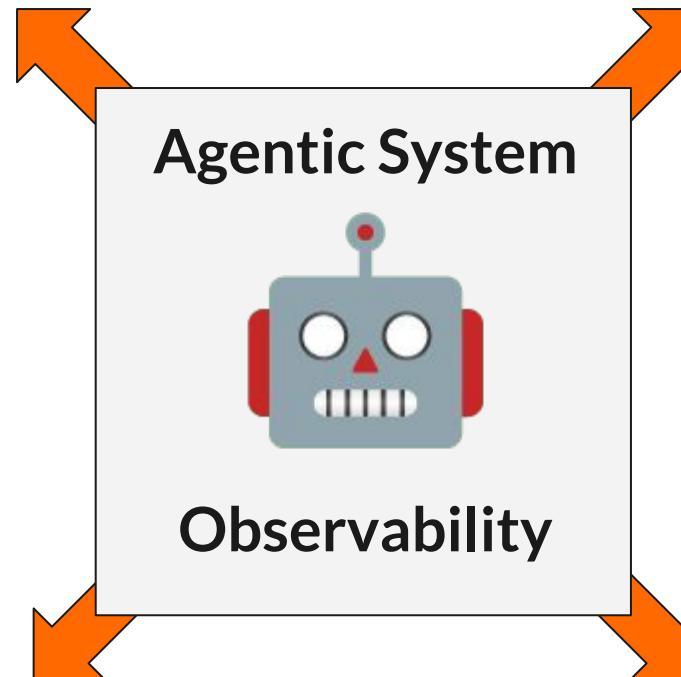
LLM Perf / Debugging



Explainability



a powerful study of loneliness sexual UNK and desperation by patient UNK up the atmosphere and pay attention to the wonderfully written script br br i praise robert altman this is one of his many films that deals with unconventional fascinating subject matter this film is disturbing but it's sure to UNK a strong emotional response from the viewer if you want to see an unusual film some might even say bizarre this is worth the time br unfortunately it's very difficult to find in video stores you may have to buy it off the internet



Outlier & Drift Detection





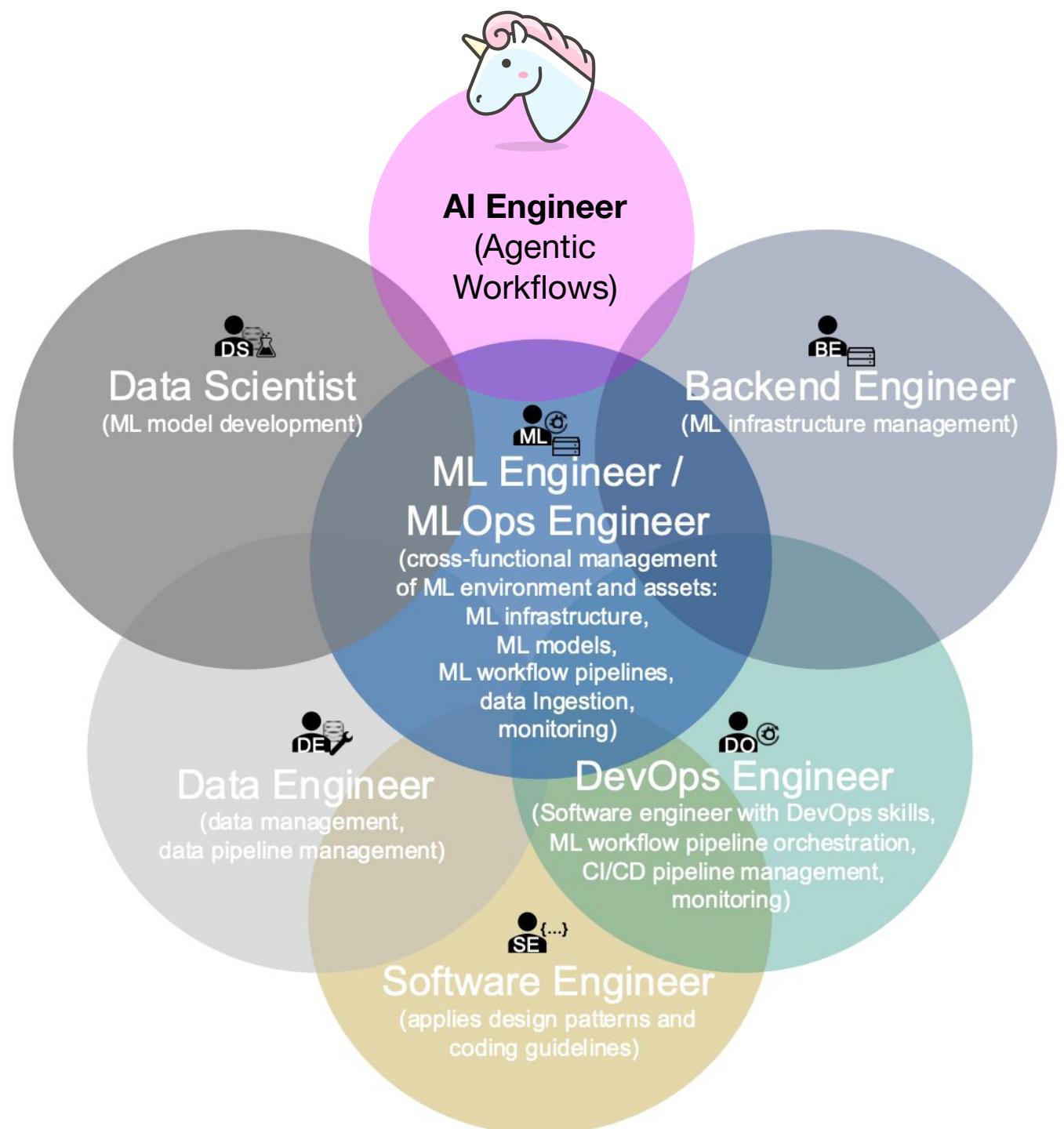
Organisational Process

From SDLC to LLM-DLC

Identifying the **personas** throughout the different **stages** of machine learning development and operation.

Defining the organisation-wide **architectural blueprint** for production machine learning operations.

Bringing together personas, blueprint and processes through the **Development Lifecycle** for large scale **agentic systems**





The State of GenAI & ML in Cloud Native Ecosystem

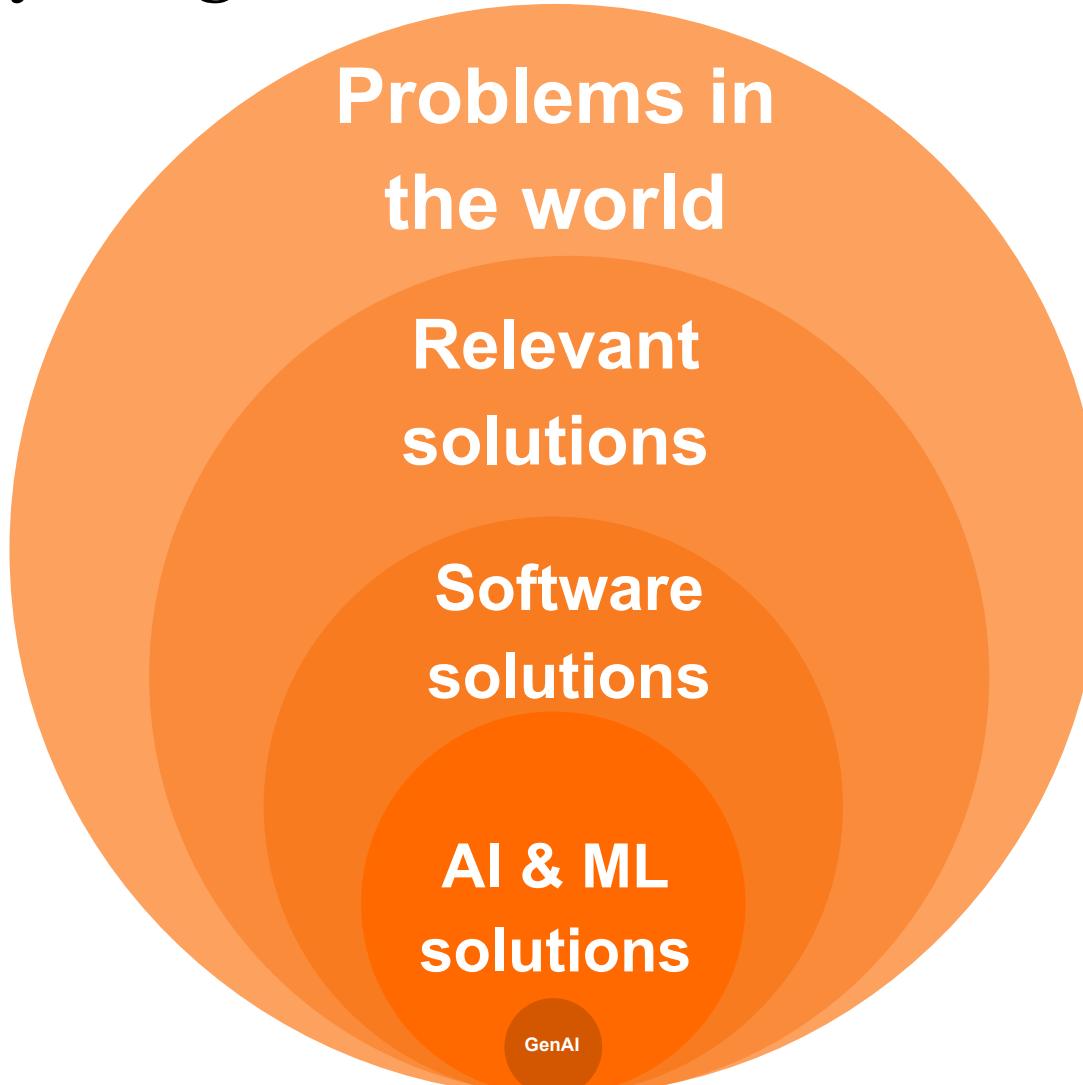
Closing Thoughts

Key Areas to Highlight
and Resources



The Reality

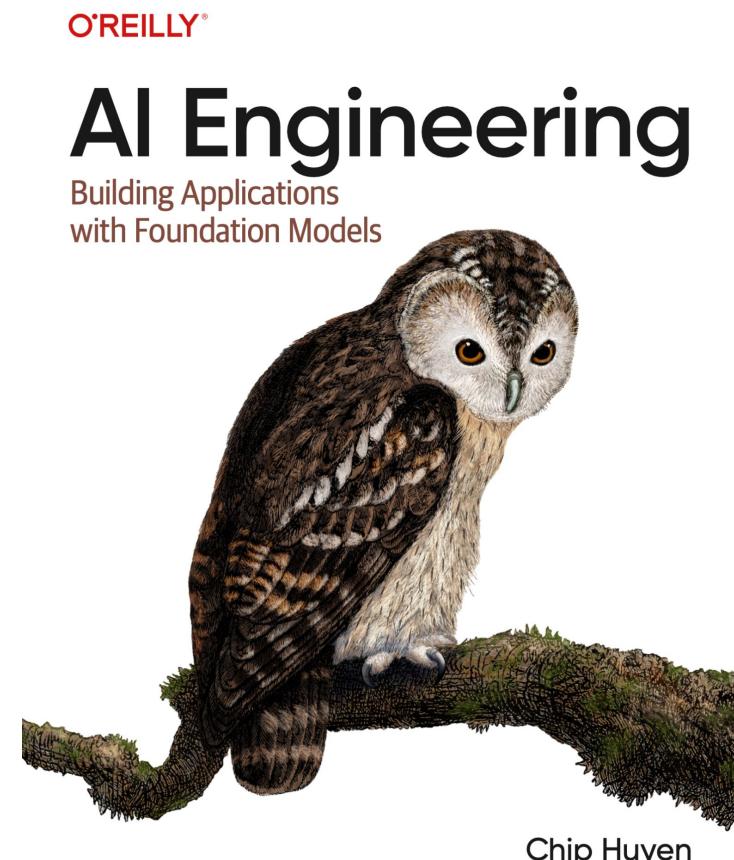
Not Everything is AI



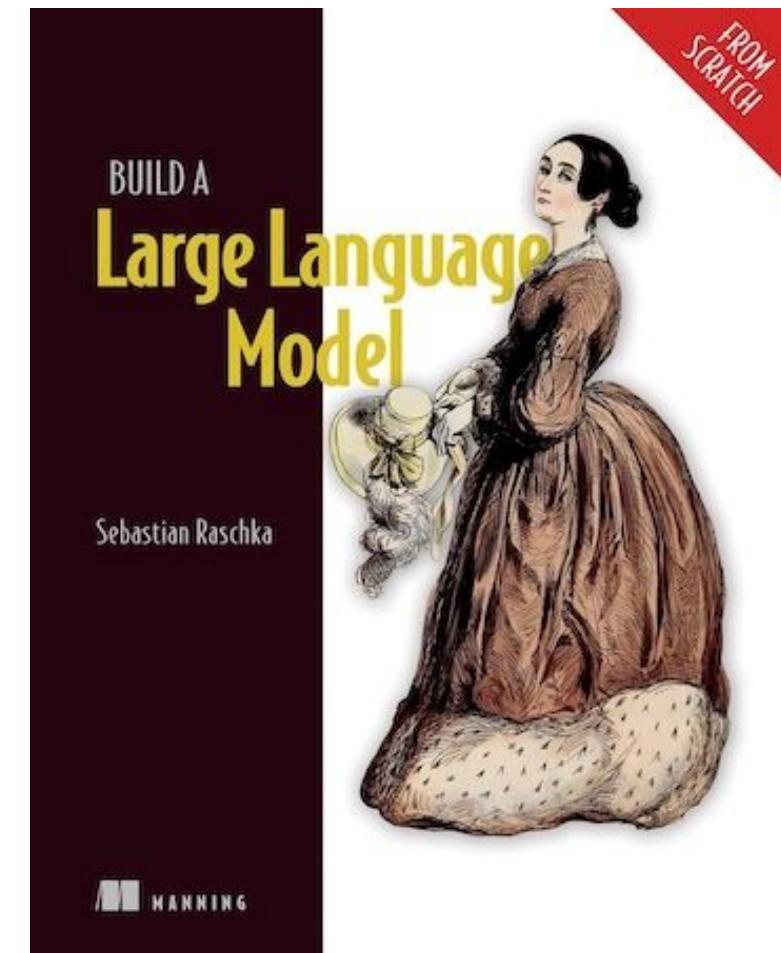
**When you run
around with a
hammer
everything may
look like a
nail**

Helpful resources

Deep Dives



<https://huyenchip.com/books/>



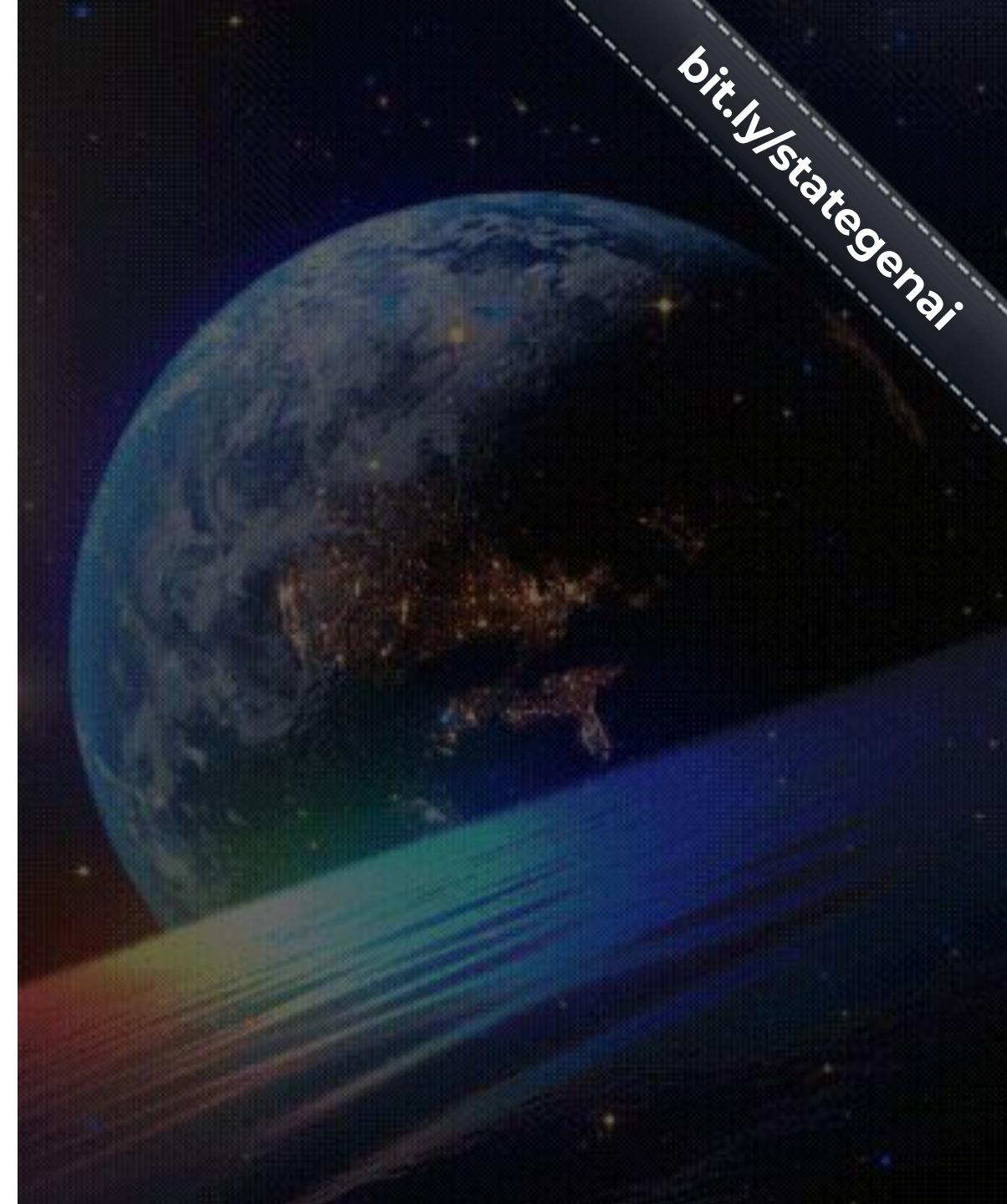
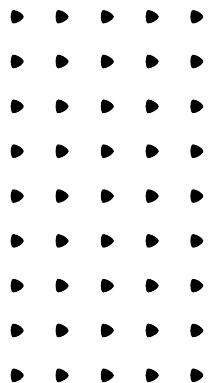
<https://sebastianraschka.com/books/>

The State of GenAI & ML in the Cloud Native Ecosystem

Alejandro Saucedo
Bartosz Ocytko



zalando





Challenges & Patterns

per agency level

Agency	Type	Challenges	Patterns
★★★	Simple processor	Hallucinations, latency, model access, prompt quality, prompt injection, <u>data poisoning</u>	Evals, Model proxies, Prompt refinement
★★★	Router	Pre-defined, static workflows	CoT, ReAct, ADaPT, Reflexion
★★★	Tool caller	Access to diverse tool libraries, wrong tool, params, invalid output	Structured Outputs, MCP (Model Context Protocol)
★★★	Multi-step agents	(non-) predictability, planning failures, <u>communication</u> , <u>large prompts</u>	workflow orchestration, visual debuggers, memory, pub/sub, channels, rendezvous (Ada '80s)
★★★	Multi-agent, autonomous agents	<u>security</u> , <u>rogue agents</u> , <u>endless loops</u>	watchdog, sandbox