

# GenAI (...LLMs) Everywhere... - AI Race!

# DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI  
research@deepseek.com

## Abstract

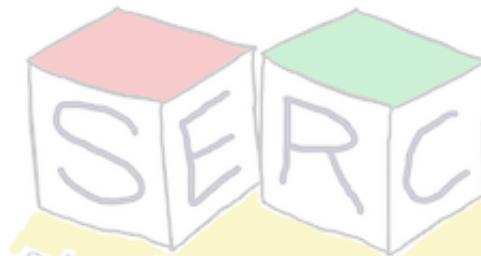
We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-tuning (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1-Zero, DeepSeek-R1, and six dense models (1.5B, 7B, 8B, 14B, 32B, 70B) distilled from DeepSeek-R1 based on Qwen and Llama.

# Introducing the SWE-Lancer benchmark

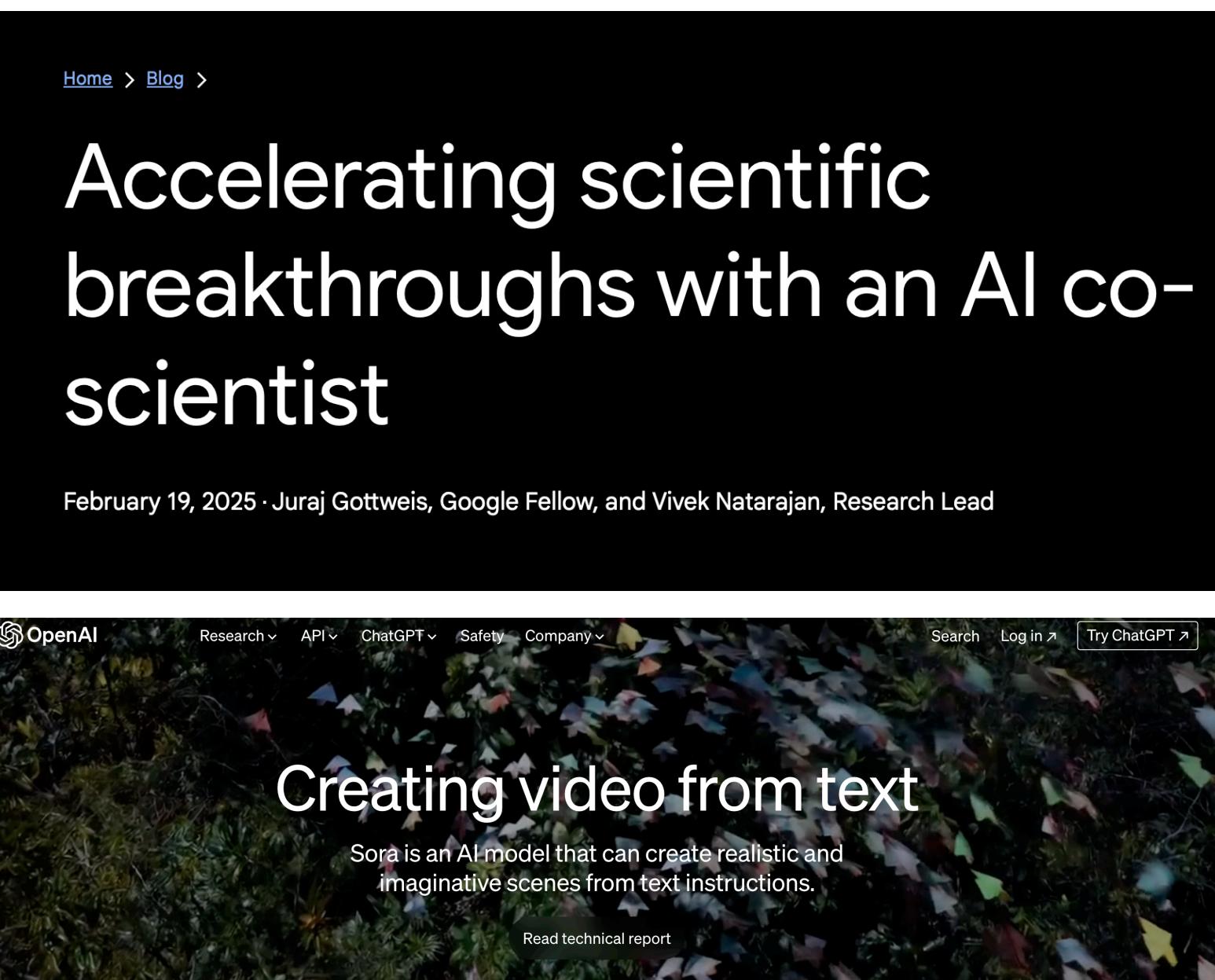
## Can frontier LLMs earn \$1 million from real-world freelance software engineering?

[Read paper ↗](#)

[Access repository ↗](#)



**Source:** openai.com, meta.com, cursor.com



index.mdx M

changelog > website > src > pages > index.mdx

```
10 }
11
12 ---
13
14 ## Terminal Debugger (v0.2.12) {{ date: '2023-05-17T00:00Z' }}
15
16 ### 🖥 In-terminal Debugging
17
18 - Press Cmd+D to auto-debug a terminal error
19 - Press Cmd+Shift+L, and the model will add the terminal context to chat
20
21 ### 📌 Activity Bar pinning
22
23 - You can Pin custom extensions to your activity bar in the top left
24 <div
25   style={{
26     'padding-top': '1rem',
27   }}
28 >
29 <CustomImage
30   src={pinnedExtensions}
31   width="200"
32   height="200"
33   style={{
34     // Place in the center
35   }}</CustomImage>
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL GITLENS

node - website

```
at renderNodeDestructiveImpl (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:610:14)
4:11) at renderNodeDestructive (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:6076:14)
) at renderIndeterminateComponent (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:5785:7)
at renderElement (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:5946:7)
at renderNodeDestructiveImpl (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:610:14)
4:11) at renderNodeDestructive (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:6076:14)
) at renderNode (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:6259:12)
at renderChildrenArray (/Users/amansanger/eversphere/changelog/website/node_modules/react-dom/cjs/react-dom-server.browser.development.js:6211:7) {
digest: undefined
}
```

%DL insert to chat, %D auto debug



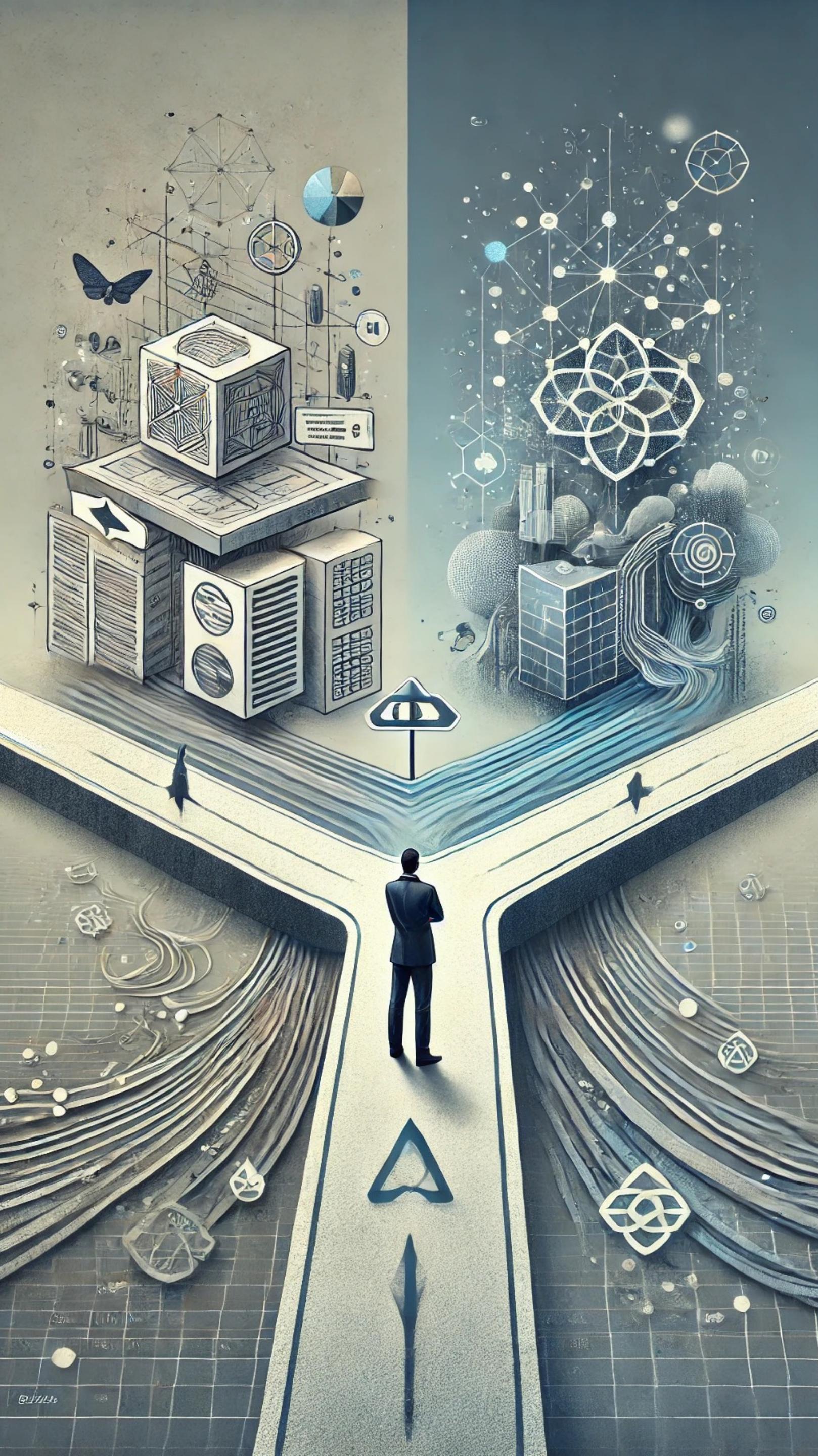
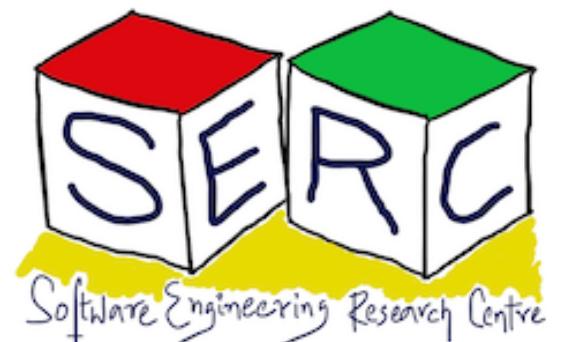
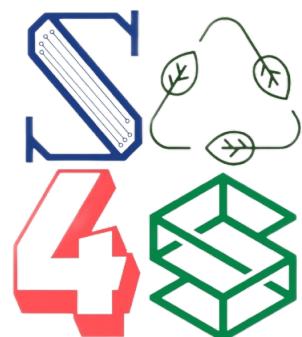
Even the car that dropped me ;)

# Navigating the Crossroads of GenAI and SE: Insights from our Research

Karthik Vaidhyanathan

European Data Intensive Software Systems Winter School 2025

Feb 24, 2025, L'Aquila





# ABOUT ME

Logic takes you from A to B, Imagination takes you elsewhere -- Albert Einstein



**Karthik Vaidyanathan**

Assistant Professor

Software Engineering Research Center and  
Leadership Member, Smart City Research Center

IIIT Hyderabad, India



## Education



Double Master Degree - Software  
Architecture and Machine Learning  
PhD from GSSI, Italy  
Postdoc, University of L'Aquila, Italy



## Fun Facts!

1. Cricket fanatic!
2. Movie buff!!
3. From God's own Country!!

**Research Interests**

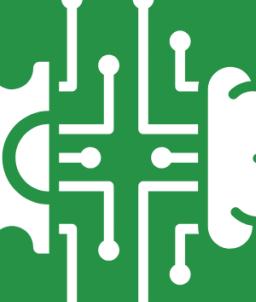
**AI4SA**

- 1. AI for Architectural Knowledge
- 2. AI for self-adaptation

**SA4AI**

- 1. Sustainable AI-enabled systems
- 2. Self-adaptive AI Systems (Edge-Cloud)



# Quest for Artificial Intelligence

## Next word prediction problem

Google

What is software

- what is software
- what is software engineering
- what is software testing
- what is software development
- what is software and hardware
- what is software as a service
- what is software piracy
- what is software architecture

Google

software engineering is

- software engineering is **dead**
- software engineering is a **layered technology**
- software engineering is **easy or hard**
- software engineering is **hard**
- software engineering is **what**
- software engineering is **primarily concerned with**

Detected language: English

how are you?

See dictionary

Italian Spanish English

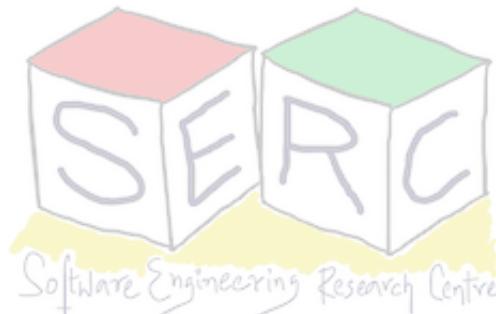
12 / 5,000

Italian Spanish English

Come stai?

Send feedback

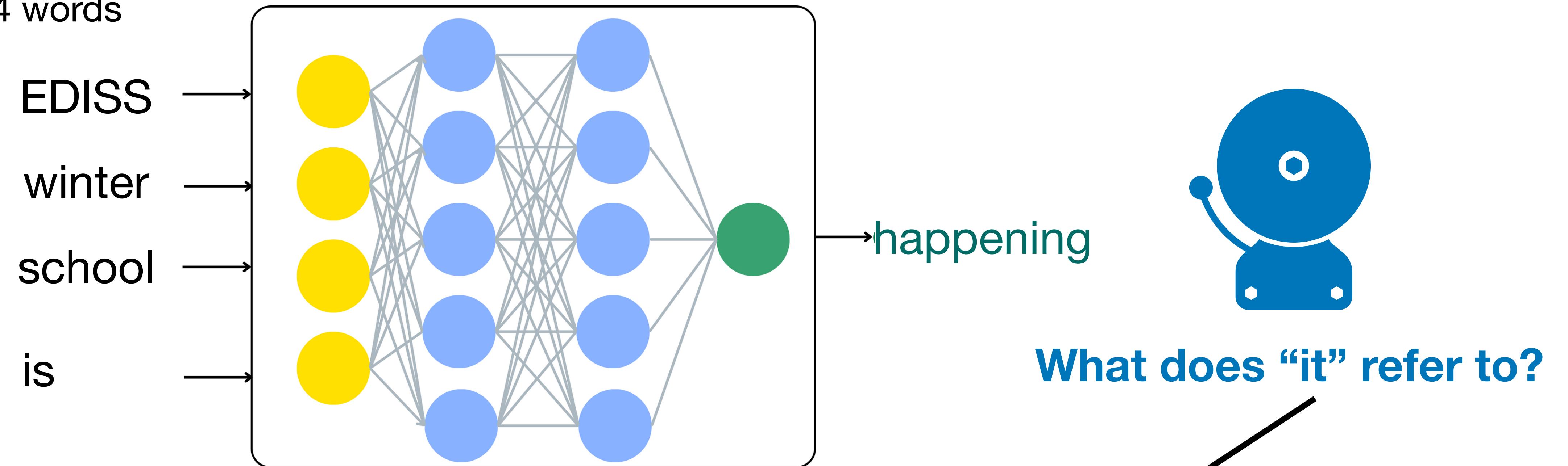
Problem since 1950's!!!



# Its been a long journey MCP, ..backprop...RNN CNN,.....

## Predict the next word in a sequence!

Eg: Context of 4 words



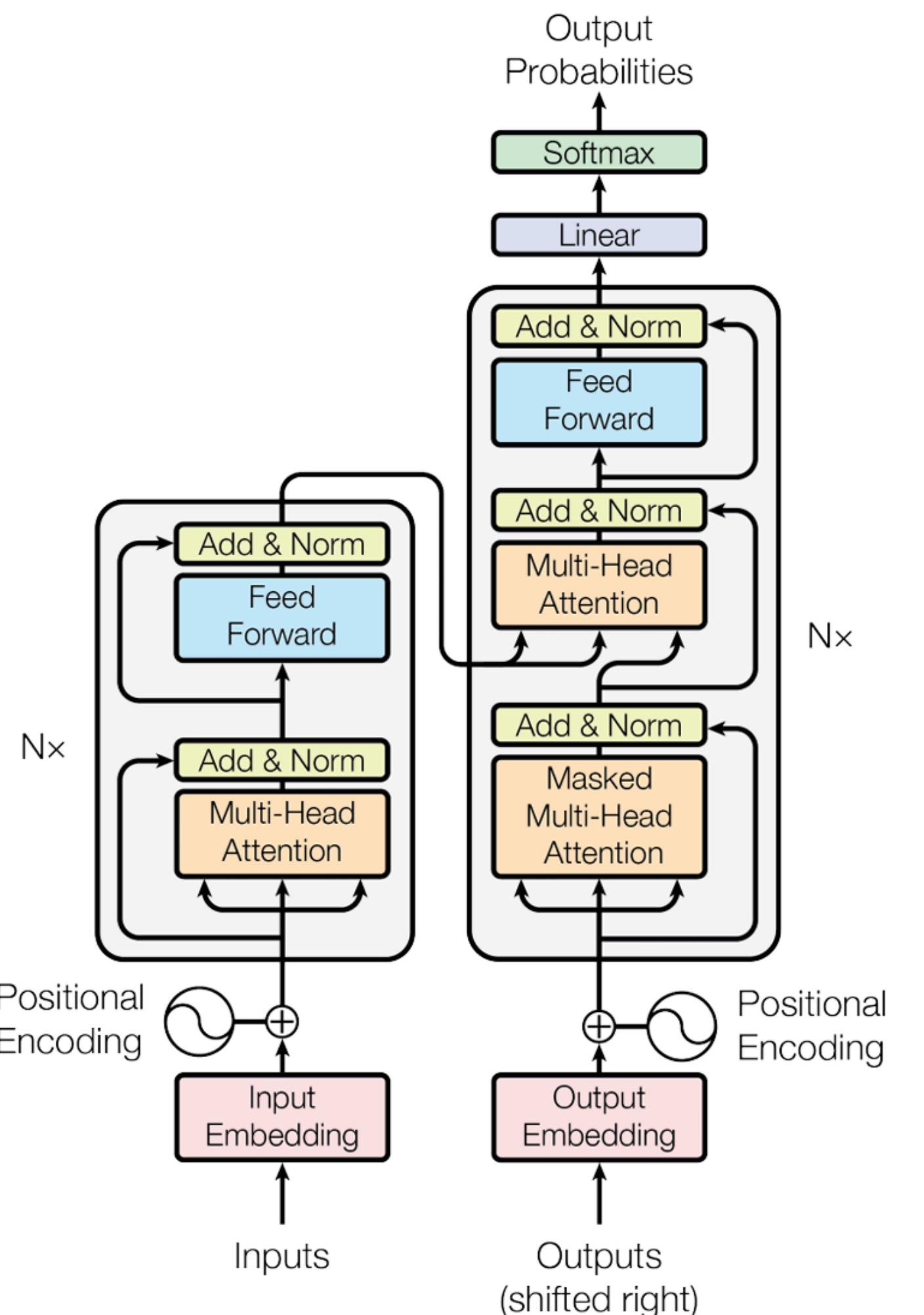
EDISS winter school is happening in

EDISS winter school is happening in L’Aquila.

EDISS winter school is happening in L’Aquila. **It**

# Attention is all you need!

## Age of Transformers



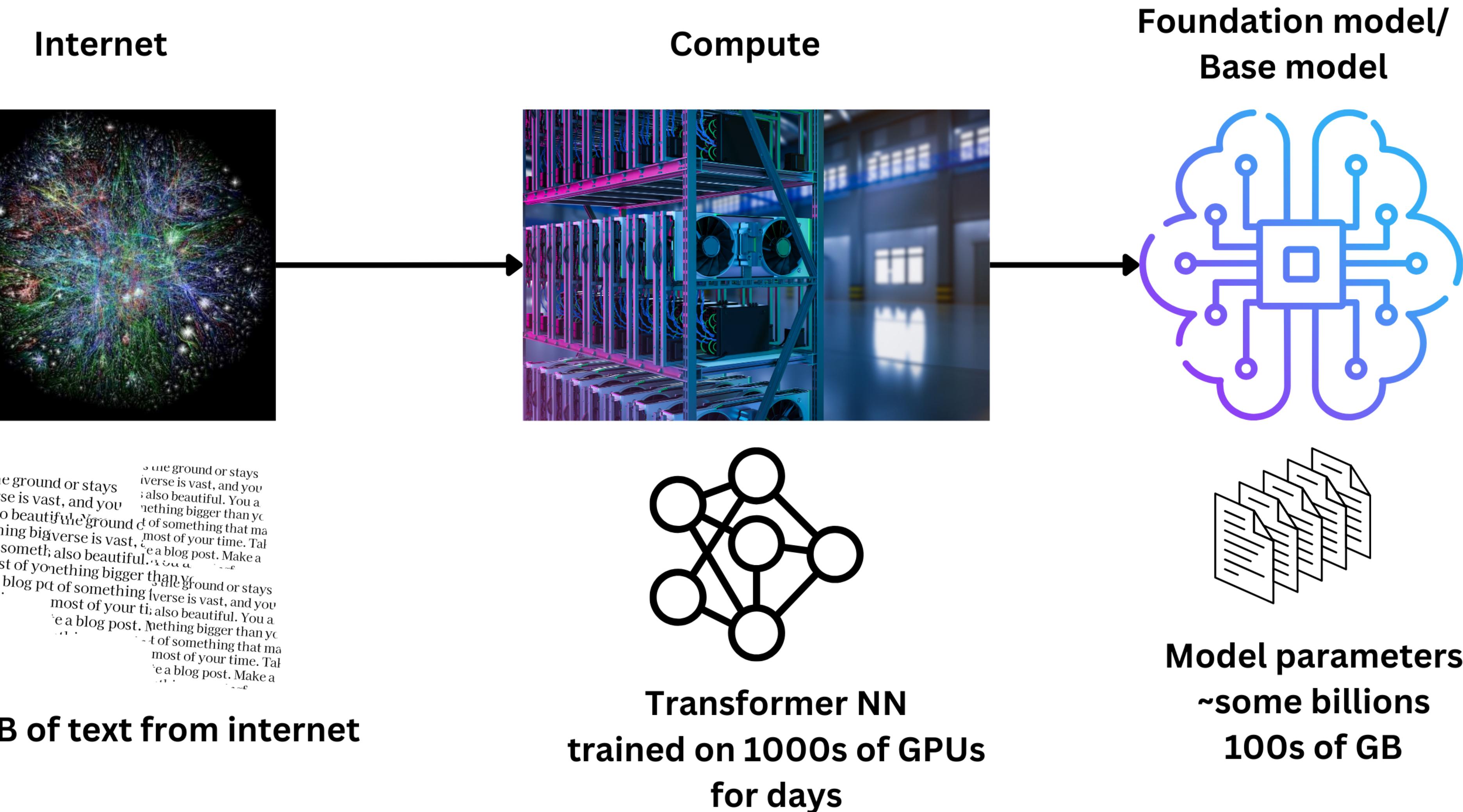
EDISS  
Winter  
School  
is  
happening  
in  
L'Aquila  
It

EDISS  
Winter  
School  
is  
happening  
in  
L'Aquila  
It



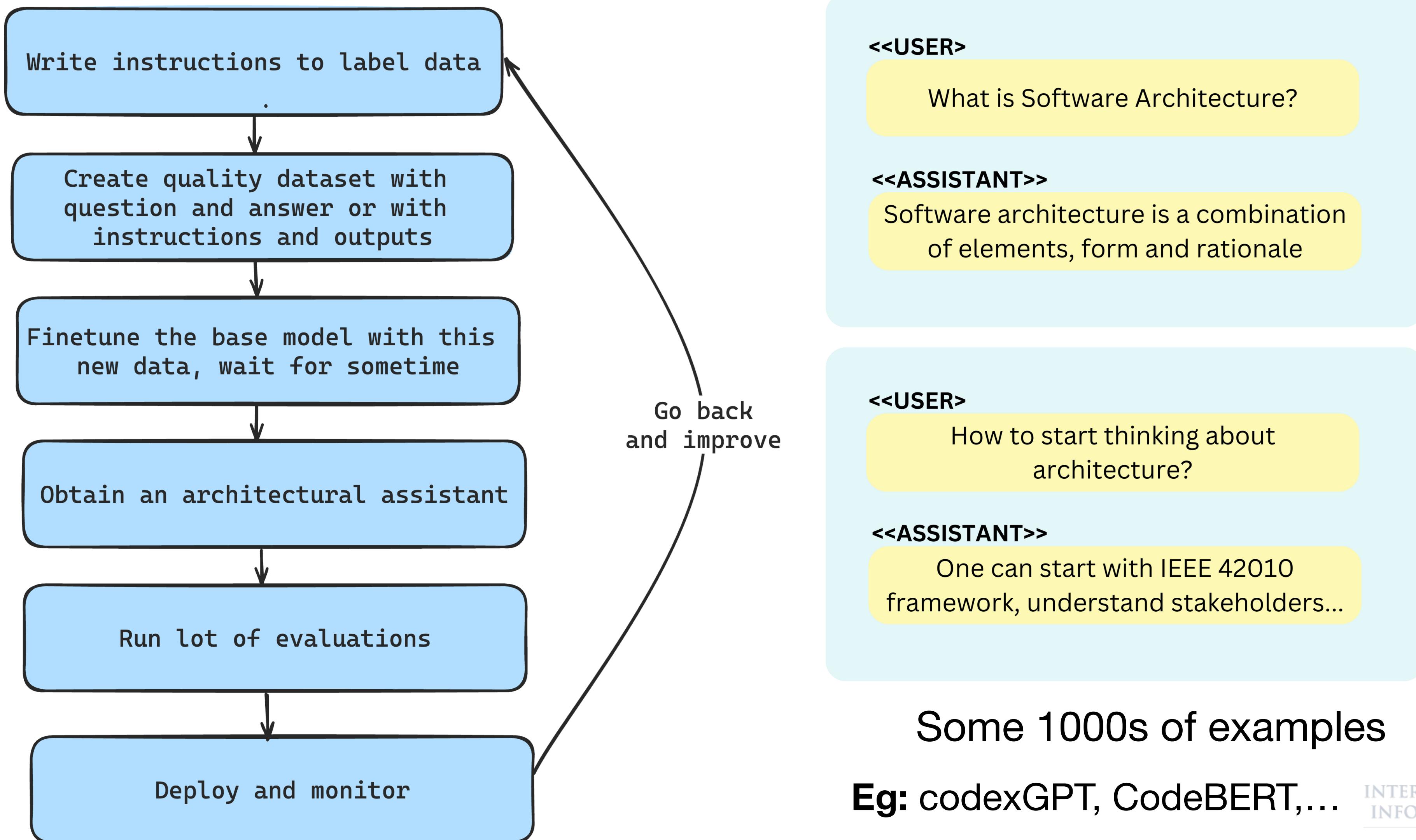
# “Large” Language Models (LLM)

Do you have a ton of text and compute power?

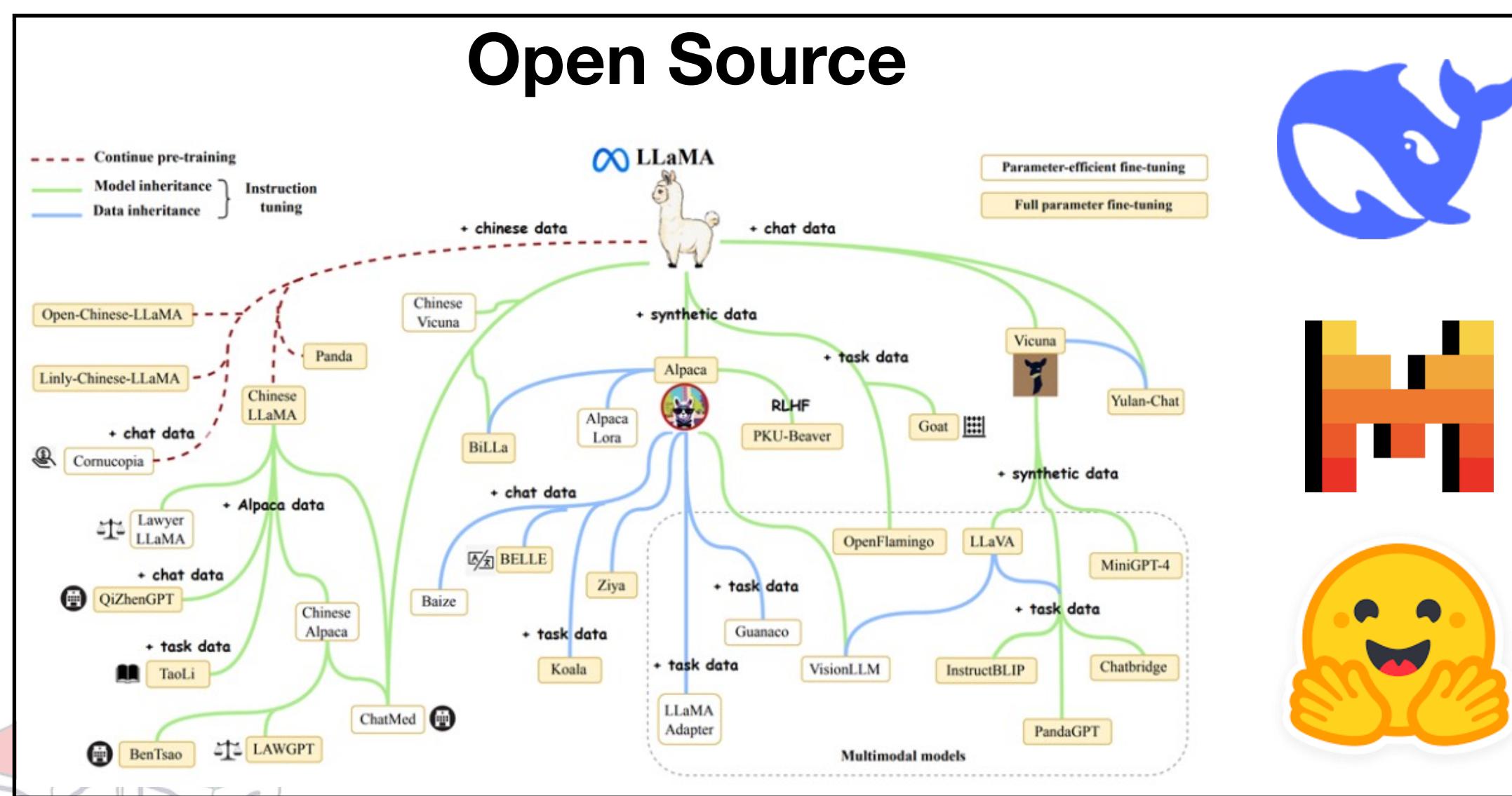
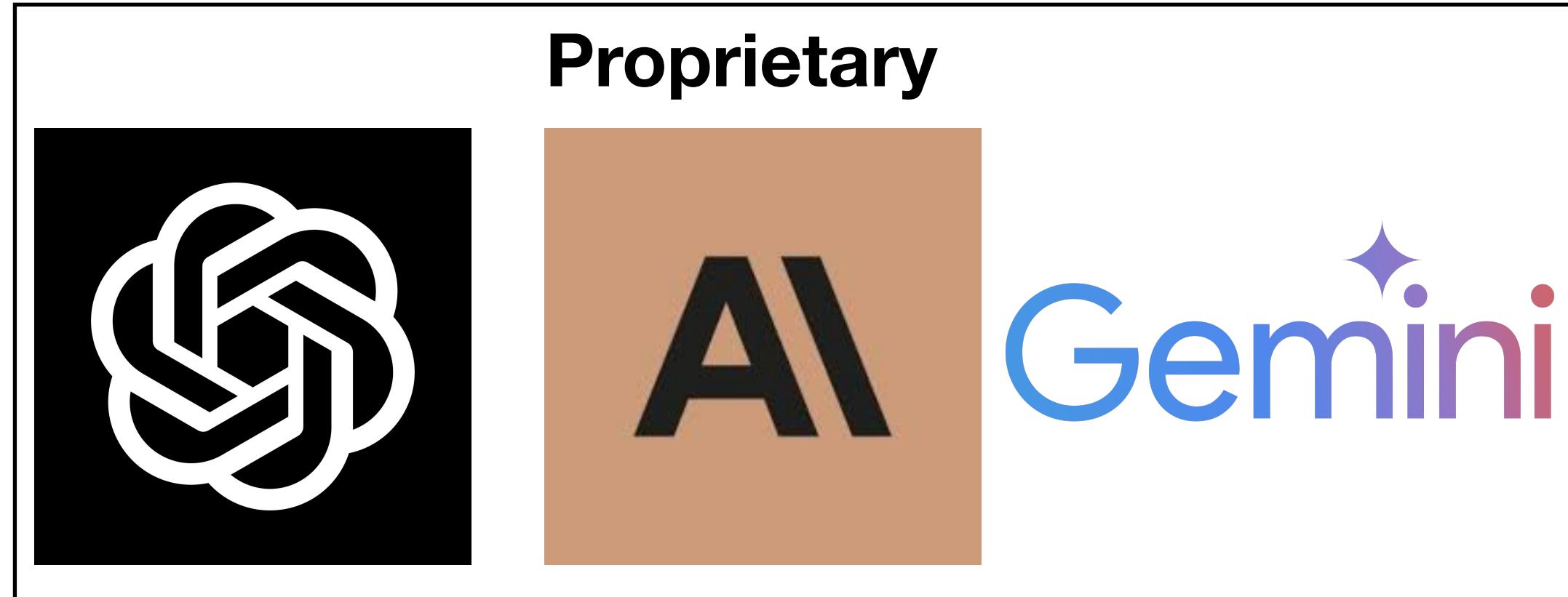


# I want something more specific to Architecture

May be you can fine-tune and create your model



# Today we have different varieties of LLMs



LMSYS Chatbot Arena Leaderboard							
Arena Elo							
<a href="#">Vote</a>   <a href="#">Blog</a>   <a href="#">GitHub</a>   <a href="#">Paper</a>   <a href="#">Dataset</a>   <a href="#">Twitter</a>   <a href="#">Discord</a>							
LMSYS <a href="#">Chatbot Arena</a> is a crowdsourced open platform for LLM evals. We've collected over 500,000 human preference votes to rank LLMs with the Elo ranking system.							
Rank	Model	Arena Elo	95% CI	Votes	Organization	License	Knowledge Cutoff
1	<a href="#">Claude 3 Opus</a>	1256	+3/-4	47589	Anthropic	Proprietary	2023/8
1	<a href="#">GPT-4-1106-preview</a>	1254	+3/-4	62657	OpenAI	Proprietary	2023/4
1	<a href="#">GPT-4-0125-preview</a>	1250	+3/-3	47631	OpenAI	Proprietary	2023/12
4	<a href="#">Bard (Gemini Pro)</a>	1208	+5/-5	12468	Google	Proprietary	Online
4	<a href="#">Claude 3 Sonnet</a>	1204	+3/-3	57740	Anthropic	Proprietary	2023/8
6	<a href="#">Command R+</a>	1194	+5/-5	17404	Cohere	CC-BY-NC-4.0	2024/3
6	<a href="#">GPT-4-0314</a>	1189	+4/-3	41292	OpenAI	Proprietary	2021/9
8	<a href="#">Claude 3 Haiku</a>	1182	+3/-4	50689	Anthropic	Proprietary	2023/8
9	<a href="#">GPT-4-0613</a>	1164	+3/-3	60213	OpenAI	Proprietary	2021/9
9	<a href="#">Mistral-Large-2402</a>	1158	+3/-4	35075	Mistral	Proprietary	Unknown
10	<a href="#">Qwen1.5-72B-Chat</a>	1153	+4/-5	27050	Alibaba	Qianwen LICENSE	2024/2

<https://chat.lmsys.org/>



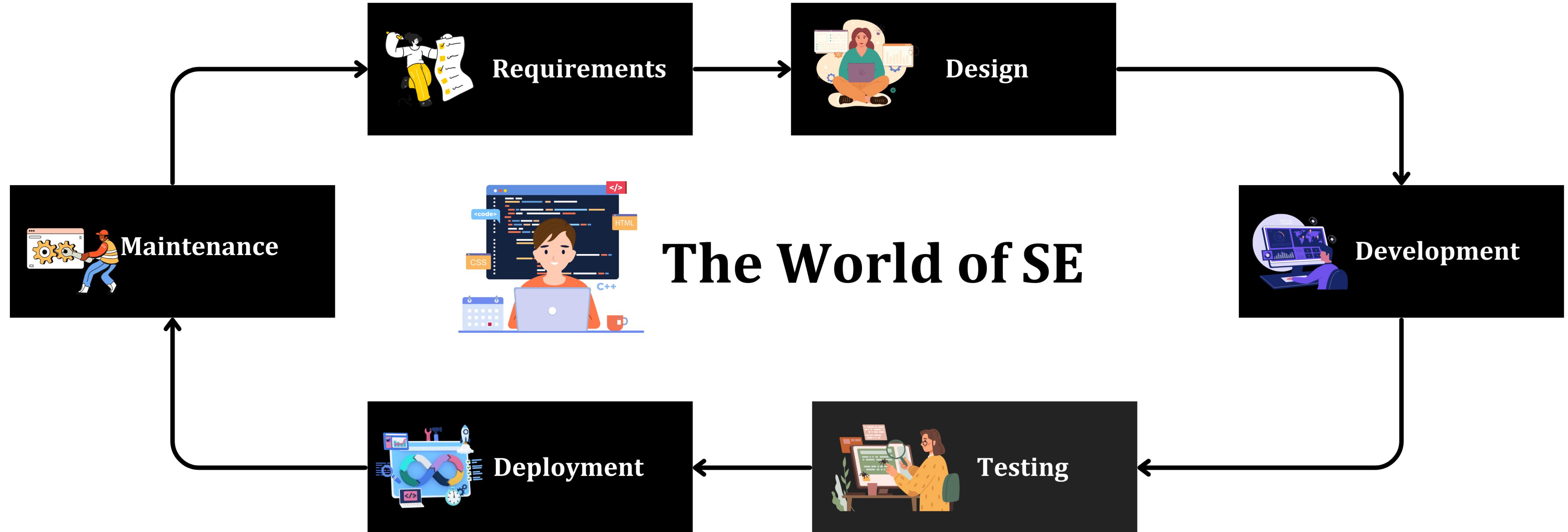
# Bottomline: Text is an Abstraction of Reality!



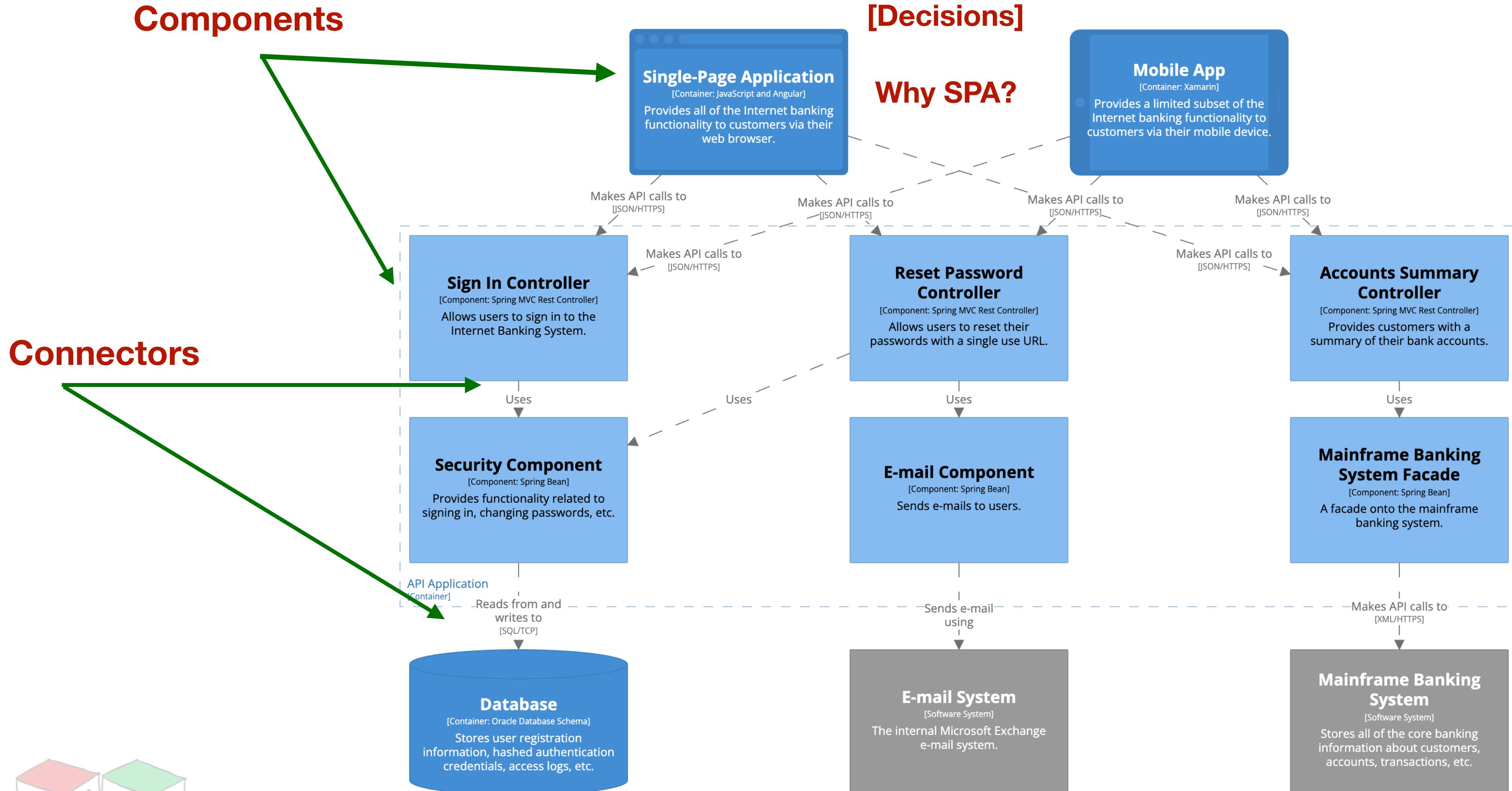
The image shows a lively public square with a central fountain featuring a statue, surrounded by historic buildings, street lamps, and people walking, set against a bright blue sky with scattered clouds.

A sunny day in a wide, paved European square with a central fountain, surrounded by classic buildings and parked cars.

# The World of Software Engineering



# Starting with Design: Architecture an Abstraction



# Design Decisions is all you need!

## Software Architecture as a Set of Architectural Design Decisions

Anton Jansen

Department of Computing Science  
University of Groningen  
PO BOX 800, 9700 AV, The Netherlands  
[anton@cs.rug.nl](mailto:anton@cs.rug.nl)

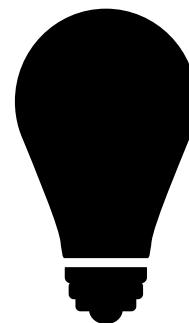
Jan Bosch

Software & Application Technologies Lab  
Nokia Research Center  
PO BOX 407, FI-00045, Finland  
[jan.bosch@nokia.com](mailto:jan.bosch@nokia.com)

### Abstract

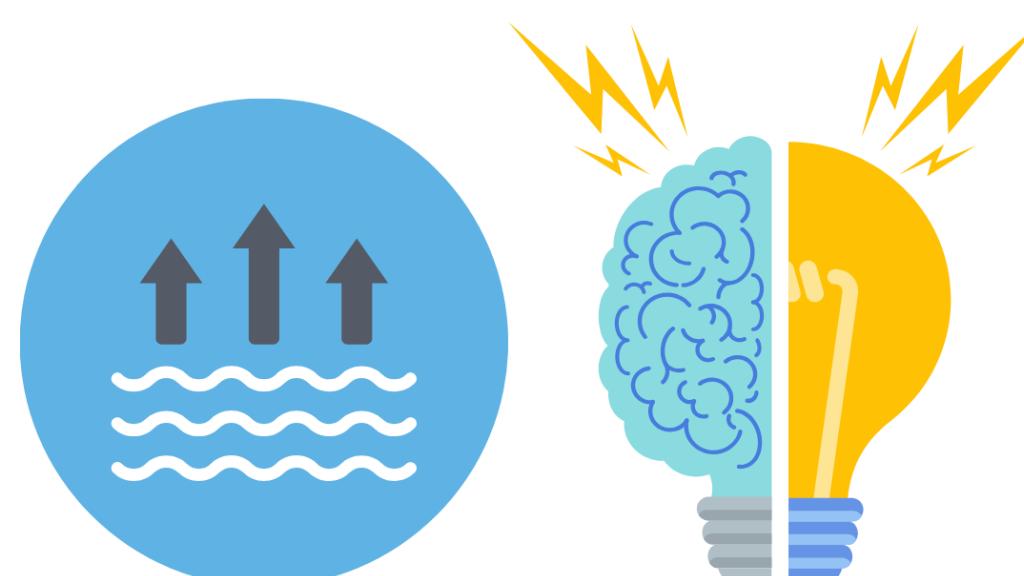
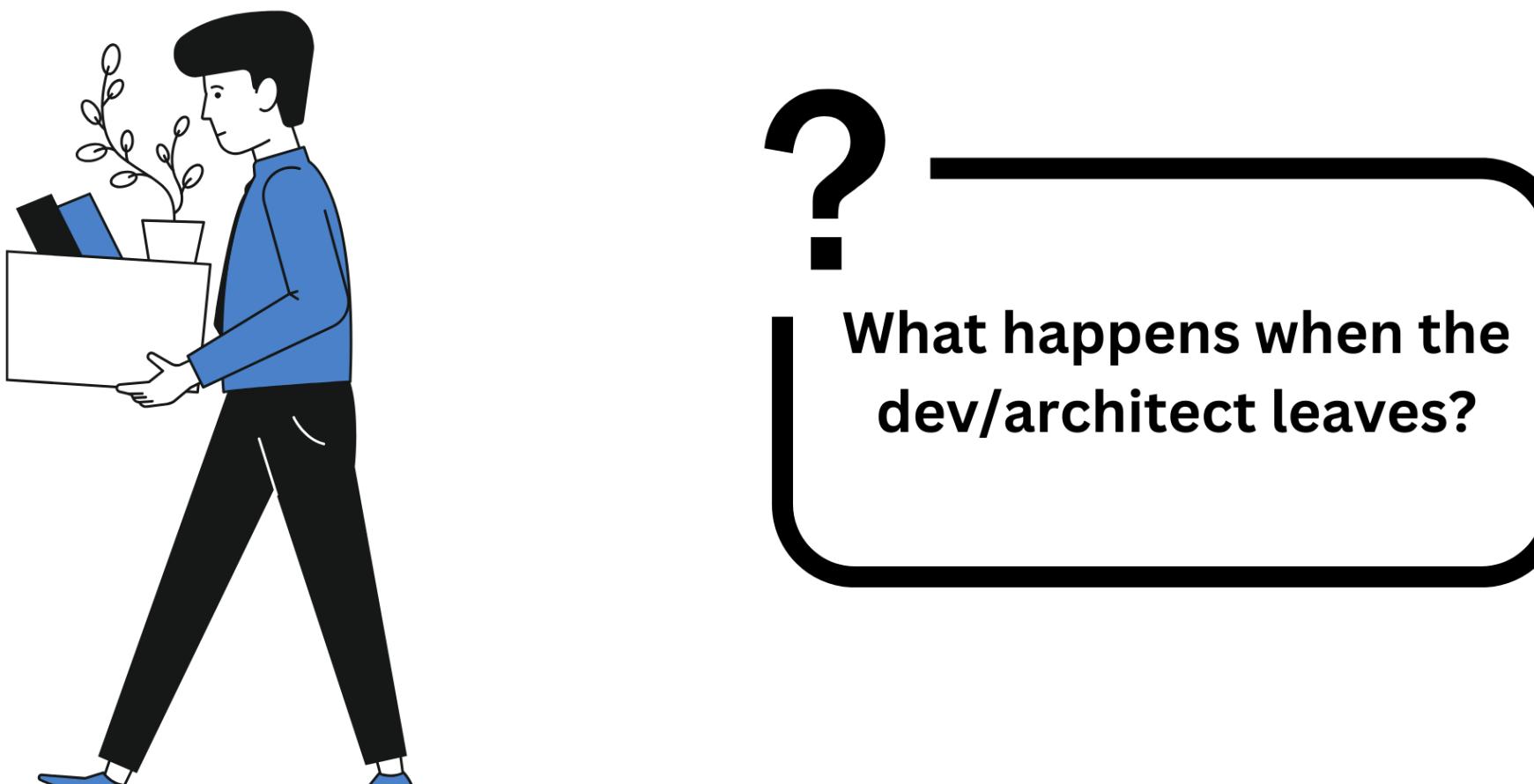
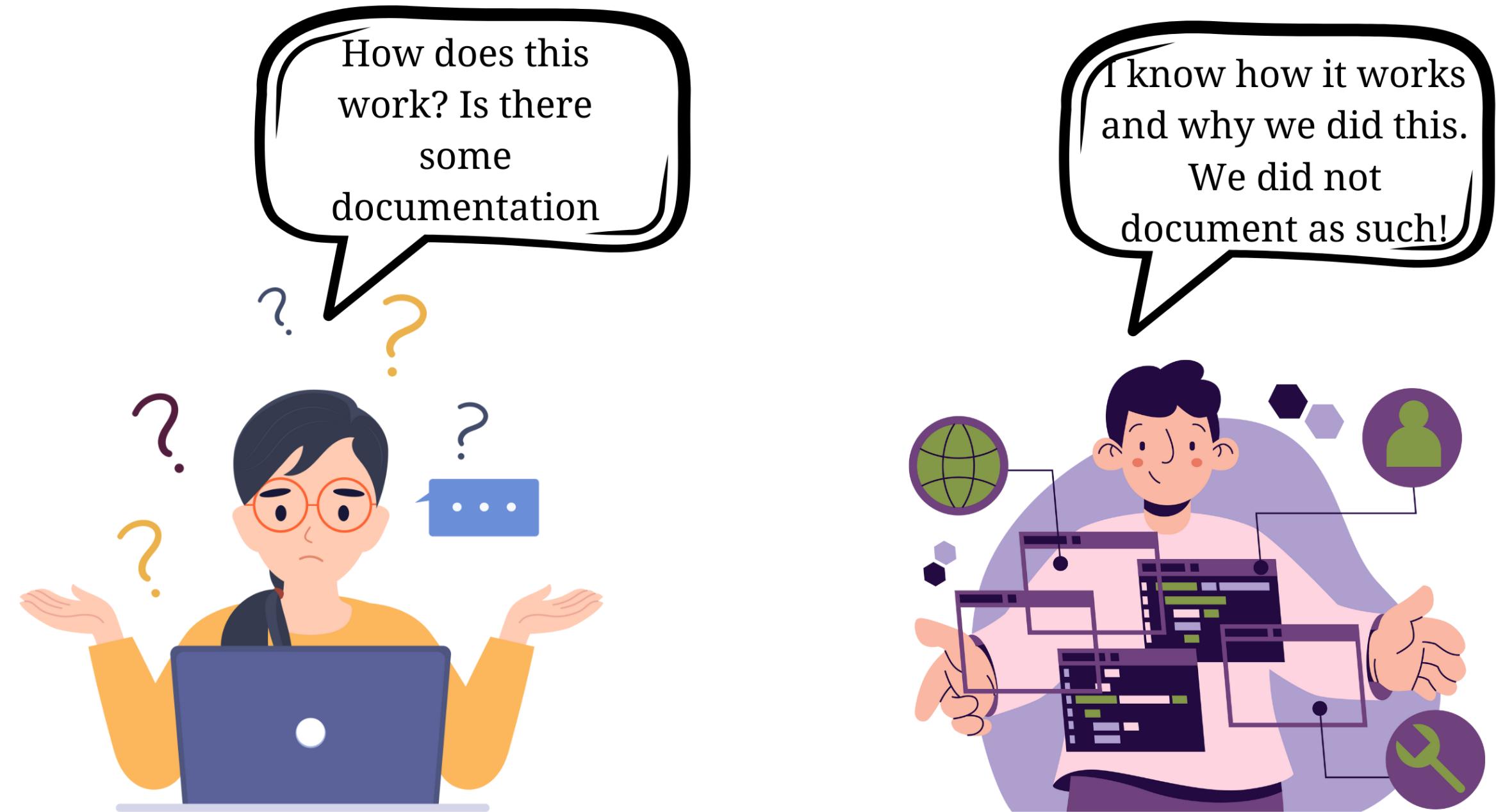
*Software architectures have high costs for change, are complex, and erode during evolution. We believe these problems are partially due to knowledge vaporization. Currently, almost all the knowledge and information about the design decisions the architecture is based on are implicitly embedded in the architecture, but lack a first-class repre-*

this notion of architectural design decisions, although architectural design decisions play a crucial role in software architecture, e.g. during design, development, evolution, reuse and integration of software architectures. In design, the main concern is which design decision to make. In development, it is important to know which and why certain design decisions have been taken. Architecture evolution is about making new design decisions or removing obso-



**Software Architecture is a set of key design decisions!!!**

# The Key Issue



Takes away the knowledge!

Knowledge Vaporisation!

# Architecture Knowledge Management

*Architecture knowledge management (AKM) aims to **codify and maintain** the Architectural knowledge of a software system in a form that can be **easily accessed** by different stakeholders*



Journal of Systems and Software

Volume 116, June 2016, Pages 191-205



## 10 years of software architecture knowledge management: Practice and future

Rafael Capilla<sup>a</sup>   , Anton Jansen<sup>b</sup>  , Antony Tang<sup>c</sup>  , Paris Avgeriou<sup>d</sup>  ,  
Muhammad Ali Babar<sup>e</sup> 

Show more ▾

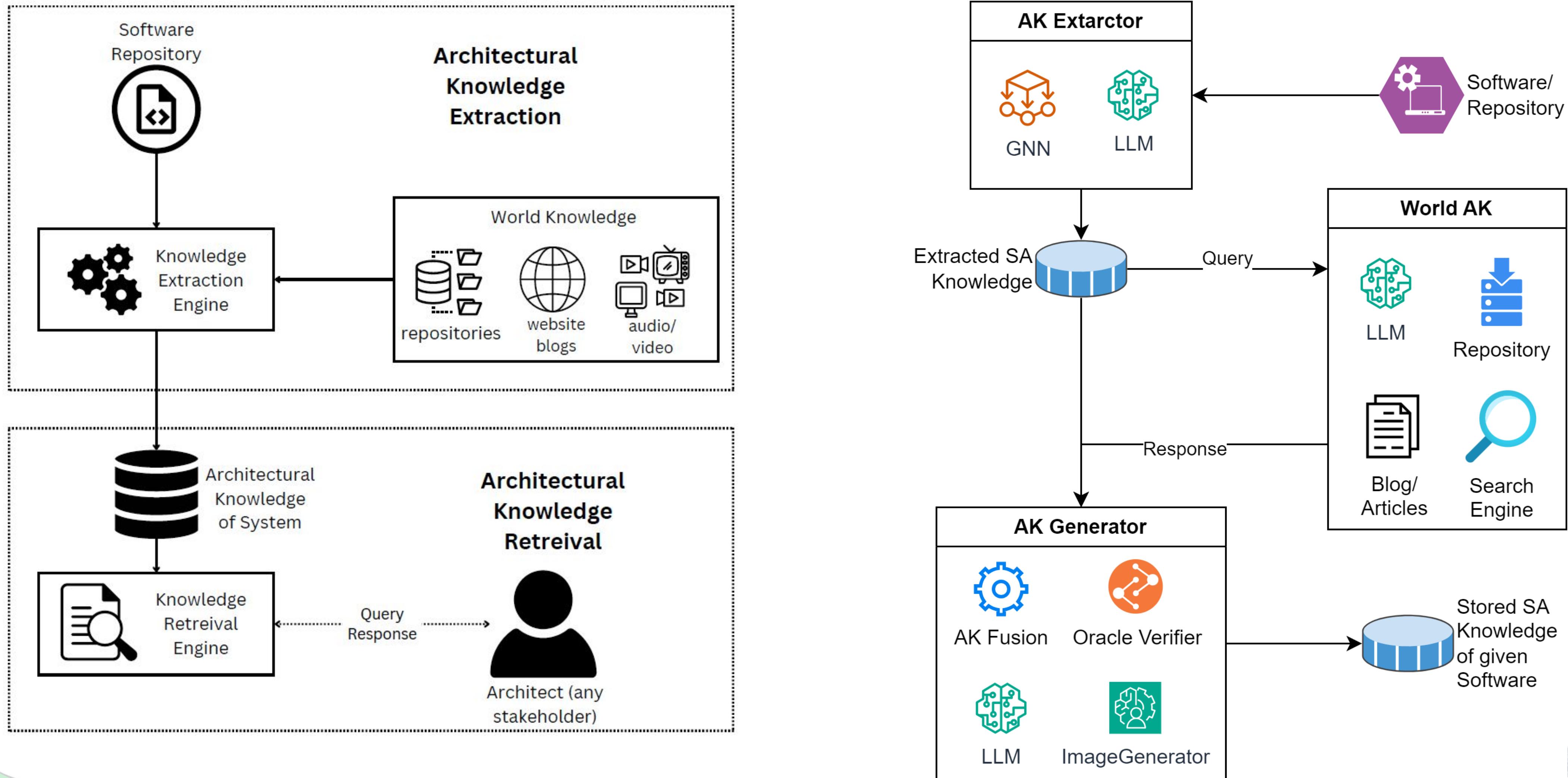
+ Add to Mendeley  Share  Cite

<https://doi.org/10.1016/j.jss.2015.08.054> ↗

[Get rights and content ↗](#)

Need for better tools => Automate using **ArchBots** or a co-pilot

# Generative AI for Architectural Knowledge Management



# Starting with Decision Records

- **Architecture Decision Records: ADR**
- Lightweight mechanism for documenting decisions
- Design decisions require careful considerations of various parameters
  - This requires broader understanding of domain as well as expertise
  - **Can we use LLMs to generate architecture design decisions?**
  - **Can LLMs be used to extract architectural information from design decisions?**

**Title:** Deciding the technology for the data analysis component

**Context**

We need to decide whether to use Python as a programming language for our project. Our project involves data analysis, machine learning, and web development.

**Decision**

We have decided to use Python as our primary programming language for our project.

**Rationale**

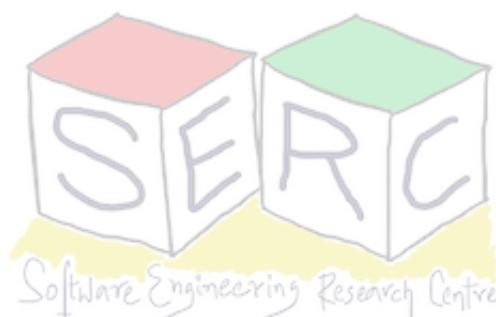
1. Python has support for various ML and data analysis
2. Team members are already familiar with Python

**Status**

Decided

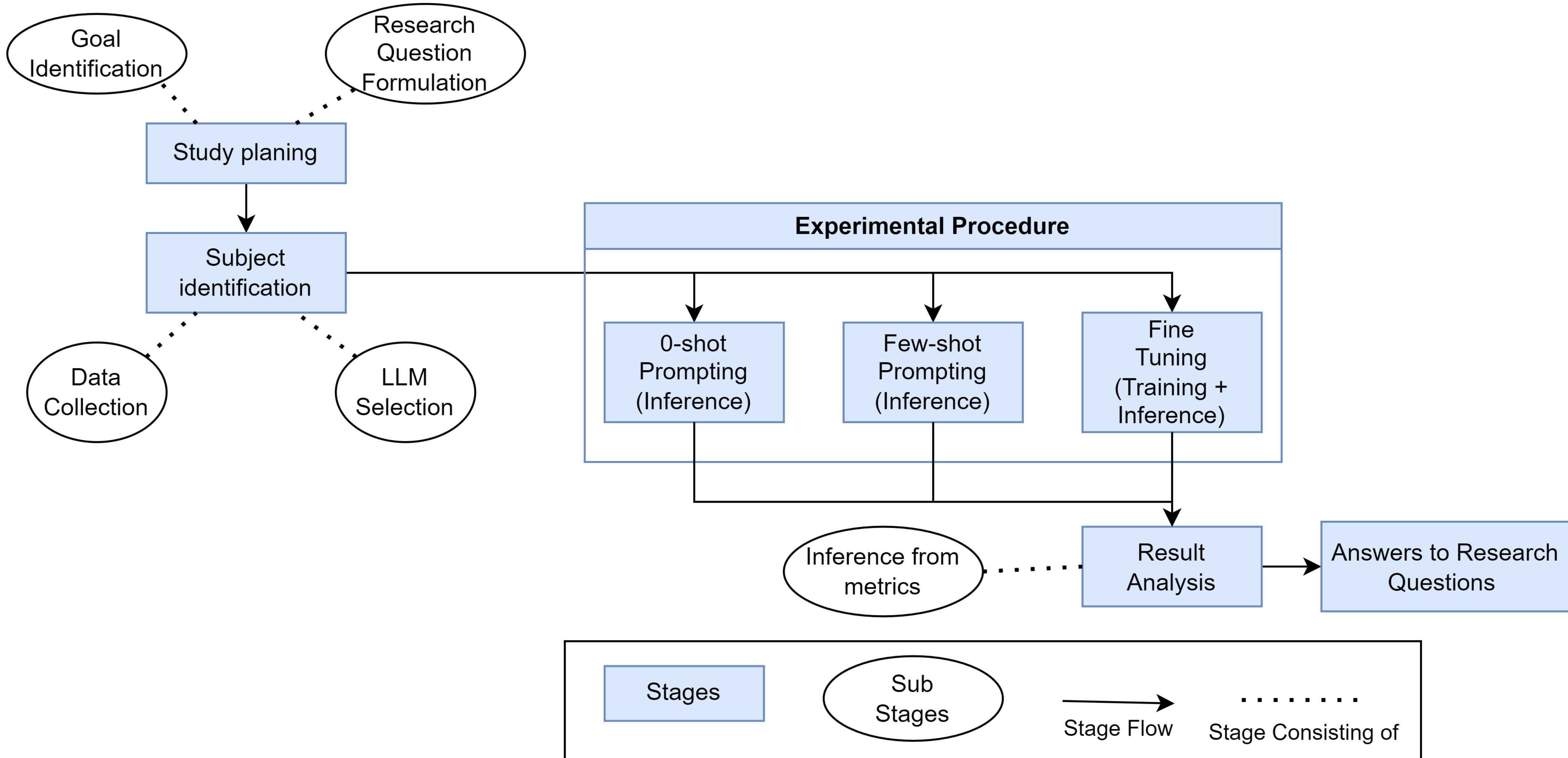
**Consequences**

1. Functionality needs to be exposed as API for integration support



<https://github.com/joelparkerhenderson/architecture-decision-record>

# Performing an Exploratory Study



# Generation Approaches

## Zero Shot Prompting - Just ask the model and expect the answer

SYSTEM

This is an Architectural Decision Record for a software. Give a ## Decision corresponding to the ## Context provided by the User

USER

#context: " We need to decide on whether to use Python as a programming language for our project. Our project involves data analysis, machine learning, and web development.

+

ASSISTANT

# Decision: Use Python as the programming language for the project.

Python is well-suited for tasks involving data analysis, machine learning, and web development due to its extensive libraries and frameworks dedicated to these areas. By leveraging the rich ecosystem of Python libraries such as NumPy, pandas, scikit-learn, TensorFlow, and Django, the development team can efficiently implement the required functionalities for the project. Additionally, Python's simplicity and readability make it a popular choice for teams working on diverse aspects of a project like data analysis, machine learning, and web development.

# Generation Approaches

## Few Shot Prompting - Give some examples in the context

**SYSTEM**  
These are architecture decision records. Follow the examples to get return Decision based on Context provided by the User.

**USER**  
## Context  
  
We need to decide on which database management system (DBMS) to use for Project X. The database will be used to store and manage large amounts of data from multiple sources. We need a DBMS that can handle transactions, offer scalability, and provide high reliability and security. Among the various options available, we are considering MySQL as a possible choice.

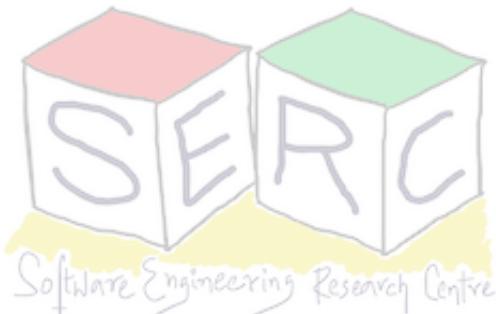
**ASSISTANT**  
## Decision  
  
After evaluating the context of our project, we have decided to use Python as the programming language.  
  
Python is a versatile language known for its simplicity and readability, making it an excellent choice for data analysis, machine learning, and web development tasks. It has a rich ecosystem of libraries and

**USER**  
## Context  
  
We need to decide on whether to use Python as a programming language for our project. Our project involves data analysis, machine learning, and web development.

**ASSISTANT**  
## Decision  
  
After evaluating the context of our project, we have decided to use Python as the programming language.

Python is a versatile language known for its simplicity and readability, making it an excellent choice for data analysis, machine learning, and web development tasks. It has a rich ecosystem of libraries and frameworks such as NumPy, Pandas, Scikit-learn for data analysis and machine learning, and Django, Flask for web development. Python's extensive community support and resources provide robust documentation, tutorials, and active user forums for problem-solving and knowledge sharing.

Give some examples



Get a response

# Generation Approaches

## Fine tuning - Fine tune base model with dataset

GPT2

Generate: {context} {decision}

Generate: {context} {decision}

Generate: {context} {decision}

T5

Read: {context}  
Generate: {decision}

Read: {context}  
Generate: {decision}

Read: {context}  
Generate: {decision}

# Selected LLMs

- **No of models:** 18

- **Model Series**

- Encoder-decoder (T5 Series)
- Decoder only (GPT series)

- **Model types**

- Chat models

- Text completion models

- From **small to large** models

family	model	size	context length	availability
GPT-2	GPT2	124M	1024	local
	GPT2-medium	335M		
	GPT2-large	774M		
	GPT2-xl	1.5B		
GPT-3	ada	350M	2048	api
	davinci	175B		
GPT-3.5	text-davinci-003	175B	4000	api
	GPT-3.5-turbo	175B		
GPT-4	GPT-4	T+	8192	api
T5	T5-small	60M	infinite	local
	T5-base	223M		
	T5-large	738M		
	T5-3b	3B		
T0	T0-3b	3B	infinite	local
Flan-T5	Flan-T5-small	77M	infinite	local
	Flan-T5-base	248M		
	Flan-T5-large	783M		
	Flan-T5-xl	3B		

# Results

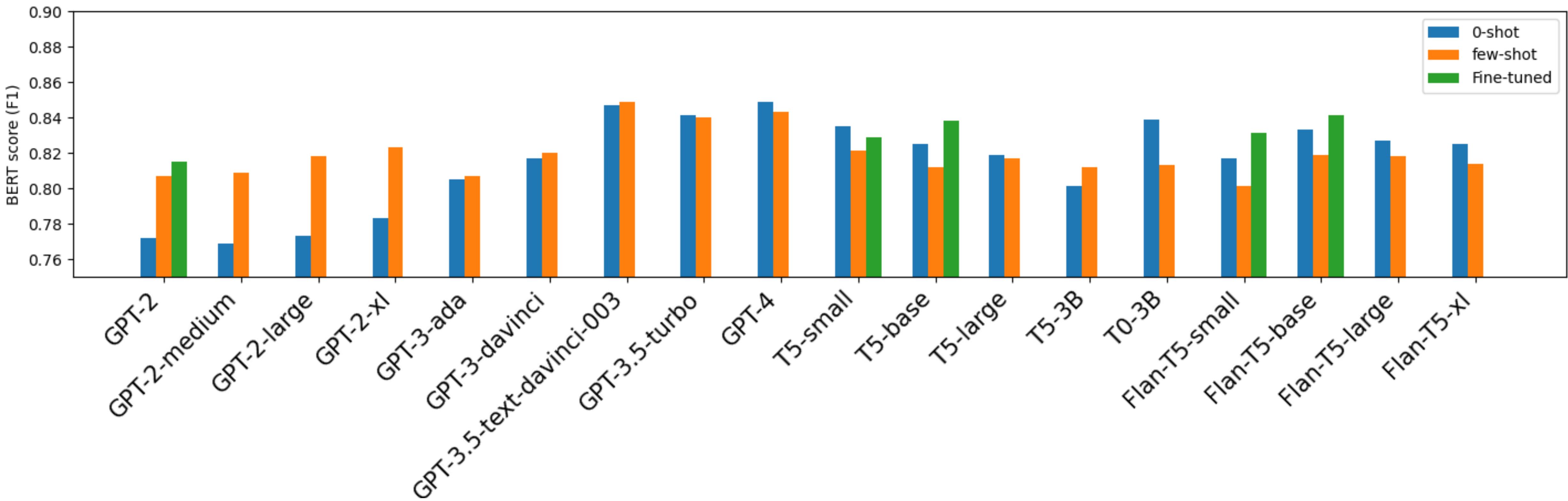
Approach	Model	Rouge-1	BLEU	Meteor	Bert-Score			Support
					Precision	Recall	F1	
0-shot	GPT3.5-text-davinci-003	0.242	0.031	0.198	0.846	0.849	0.847	91
	GPT-4	<b>0.259</b>	0.028	<b>0.219</b>	0.847	<b>0.851</b>	<b>0.849</b>	95
	T0-3b	0.187	0.005	0.122	<b>0.856</b>	0.823	0.839	91
Few-shot	GPT2-xl	0.122	0.018	0.18	0.794	<b>0.857</b>	0.823	13
	GPT3.5-text-davinci-003	<b>0.245</b>	0.028	0.207	<b>0.849</b>	0.851	<b>0.849</b>	91
	GPT-3.5-turbo	0.226	0.027	<b>0.219</b>	0.832	0.85	0.84	92
	GPT-4	0.23	<b>0.03</b>	0.21	0.836	0.851	0.843	93
Fine-Tuning	Flan-T5-base	<b>0.231</b>	<b>0.028</b>	0.171	<b>0.842</b>	<b>0.841</b>	<b>0.841</b>	18

Smaller models performing well has lot of importance in software architecture domain!



Study repo

# The Overall Results



**Smaller models perform well!...there is a need for more data and context!**



# Generated Decisions

## GPT4 zero shot

We will use Python as our programming language for this project. Python has emerged as a reliable, versatile language with a rich ecosystem of libraries and frameworks that cater to our specific needs including data analysis (Pandas, NumPy), machine learning (scikit-learn, TensorFlow), and web

## GPT3.5 few shot

After considering the various criteria discussed, we have decided to go ahead and use Python for our project.

Python is a versatile, powerful language with a wide range of features, .. community of developers who can provide support.. Python is also open-source, free to use

Python is also a great choice in terms of performance and scalability...., Python is an interpreted language..easier to debug and maintain than compiled languages.

Overall, Python is a great choice for our

It will have a large dataset, which will display a lot of data. This can help us in preparing the application.

We need to choose Python as a programming language, because all of our projects will require Python to use it.

## Flan-T5 fine-tuned

# Why not RAG?

- Software architecture is set of components and connectors.
- Software architecture is about forms, elements and rationale
  - [software, architecture, is, set, of, components, and, connectors]
  - Vector embedding: [1.5, 0.8, 0.6, 0.9. 0.4. 0.2, 0.5] - v1
  - [1.8, 0.2, 0.5, 0.9. 0.4. 0.6, 0.9] - v2
- Store vector embeddings in a Vector Data Base (VDB)
- **Retriever**
  - “What is software architecture”?: Matching documents V1 + V2
- **Generation:** Query + retrieved -> LLM -> “software architecture is about components and connectors. It is also about forms, elements and rationale”

# Introducing the ArchCopilot (Support RAG)

The screenshot shows a dark-themed application window titled "Architecture Decision Records". On the left, a sidebar lists navigation options: Home, Explore, New Chat (with a pencil icon), Choosing Java for Financial Payment Gateway (highlighted in a grey box), CSS Framework Recommendations, Project Decision Recording, Managing Architecture Decision Records, Categorizing Architecture Decisions, and Data Analysis Project Decisions.

The main content area displays a user query under the heading "USER":

I want to decide between Python and Java for building my financial payment gateway module. My organization is an Oracle customer and has a lot of java developers. The primary focus of the system should be security, Privacy and Scalability

Below it, the "ASSISTANT" section provides a response:

Based on the information provided in the documents, Java seems to be a suitable choice for building your financial payment gateway module. Java is known for its strong security features, which is crucial for handling financial transactions. Additionally, Java is widely used in enterprise environments and is supported by Oracle, which is beneficial since your organization is an Oracle customer with many Java developers.

In terms of scalability, Java is known for its ability to handle large-scale applications effectively. It is used in many high-traffic and mission-critical systems, making it a reliable choice for a payment gateway module that needs to handle a high volume of transactions.

While Python is also a popular programming language with its own strengths, based on the specific requirements of security, privacy, and scalability for a financial payment gateway module, Java seems to be the more suitable choice in this case.

At the bottom, there is a search bar with the placeholder "Enter Question" and a "Search" button with a magnifying glass icon.

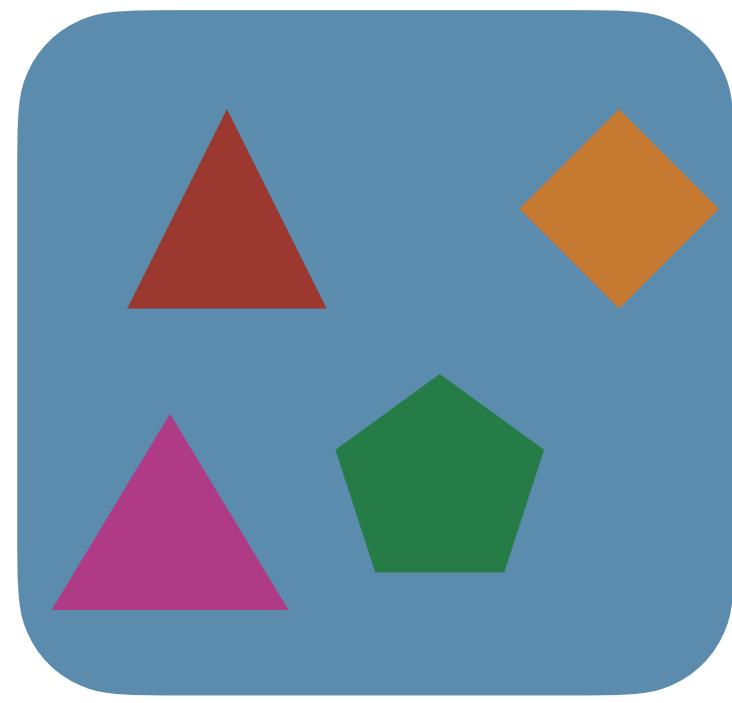
1. Add knowledge - decisions..
2. Search for similar decisions
3. Retrieve knowledge
4. Generate decision records for a given context

A companion for architect!

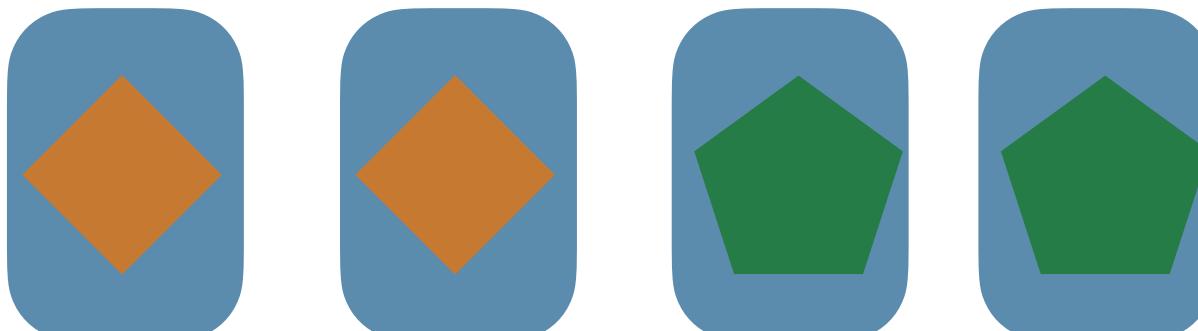
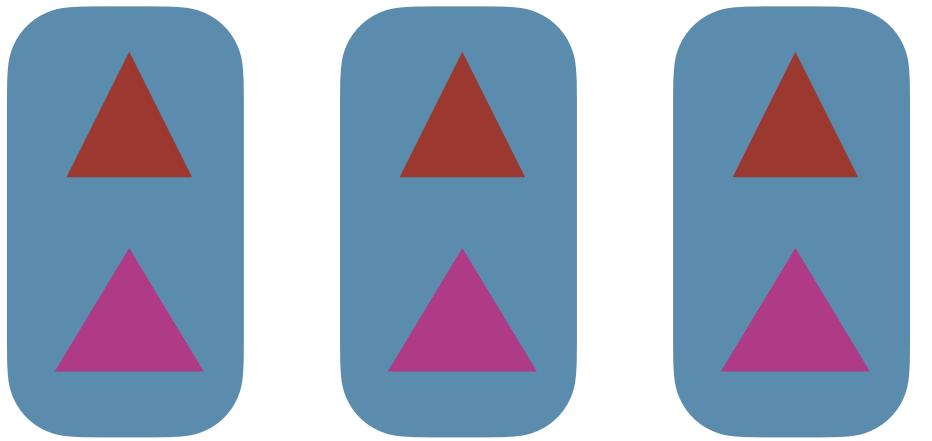
# Going towards Code Generation

Can we automate generation of architectural components?

**Monolith**



**Microservice**



**Serverless**

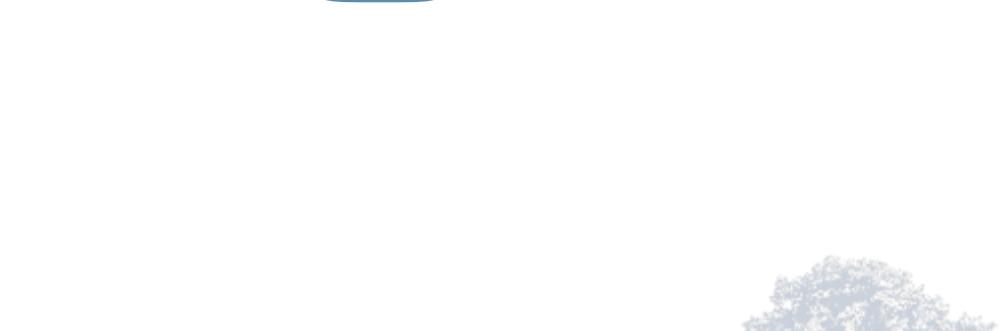
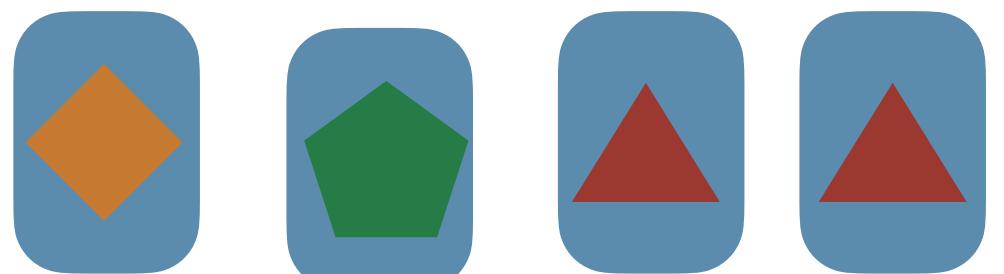
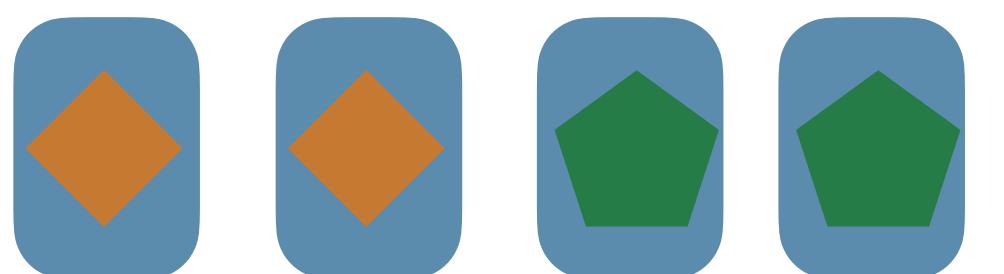
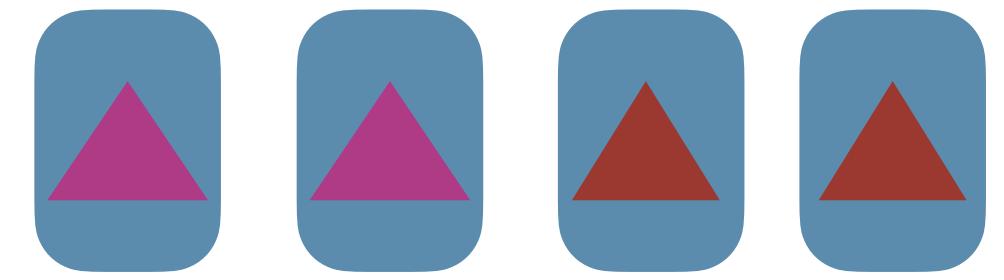


Image credits: Shrikara A



# Generating Serverless Functions!

Inspired from the concept of masked language modeling

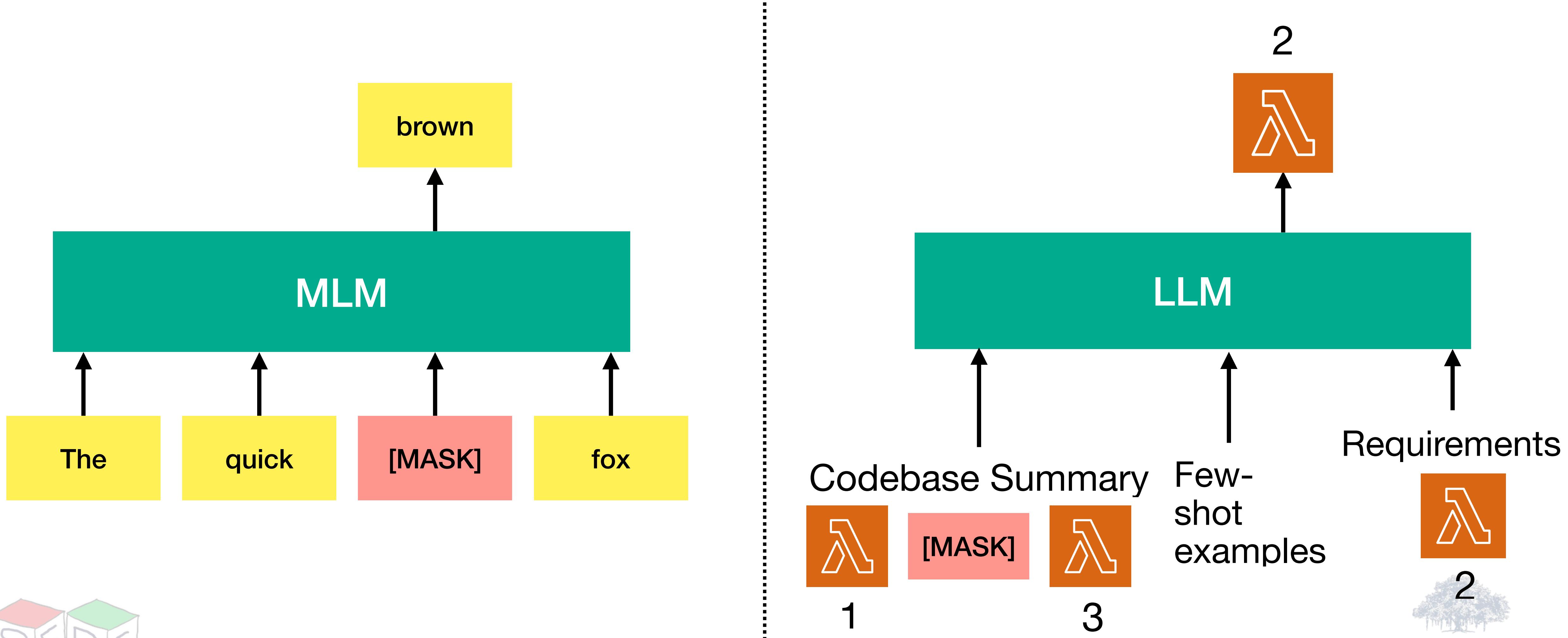
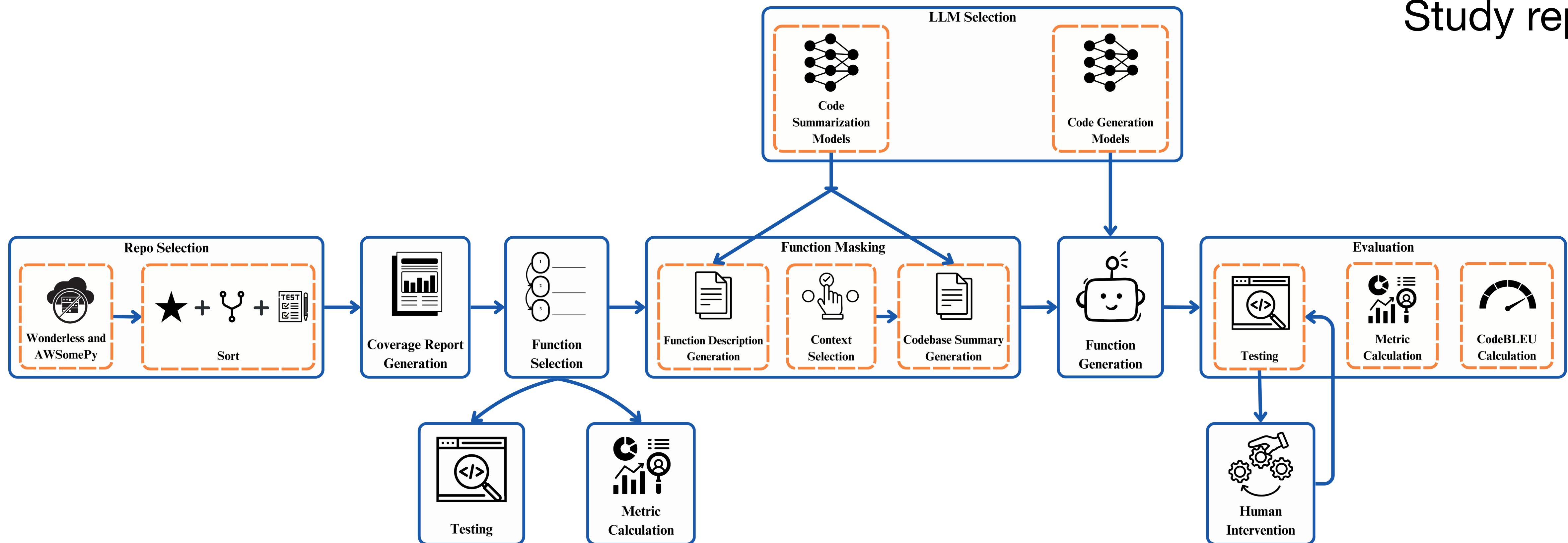


Image credits: Shrikara A

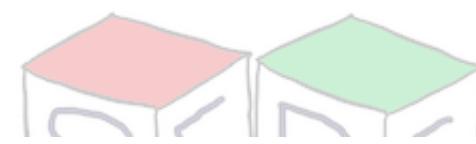


# Study Design

Study repo

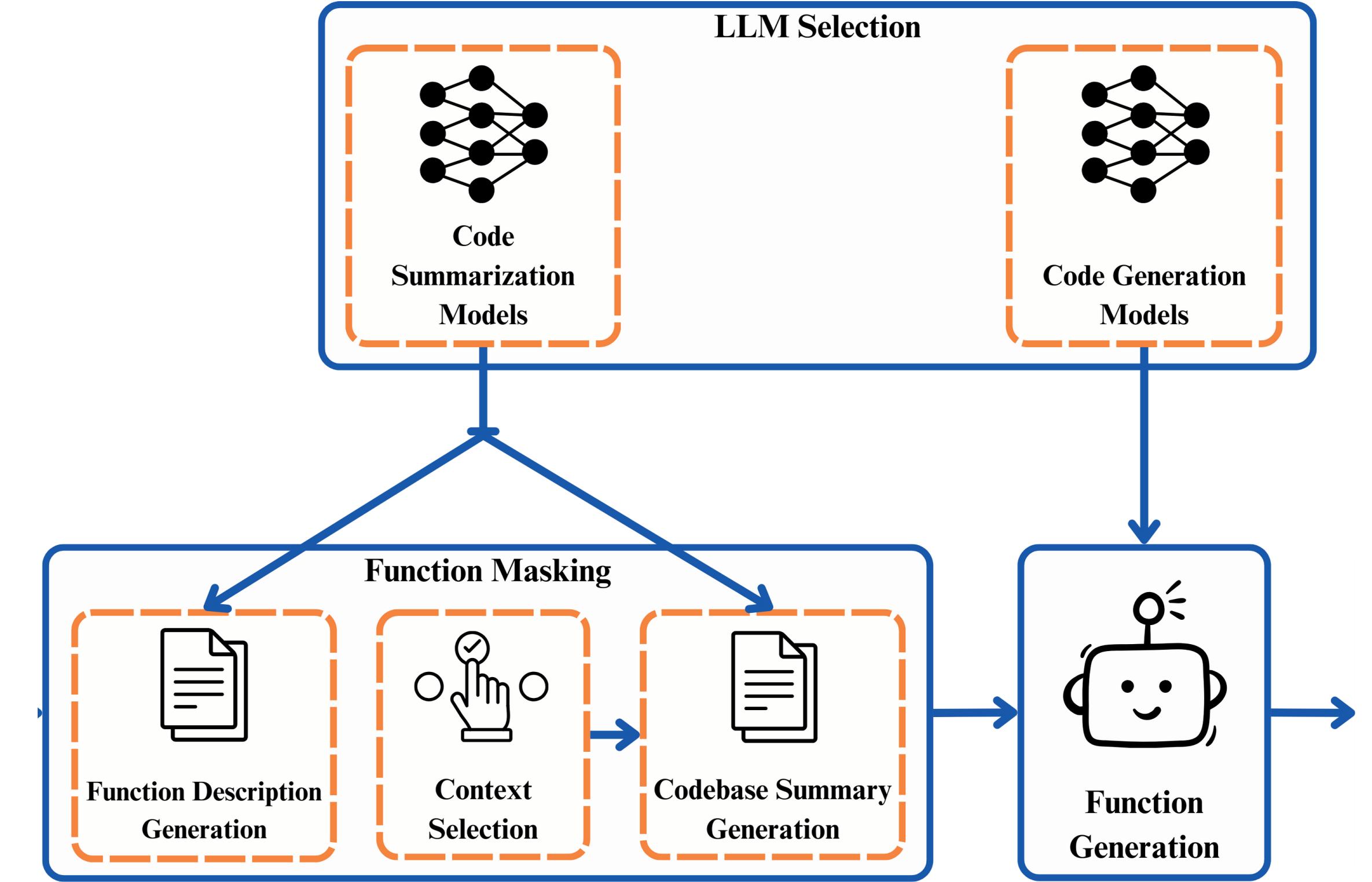


3 types of abstraction in prompt, 4 serverless repo, 5 models, 145 functions tested



# Study Design

- Choose top LLMs from LM Arena and EvalPlus leaderboards
- Kinds of prompts:
  1. Zero-shot with Readme
  2. Zero-shot with codebase summarization
  3. Few-shot with codebase summarization

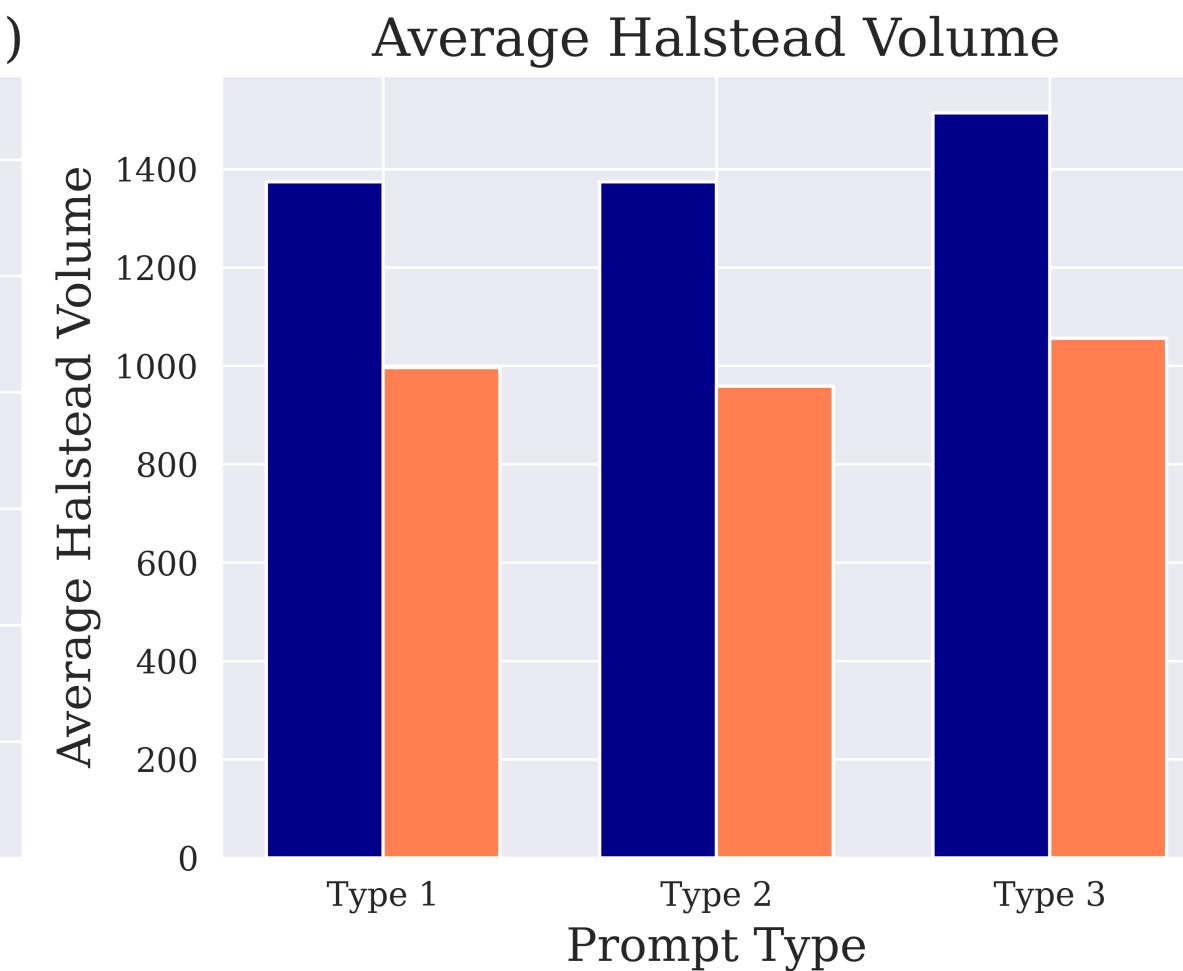
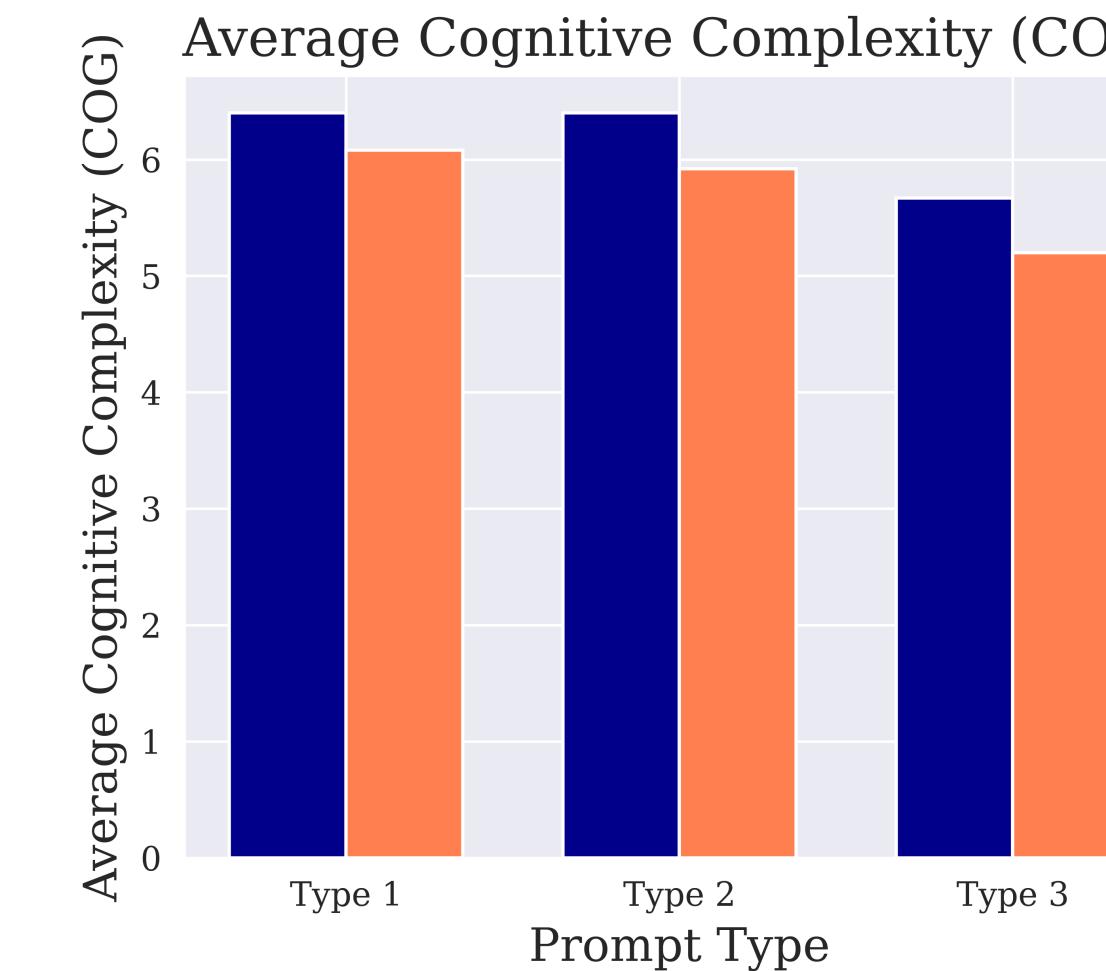
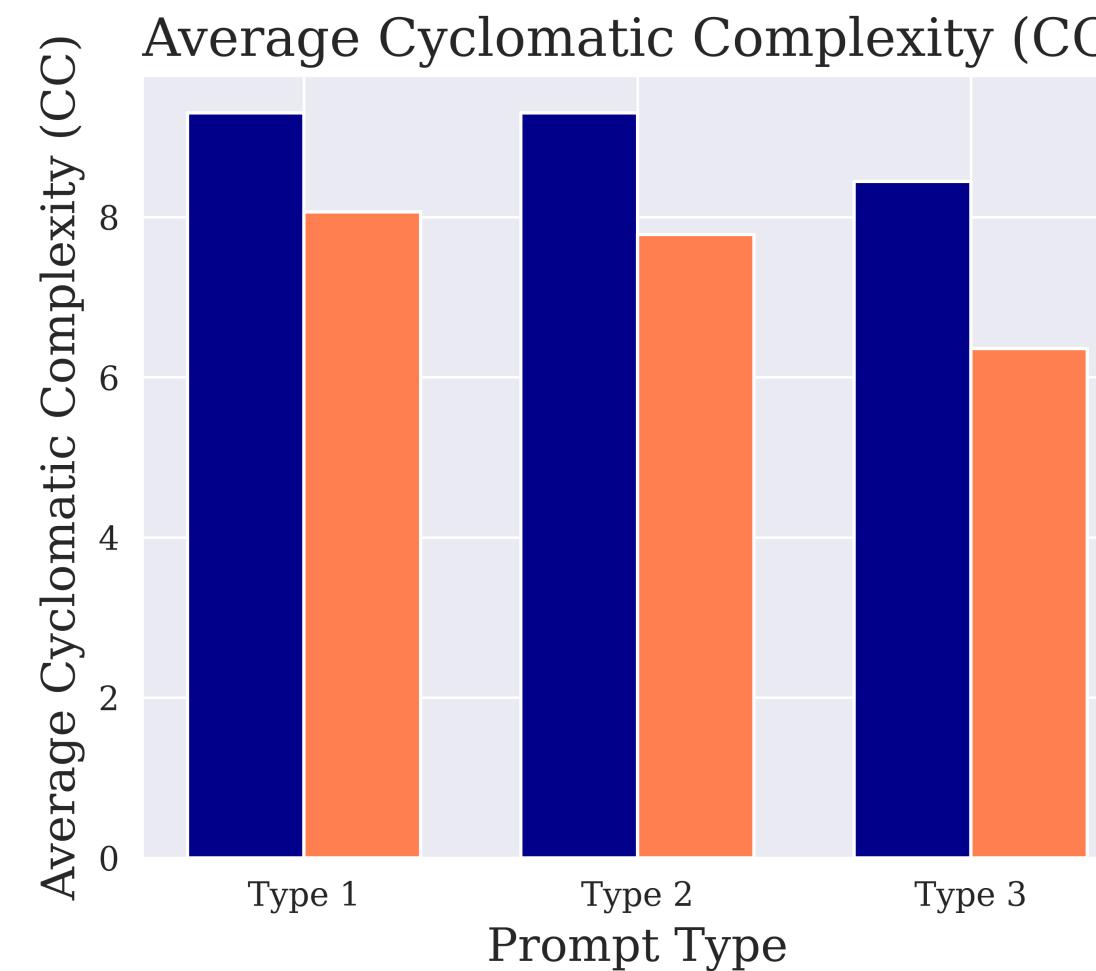
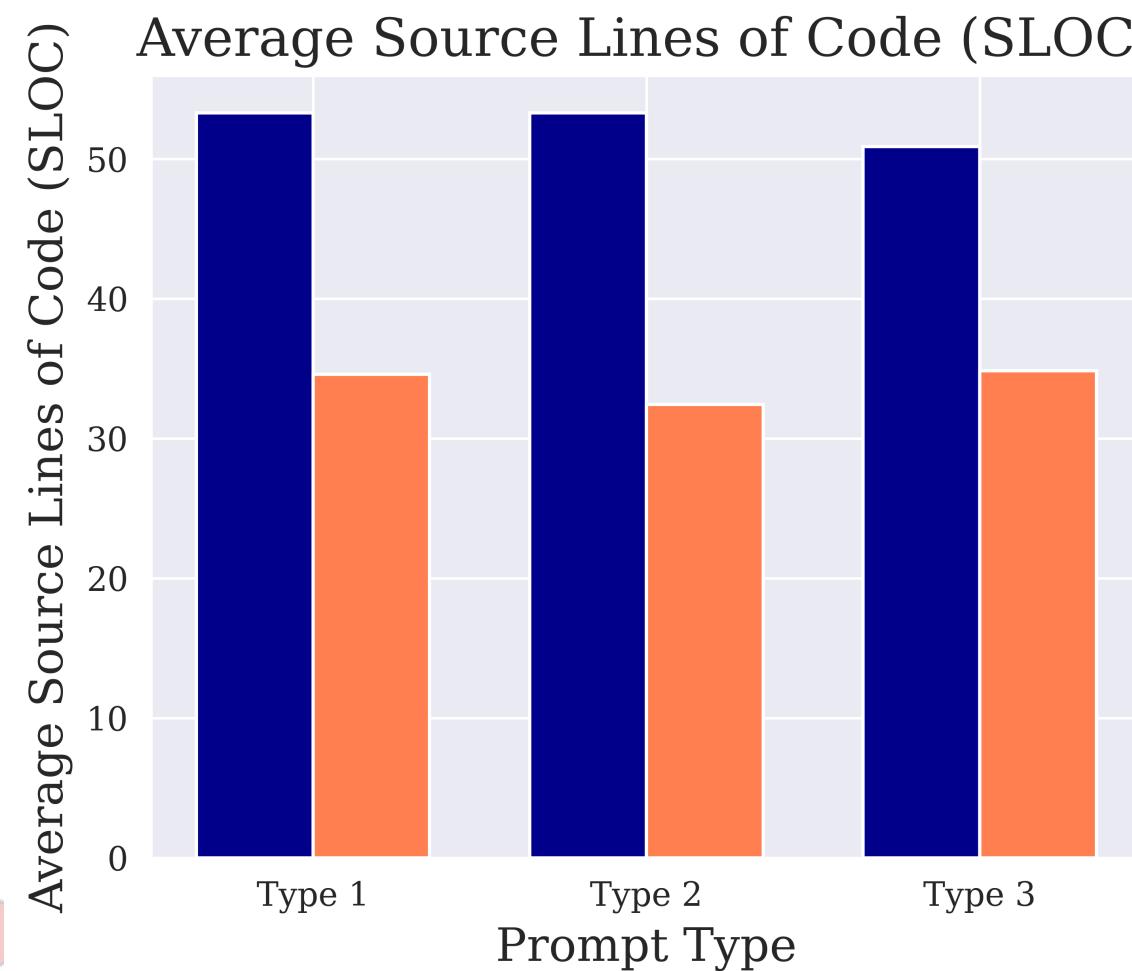
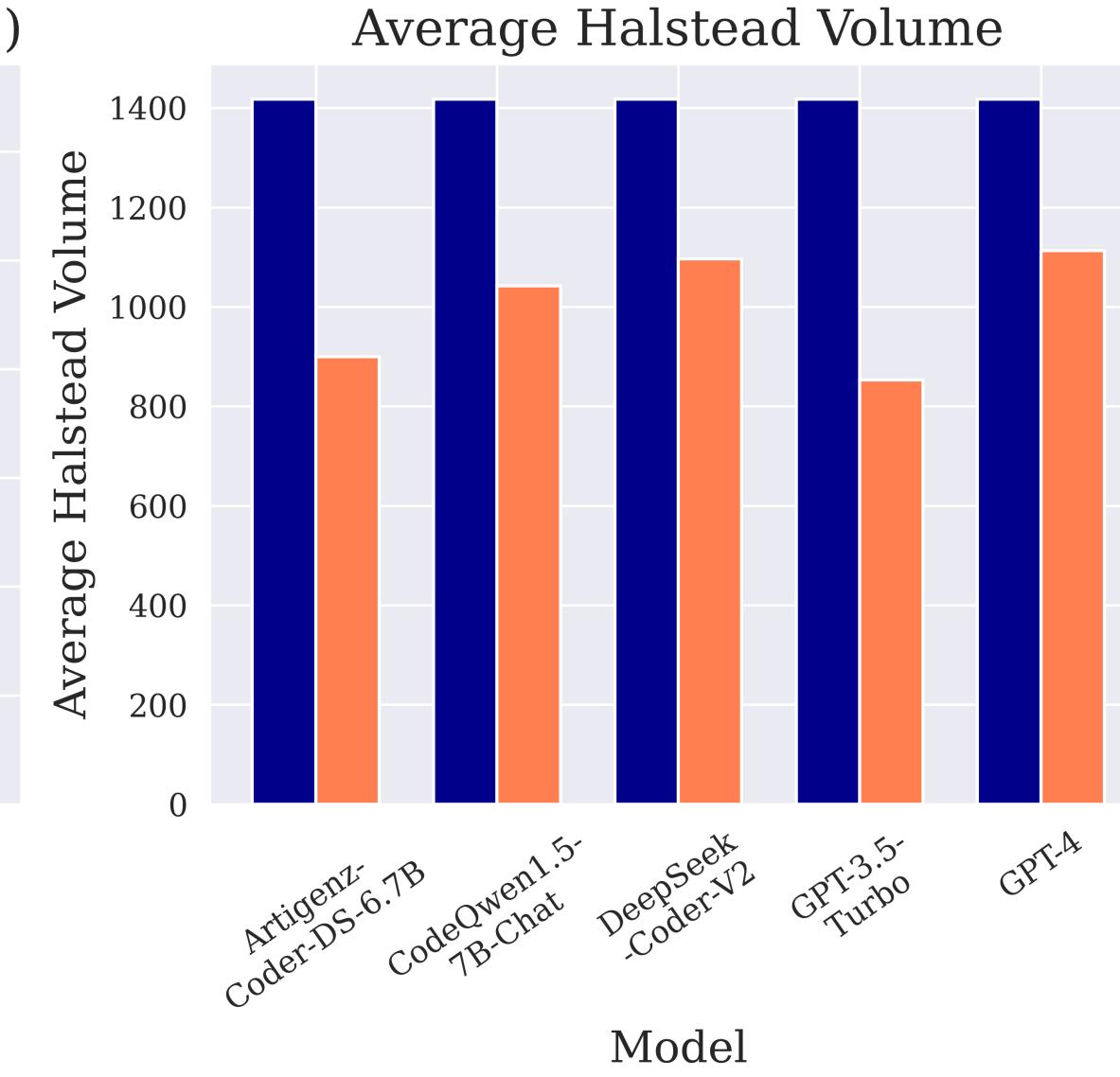
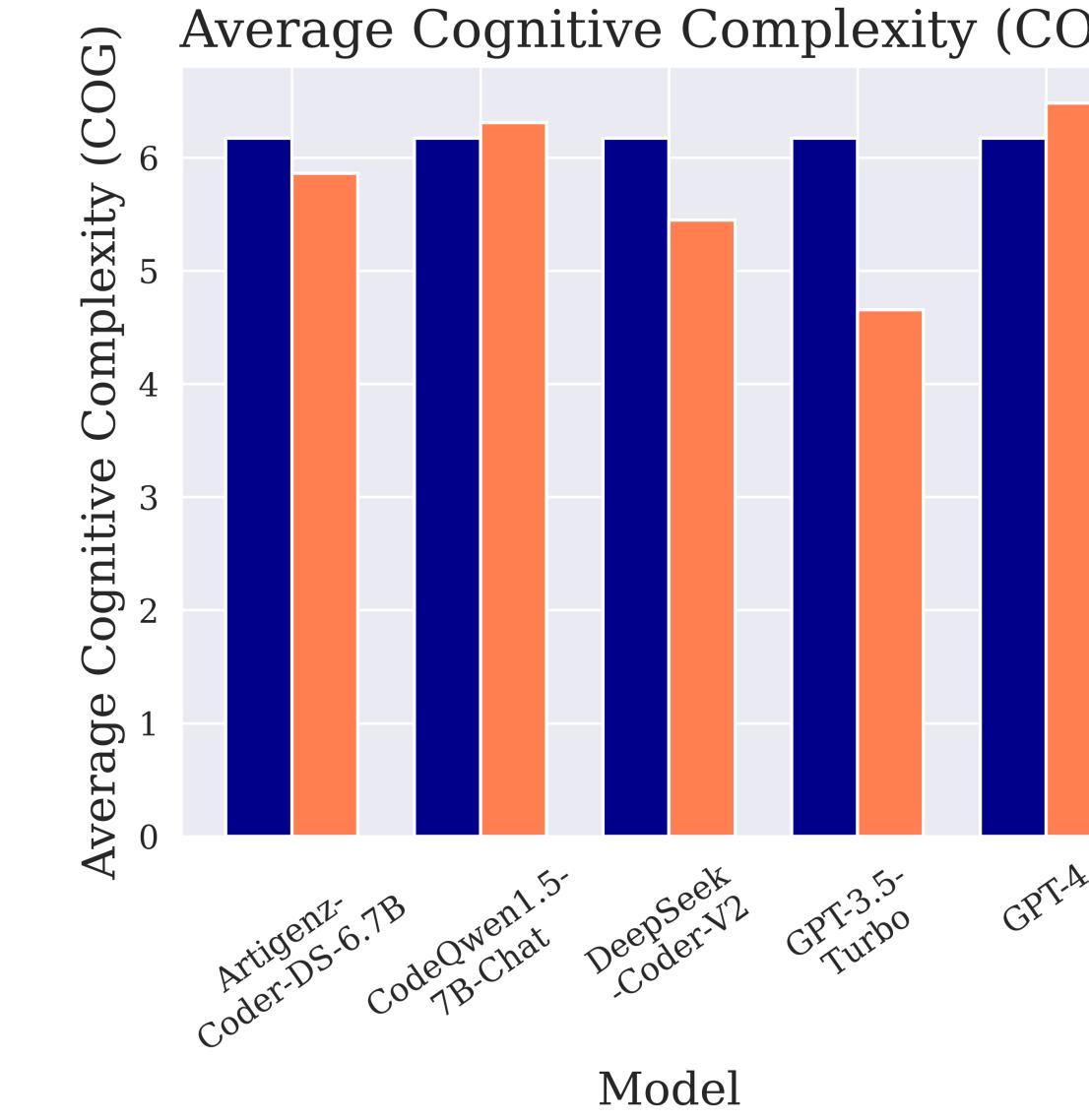
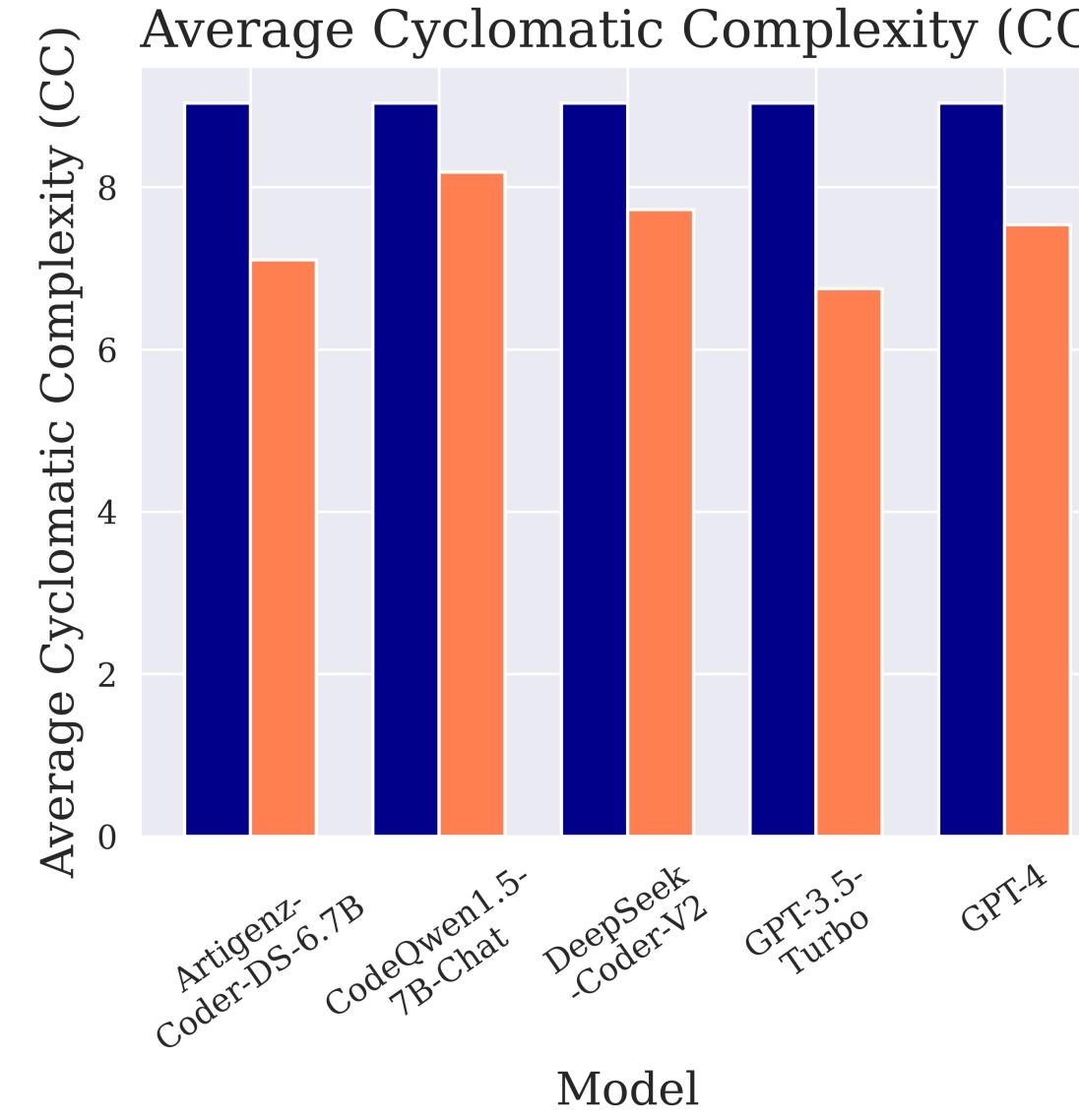
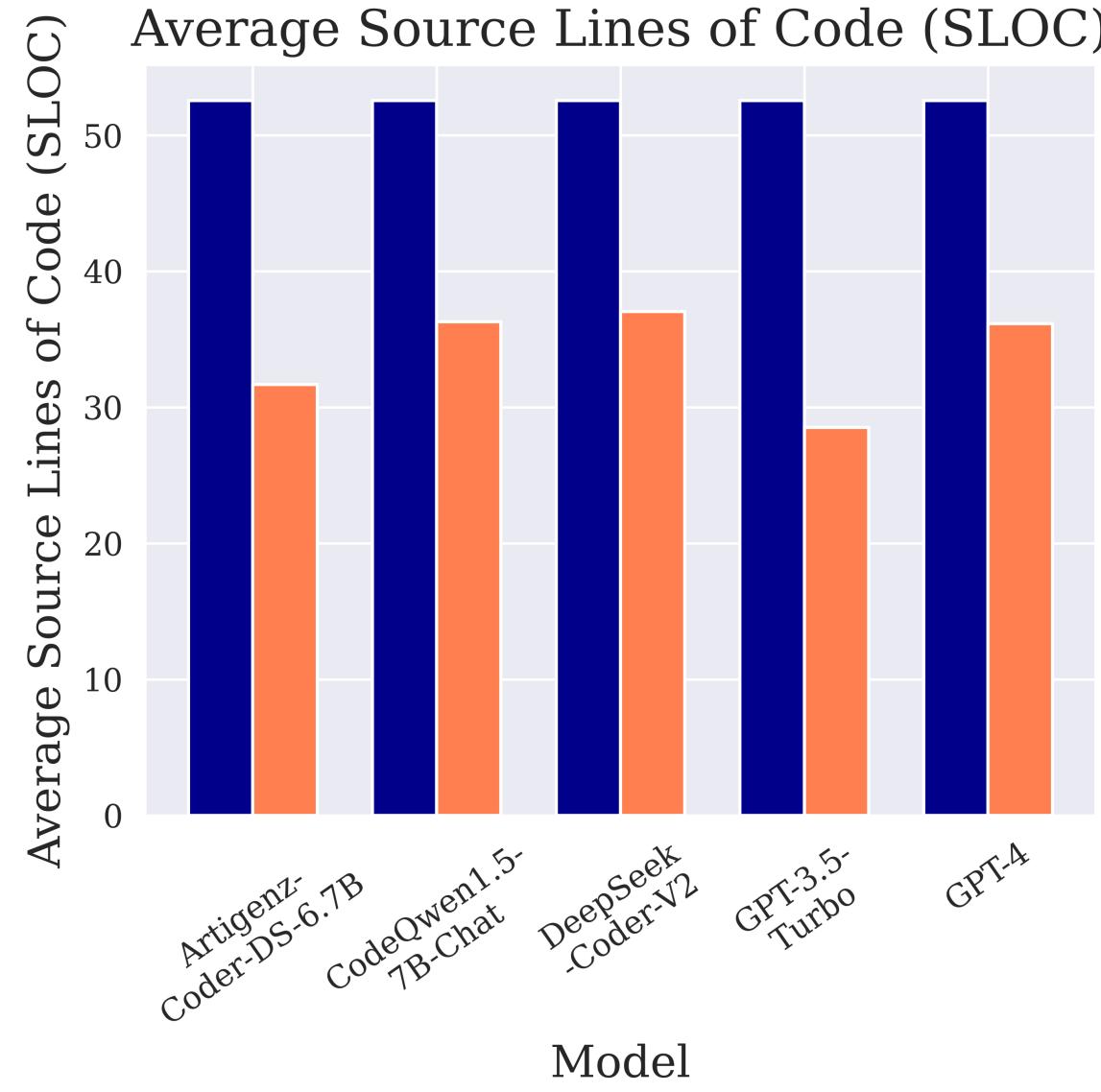


Chosen Models: Gemini 1.5 Pro for summarization, [Artigenz-Coder-DS-6.7B](#), [CodeQwen-1.5-7B](#), DeepSeek-V2, GPT-3.5-Turbo, GPT-4 for generation

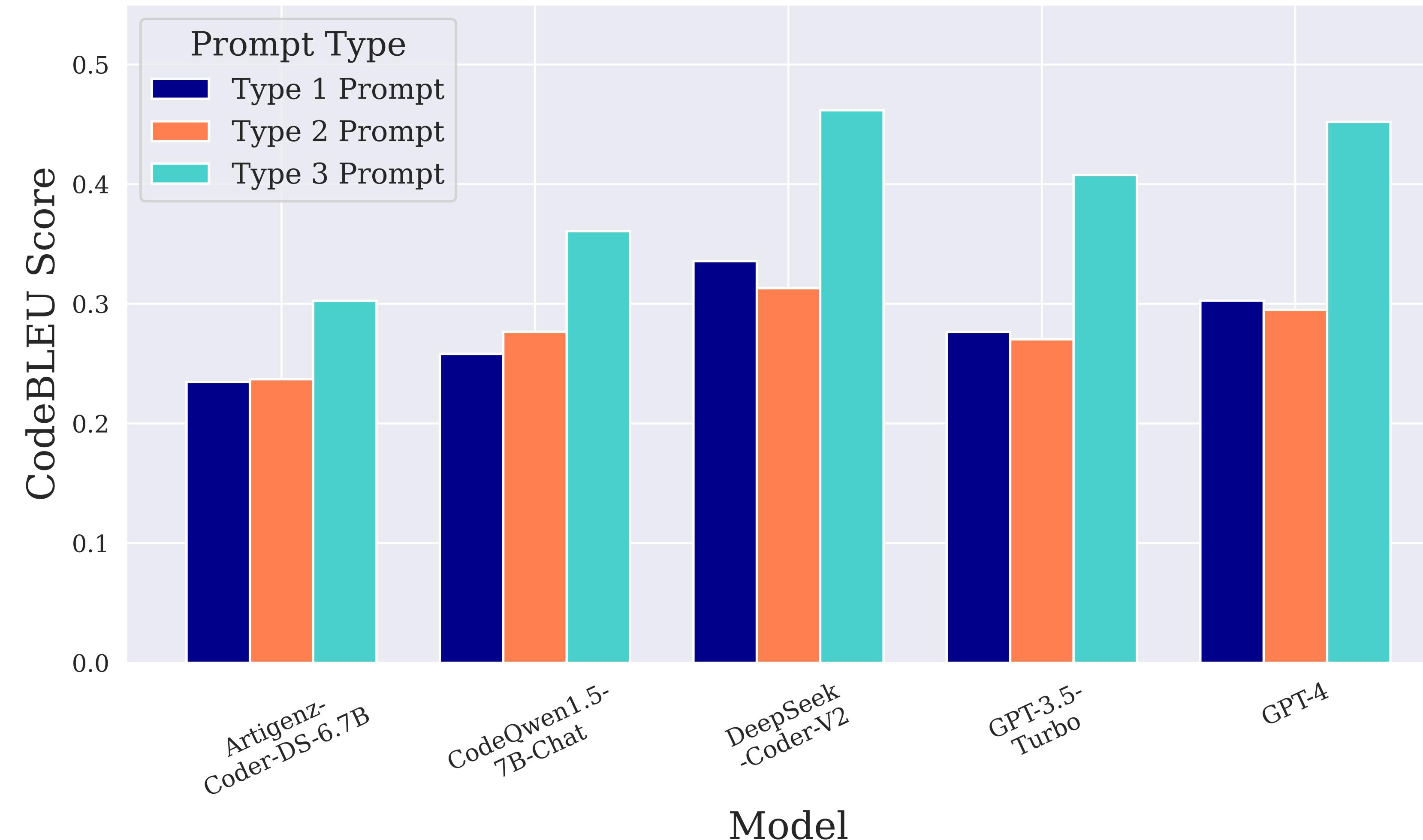
# Some Key Results

Model	Initial Test Pass Rate	Type 1 (No intervention)	Type 1 (intervention)	Type 2 (No intervention)	Type 2 (intervention)	Type 3 (No intervention)	Type 3 (intervention)
<b>Artigenz-Coder-DS-6.7B</b>	100	0	0	0	0	0	0
<b>CodeQwen1.5-7 B-Chat</b>	100	3	3	0	22	7	33
<b>DeepSeek-Coder-V2</b>	100	0	18	0	39	13	71
<b>GPT-3.5-Turbo</b>	100	0	3	0	50	4	64
<b>GPT-4</b>	100	0	24	0	61	10	50
<b>Average</b>	100	1	10	0	34	7	44

# Some Key Results - Code Quality



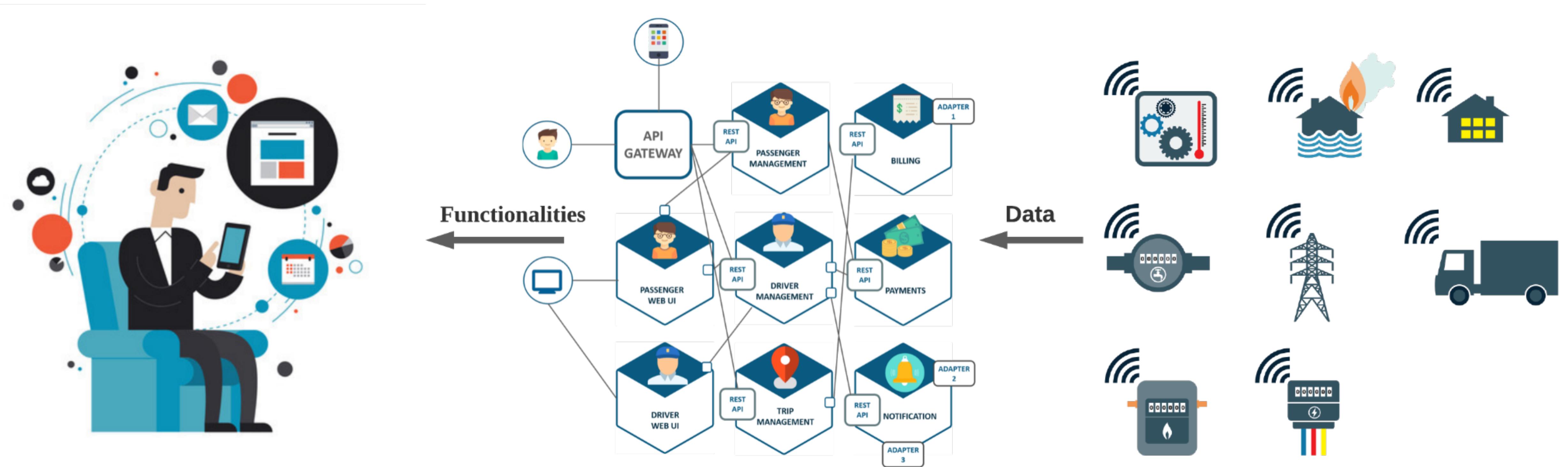
# Some Key Results - Code Quality



**Human Architects + Models + devs => Great combination!!**

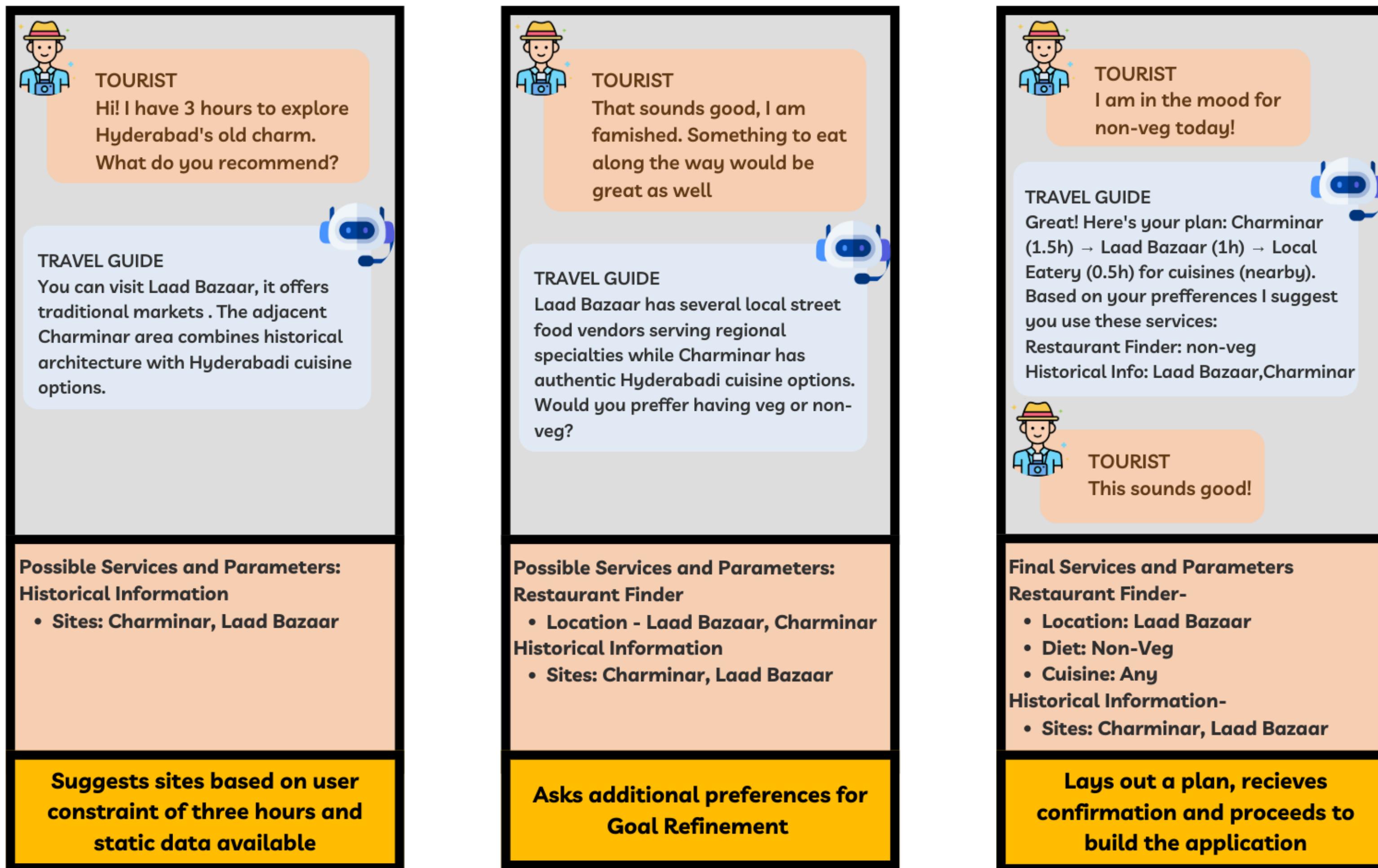
# Towards deployment: Generate + Deploy

## Applying to IoT Systems

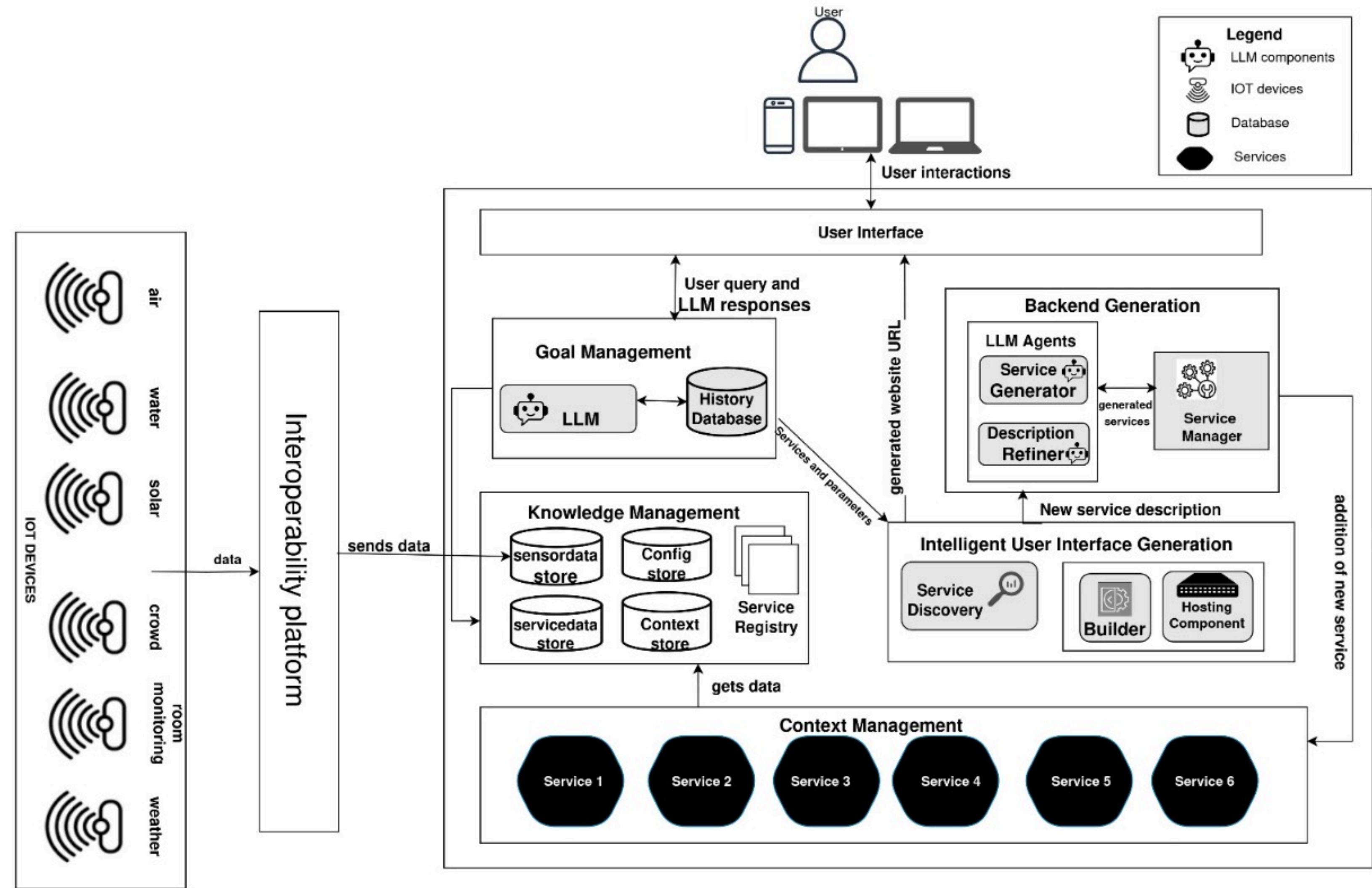


# Applying to IoT Systems

## Dynamically generate services



# IoT-Together: Mixed Initiative Interactions





Website

# Some Results

GOAL PARSER PERFORMANCE BY CATEGORY

<b>Model</b>	<b>Category</b>	<b>Precision</b>	<b>Recall</b>	<b>F1</b>	<b>Parameter Accuracy</b>
CodeQwen1.5-7B	Ambiguous	0.450	0.806	0.553	0.116
	Concrete	0.206	0.609	0.288	0.051
	<b>Overall</b>	0.282	0.670	0.370	0.071
GPT-4o-mini	Ambiguous	0.683	0.795	0.730	0.549
	Concrete	0.467	0.773	0.559	0.739
	<b>Overall</b>	0.523	0.778	0.603	0.690
DeepSeek-V2.5	Ambiguous	0.681	0.788	0.725	0.585
	Concrete	0.492	0.830	0.591	0.743
	<b>Overall</b>	0.554	0.816	0.635	0.691

USER SATISFACTION METRICS

<b>Metric</b>	<b>Average Rating (out of 5)</b>
Application Rating	4.0
Accuracy Rating	4.1
Relevance Rating	4.2

Preliminary study with 15 users

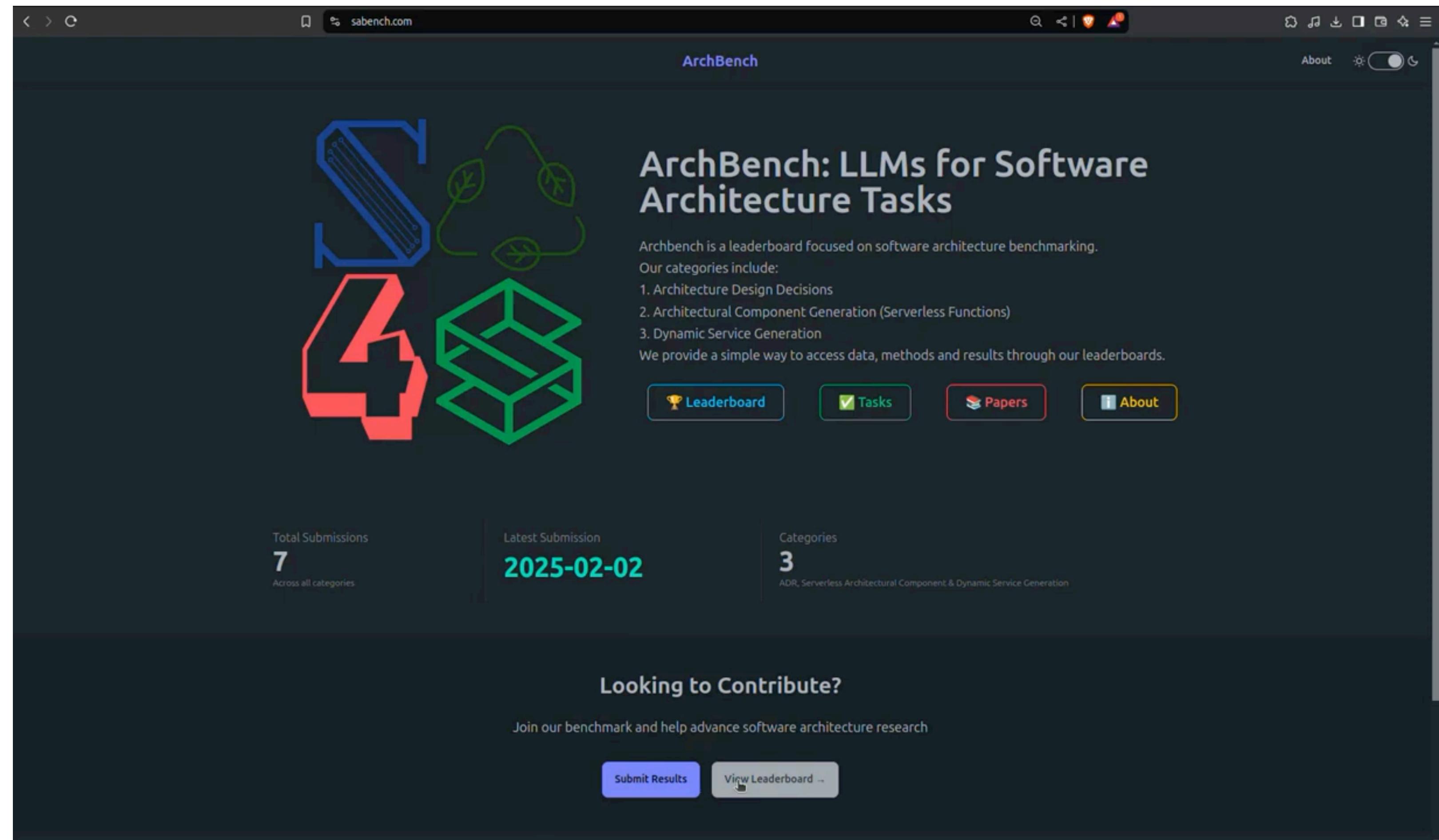
APPLICATION GENERATION PERFORMANCE METRICS

<b>Metric</b>	<b>Mean ± SD</b>	<b>Min</b>	<b>Max</b>
Total Duration (s)	$23.10 \pm 6.47$	13.46	33.08
Total Token Usage	$8164.90 \pm 2718.89$	5531	13991
Build Time (ms)	$4.85 \pm 1.98$	3.50	10.49

# Introducing sabench.com

## LLMs for Software Architecture Tasks

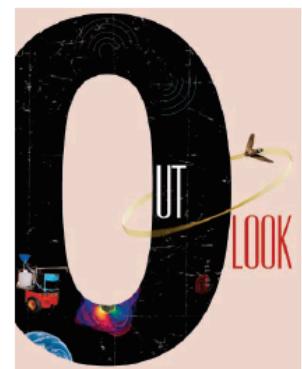
- Collection of three tasks
  - ADR
  - Serverless function generation
  - Dynamic service generation
- Datasets and leaderboards
- You can contribute as well!
- **Feedbacks welcome!**



# Extending Beyond Deployment: To Maintenance!

## What if systems could adapt like human cells?

### The Vision of Autonomic Computing



COVER FEATURE

Systems manage themselves according to an administrator's goals. New components integrate as effortlessly as a new cell establishes itself in the human body. These ideas are not science fiction, but elements of the grand challenge to create self-managing computing systems.

Jeffrey O.  
Kephart

David M.  
Chess

IBM Thomas J.  
Watson Research  
Center

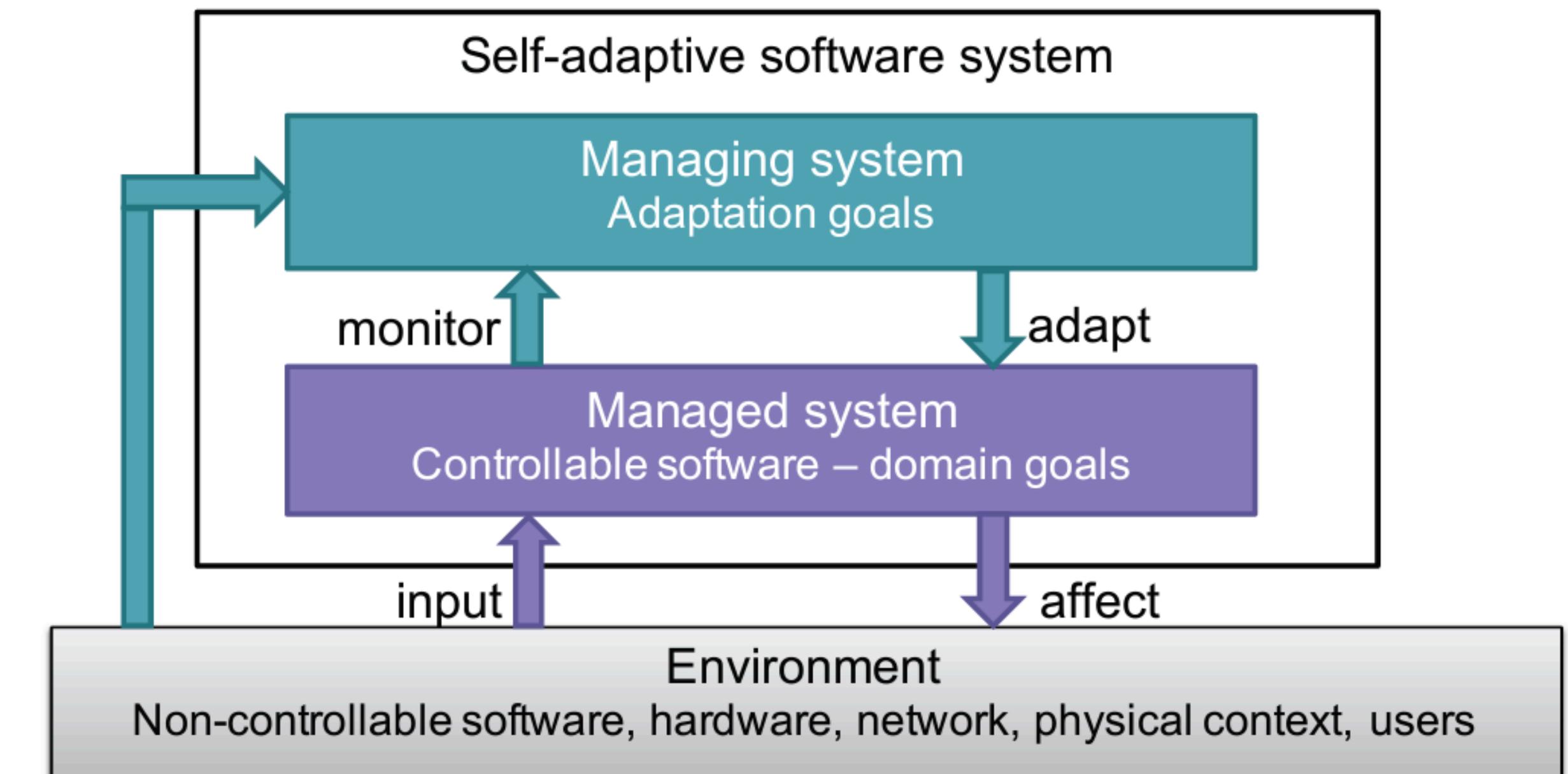
In mid-October 2001, IBM released a manifesto observing that the main obstacle to further progress in the IT industry is a looming software complexity crisis.<sup>1</sup> The company cited applications and environments that weigh in at tens of millions of lines of code and require skilled IT professionals to install, configure, tune, and maintain.

The manifesto pointed out that the difficulty of managing today's computing systems goes well beyond the administration of individual software environments. The need to integrate several heterogeneous environments into corporate-wide computing systems, and to extend that beyond company

figure, optimize, maintain, and merge. And there will be no way to make timely, decisive responses to the rapid stream of changing and conflicting demands.

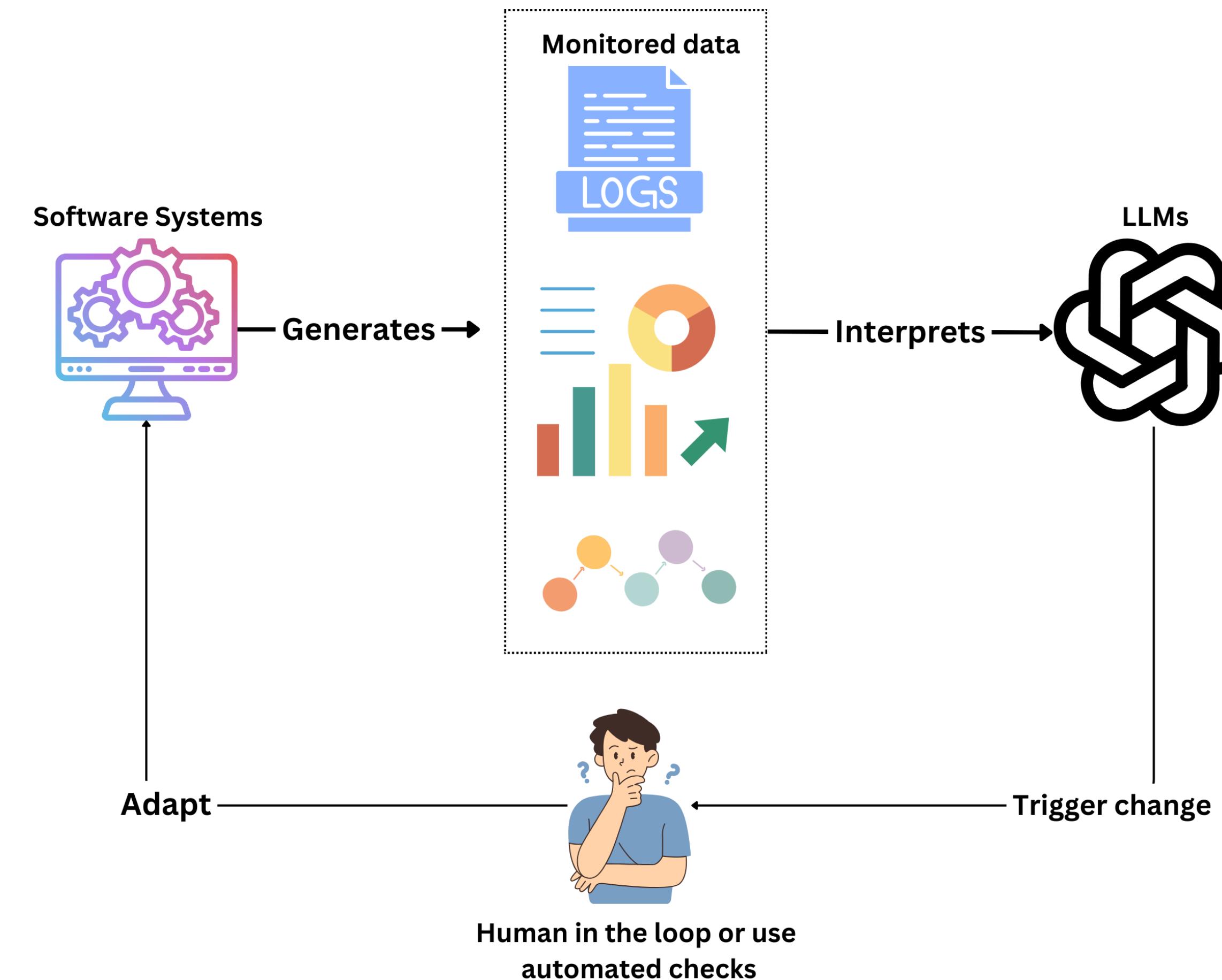
#### AUTONOMIC OPTION

The only option remaining is *autonomic computing*—computing systems that can manage themselves given high-level objectives from administrators. When IBM's senior vice president of research, Paul Horn, introduced this idea to the National Academy of Engineers at Harvard University in a March 2001 keynote address, he deliberately chose a term with a biological conno-

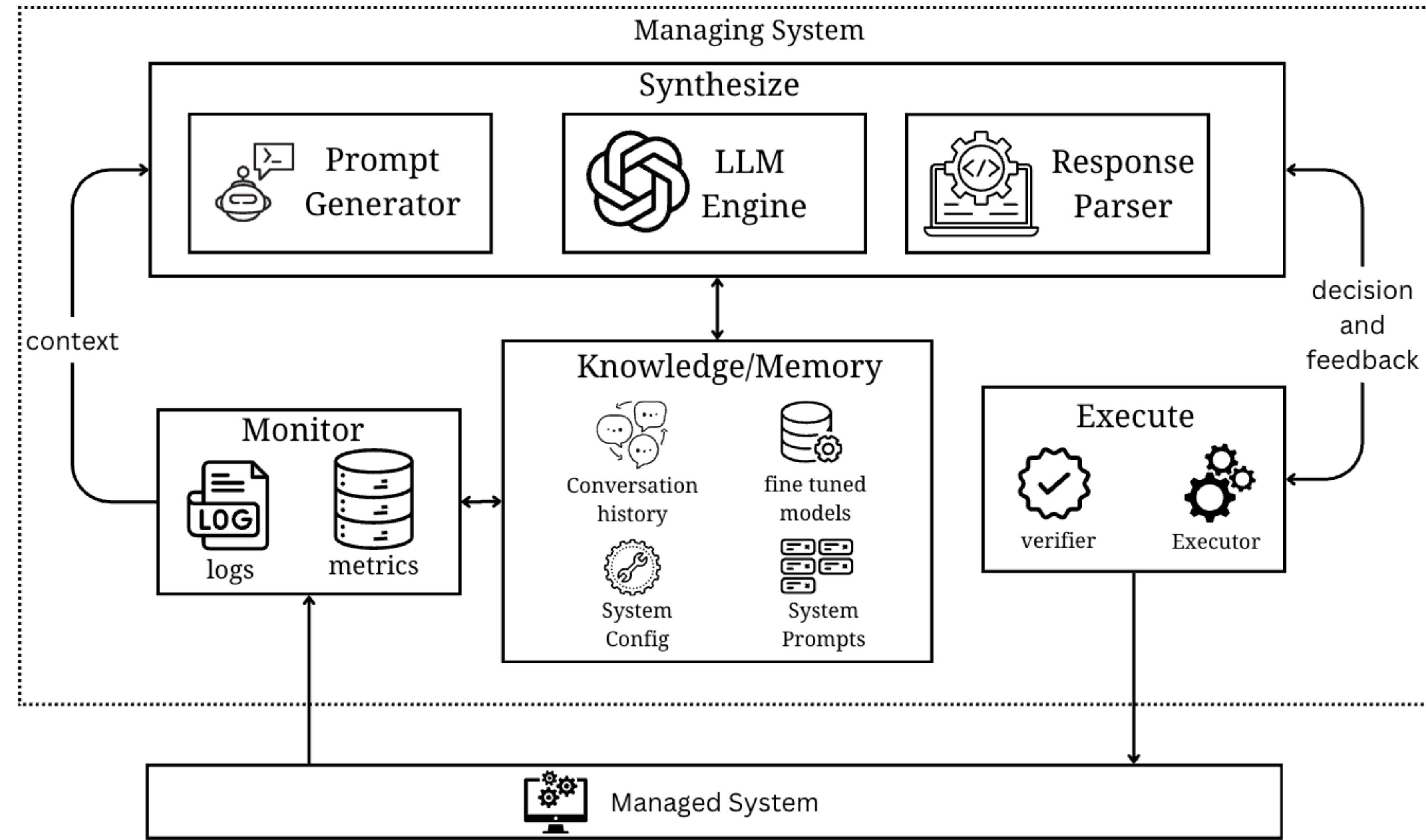


# Extending beyond

## Design time to run-time adaptation - Can LLMs help?

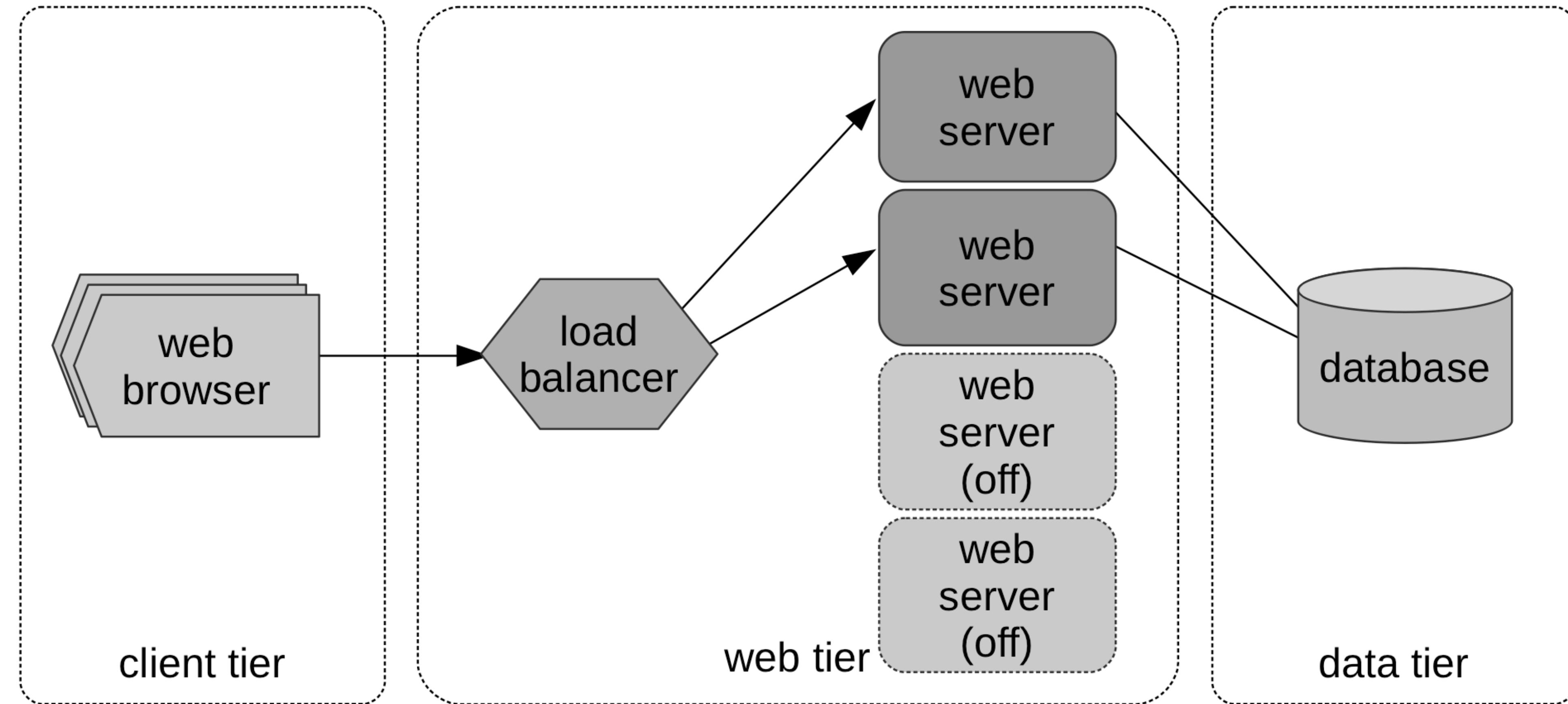


# Reimagining Self-adaptation loop



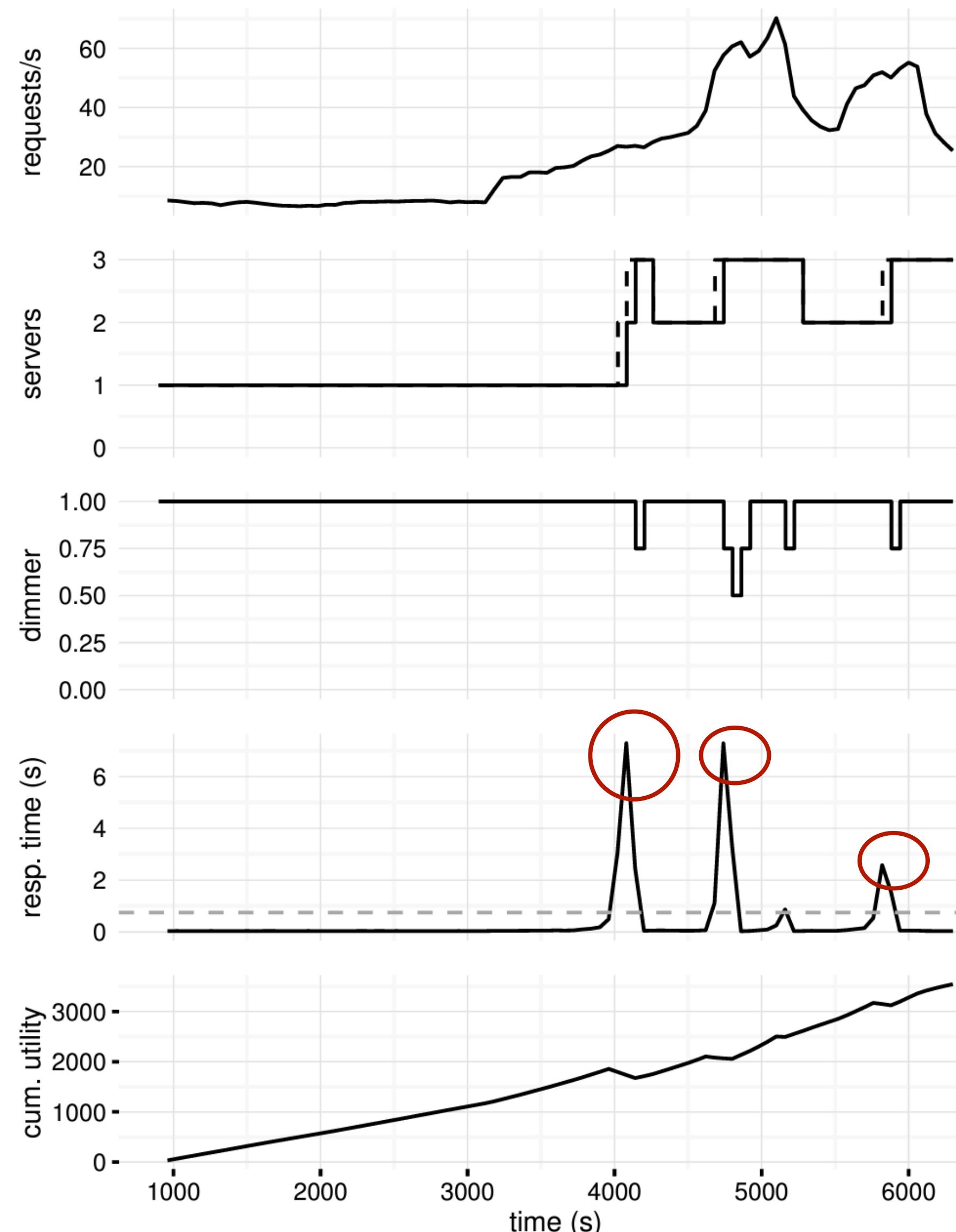
# SWIM case study

## Web Infrastructure Simulator

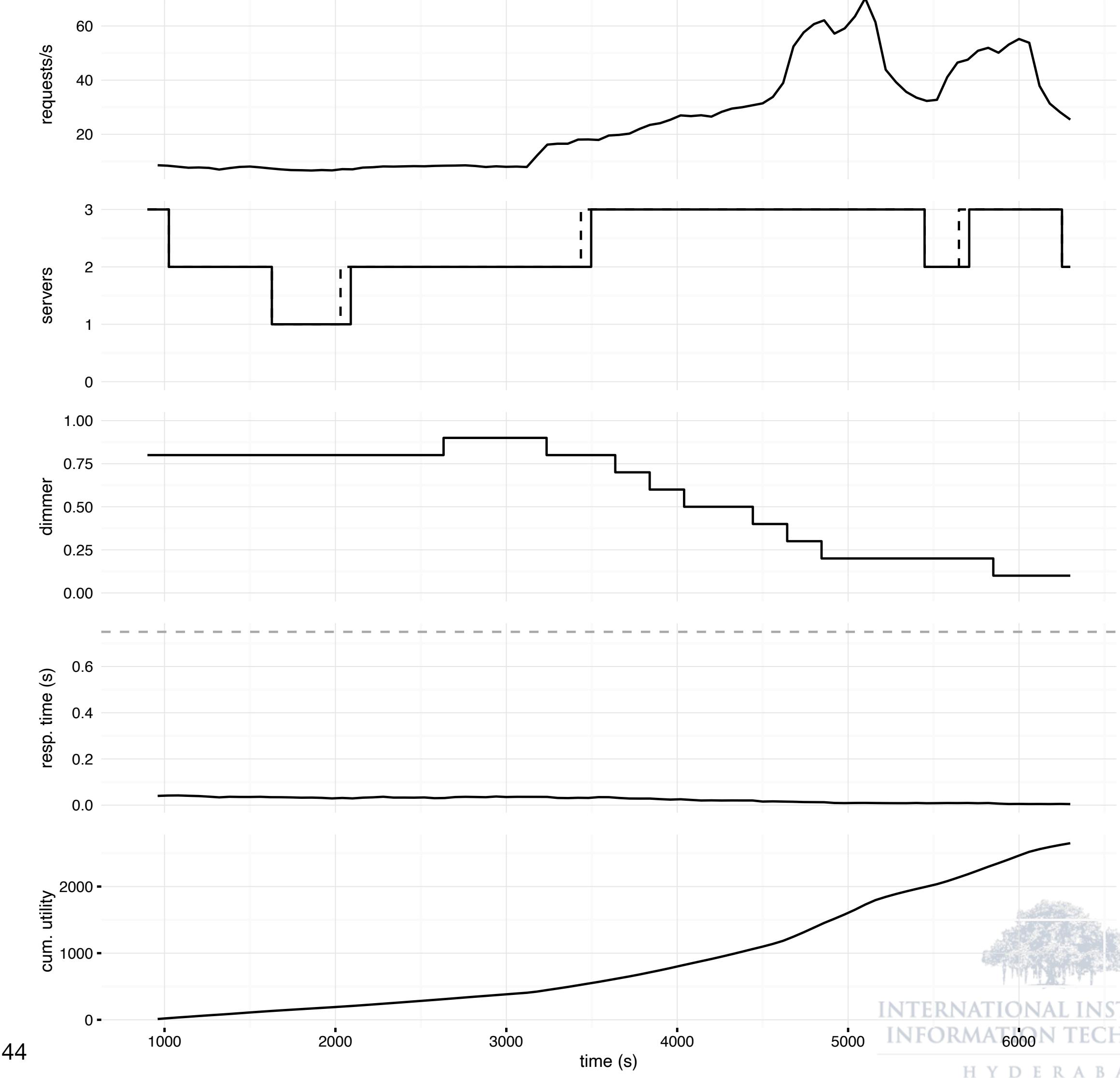


# Some Initial Results

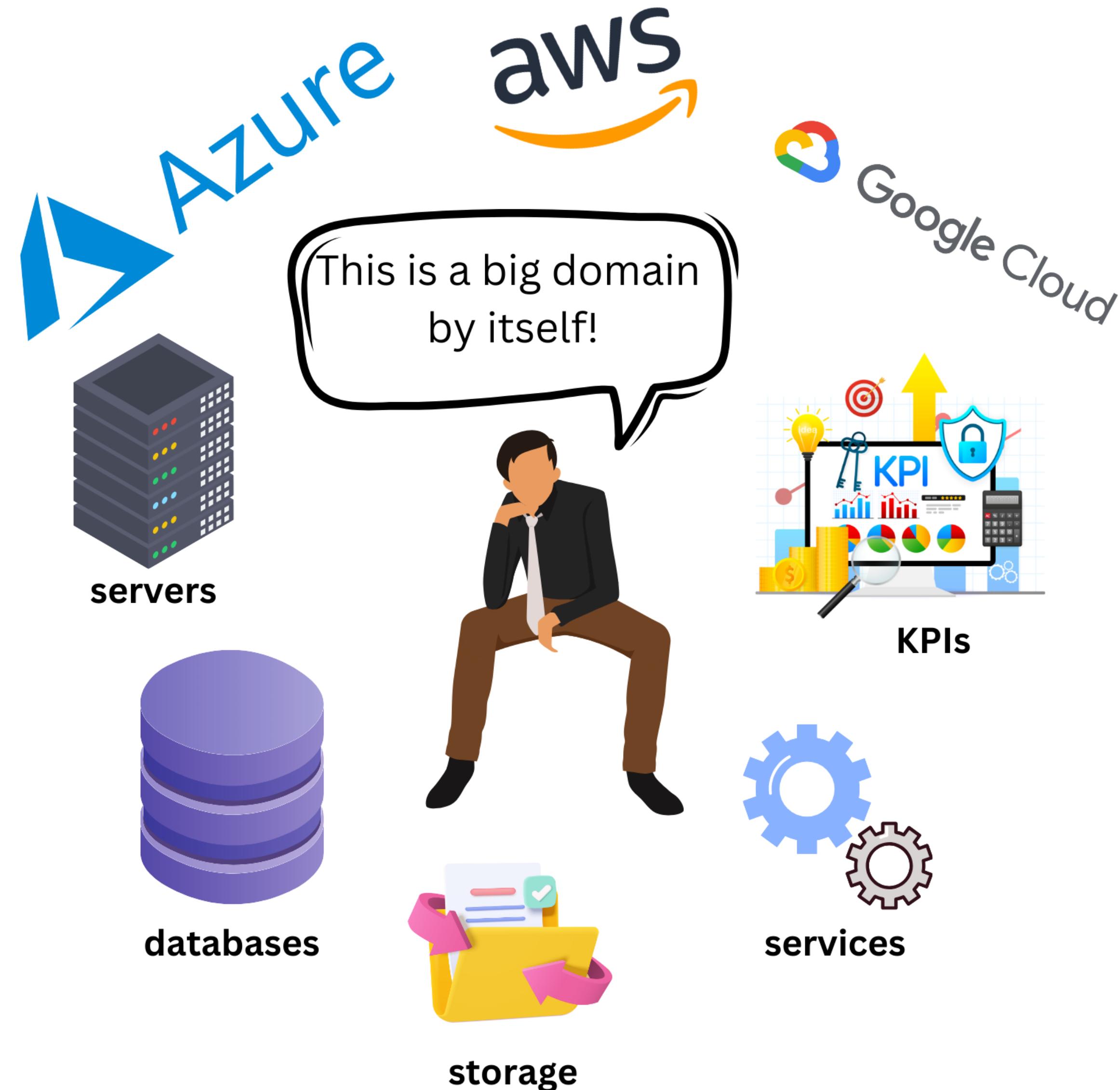
Using SWIM reactive adaptation



Using GPT-4, Promising but..



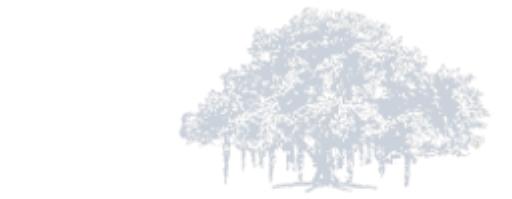
# Making it more Concrete - CloudOps domain



## AWS Well Architected Framework

Helps cloud architects build resilient, secure and high performing infrastructure

- **Build around six pillars**
  - Operational Efficiency
  - Security
  - Reliability
  - Performance Efficiency
  - Sustainability
  - Cost



# Ideas into Production: CloudOps CoPilot

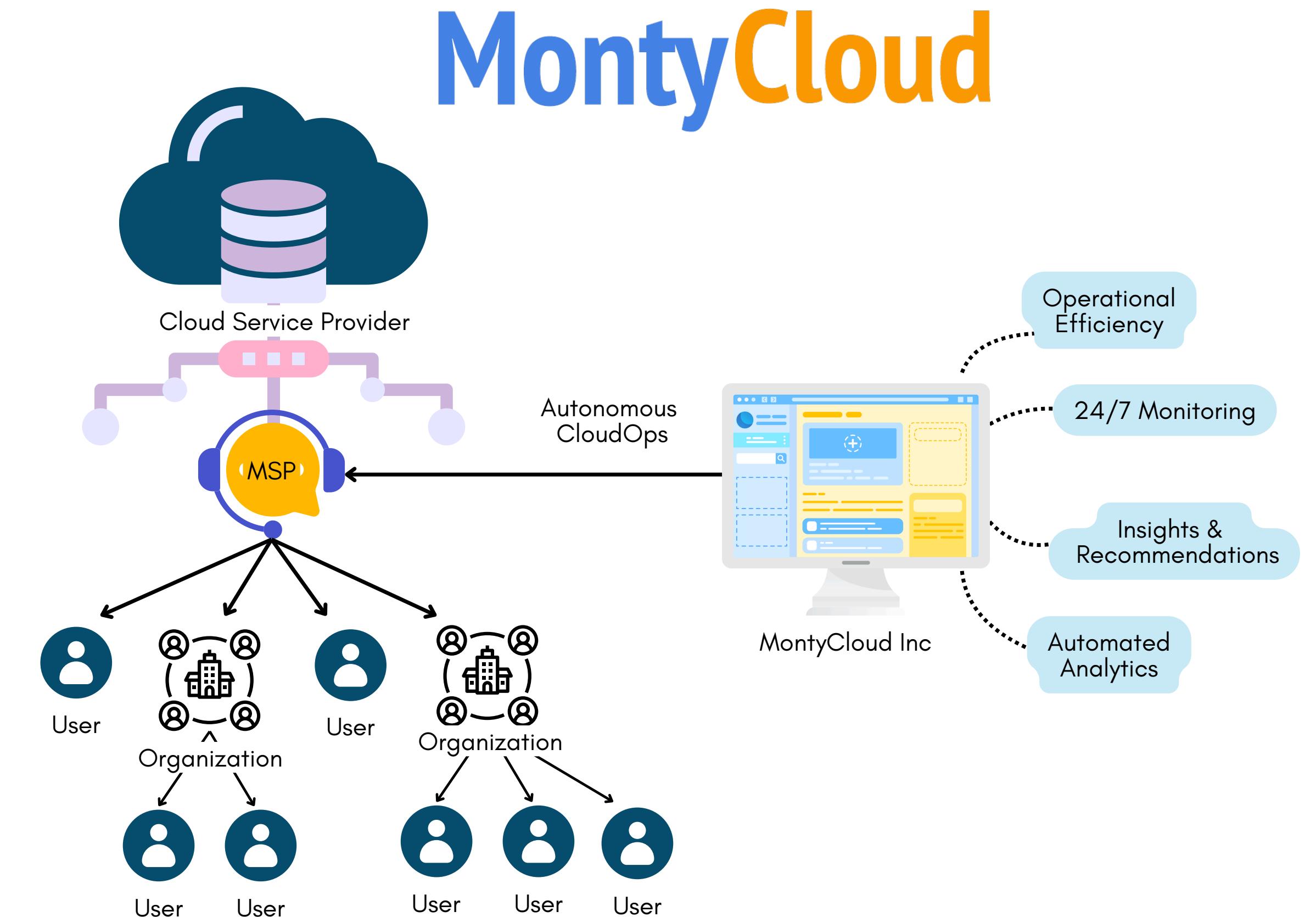
The screenshot displays the CloudOps Copilot dashboard, which integrates multiple cloud management tools. On the left, a sidebar lists navigation options such as Home, Projects, ASSESSMENTS, INVENTORY, GOVERNANCE, DAY2 CLOUDOPS, and Reports. The main dashboard features several key components:

- Open Ops Issues:** 190 total issues, including 80 Remediations By DAY2™, 25 Security, and 15 Compliance.
- Security Posture:** Last run at 21-Nov-23 09:41. It shows a Security Bot (ACTIVE) and three top violations: S3 Buckets should have a bucket policy configured (10), RDS instances should have encrypted storage (8), and IAM users should not have attached in-line policies (8).
- Compliance Assessment:** Last run at 21-Nov-23 09:59. It shows a Compliance Bot (ACTIVE) and three top violations: S3 Account Level Public Access Blocks (12), IAM Root user access key check (9), and EBS Volumes should be encrypted (8). Industry standards compliance is shown: HIPAA (83%), CIS (67%), FedRAMP (91%), NIST (83%), and PCI (100%).
- AWS Costs:** Total Spend is \$13,798.33. Cost by View shows HR Department (\$1456.27) and Dev Resources (\$400.22). Top Services include Instance (\$9,081.00), Volume (\$4,081.00), and Snapshots (\$3,234.00).
- Cost Optimization:** Last run at 21 Nov 2023 10:13 AM. It tracks potential cost savings of \$1425, with 0 Over Provisioned and 0 Under Provisioned resources. It also monitors 64 Abandoned Resources and 0 Needs Optimization.
- Cloud Footprint:** A world map showing the distribution of 3399 total resources across various AWS regions, with counts ranging from 61 to 147.
- Top Resources:** A summary of the most abundant resource types:
  - Compute Instance: 3005
  - Image: 200
  - EBS Volume: 35
  - VPC Endpoint: 30
  - Virtual Private Cloud: 30
  - EBS Snapshot: 20
  - SNS Topic: 15

Work done in collaboration with MontyCloud Inc.

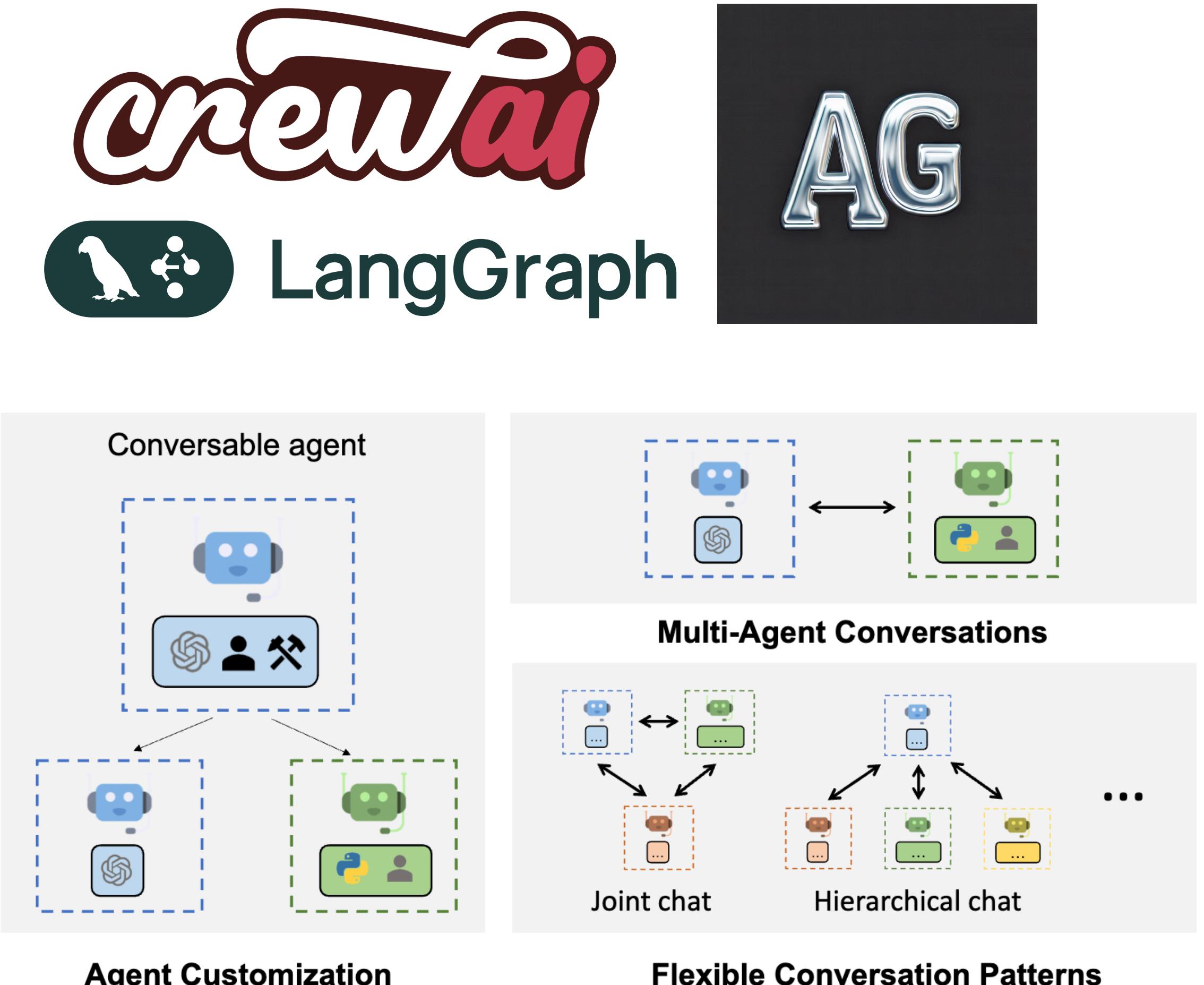
# Engineering Challenges (SE for GenAI)

- **Managing Distributed Data**
  - Diverse data sources
- **Maintainability**
  - Large code base, time for updates
- **Extensibility and Modularity**
  - Single vendor, ease of extensions!

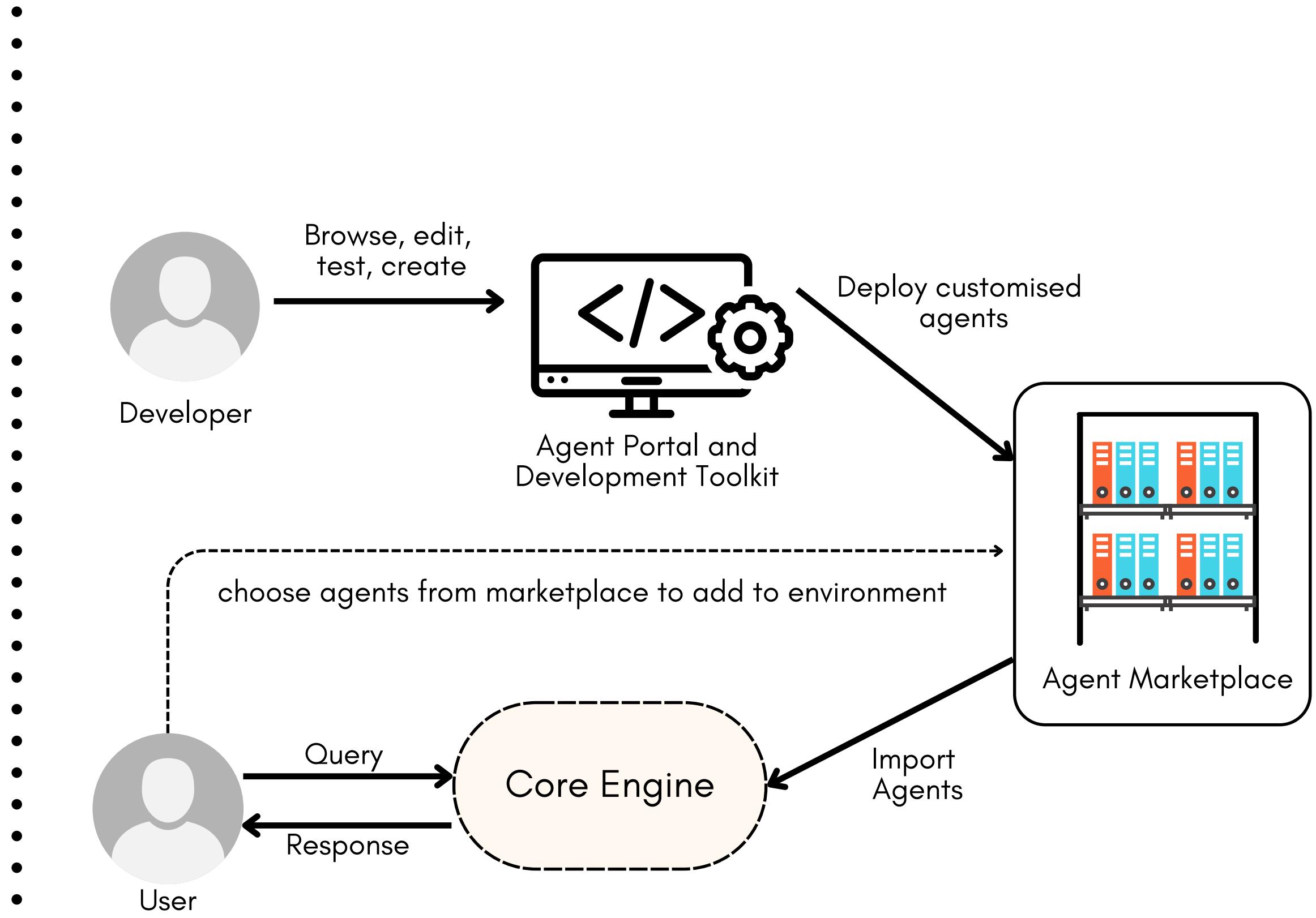
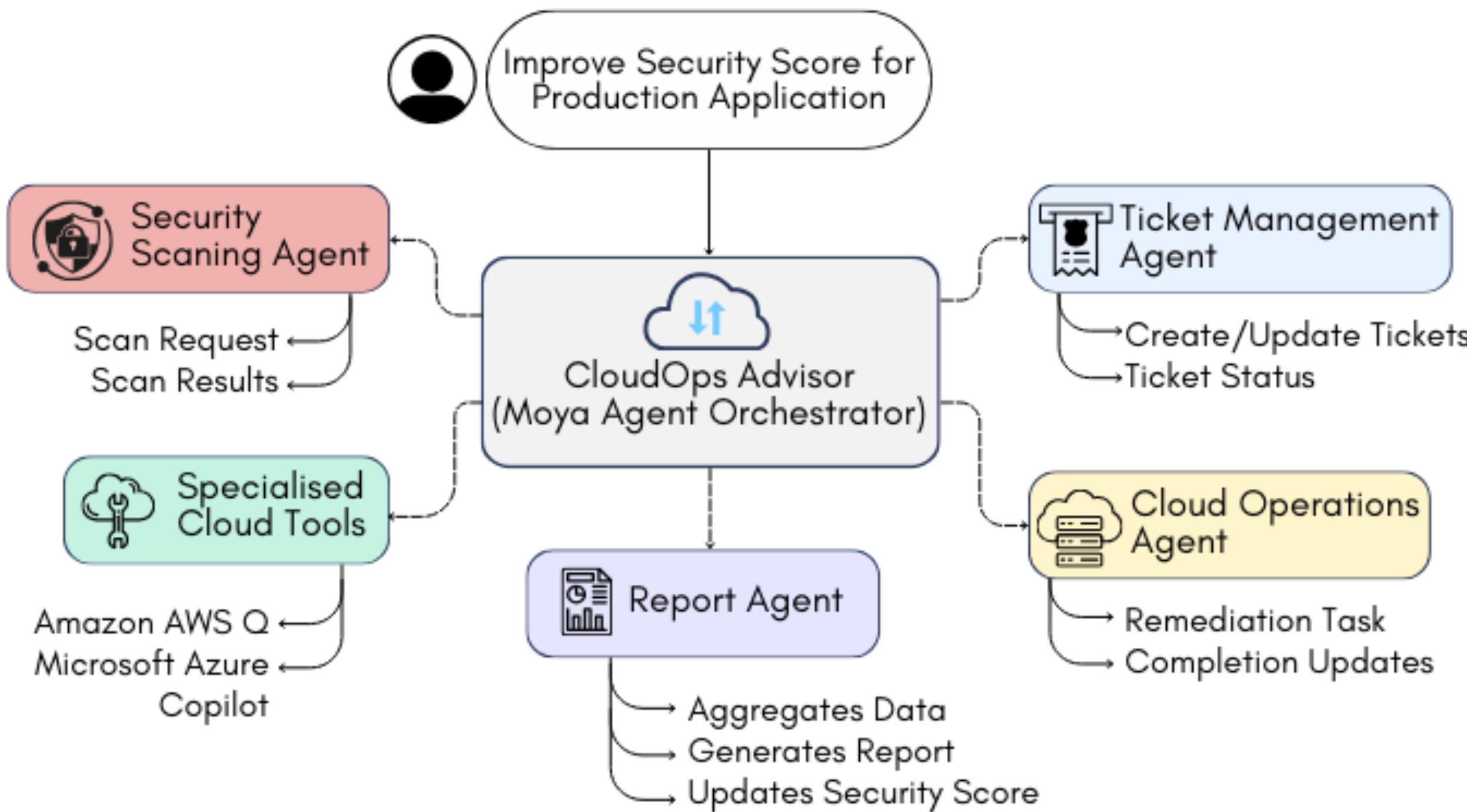


# Moving to the world of Agentic AI

- **Agent:** Any system or program that can autonomously perform a task on behalf of user or another system
  - Each agent has role and responsibilities
  - Can leverage tools at their disposal to achieve their tasks
  - Have their own memory. Converse with other agents
  - Move from API to text!



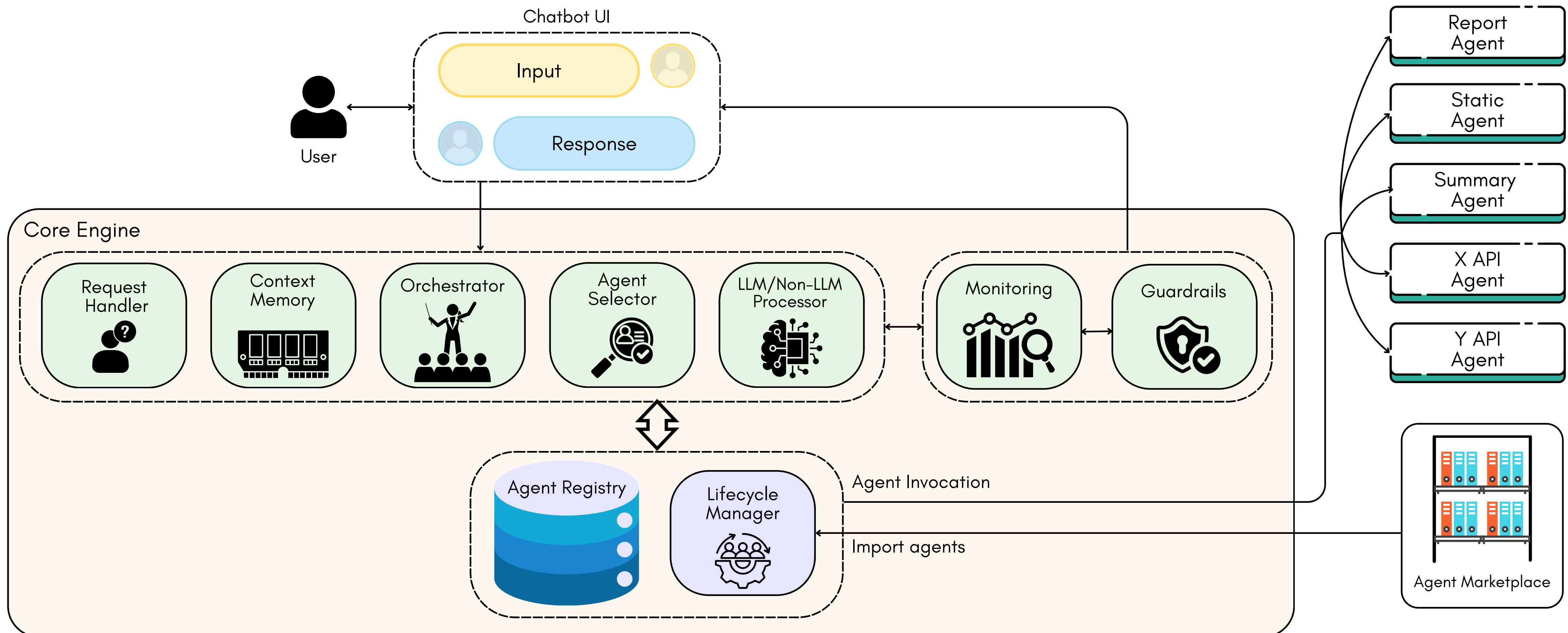
# Can we go Multi-agent?



Meta orchestration Framework



MOYA repo



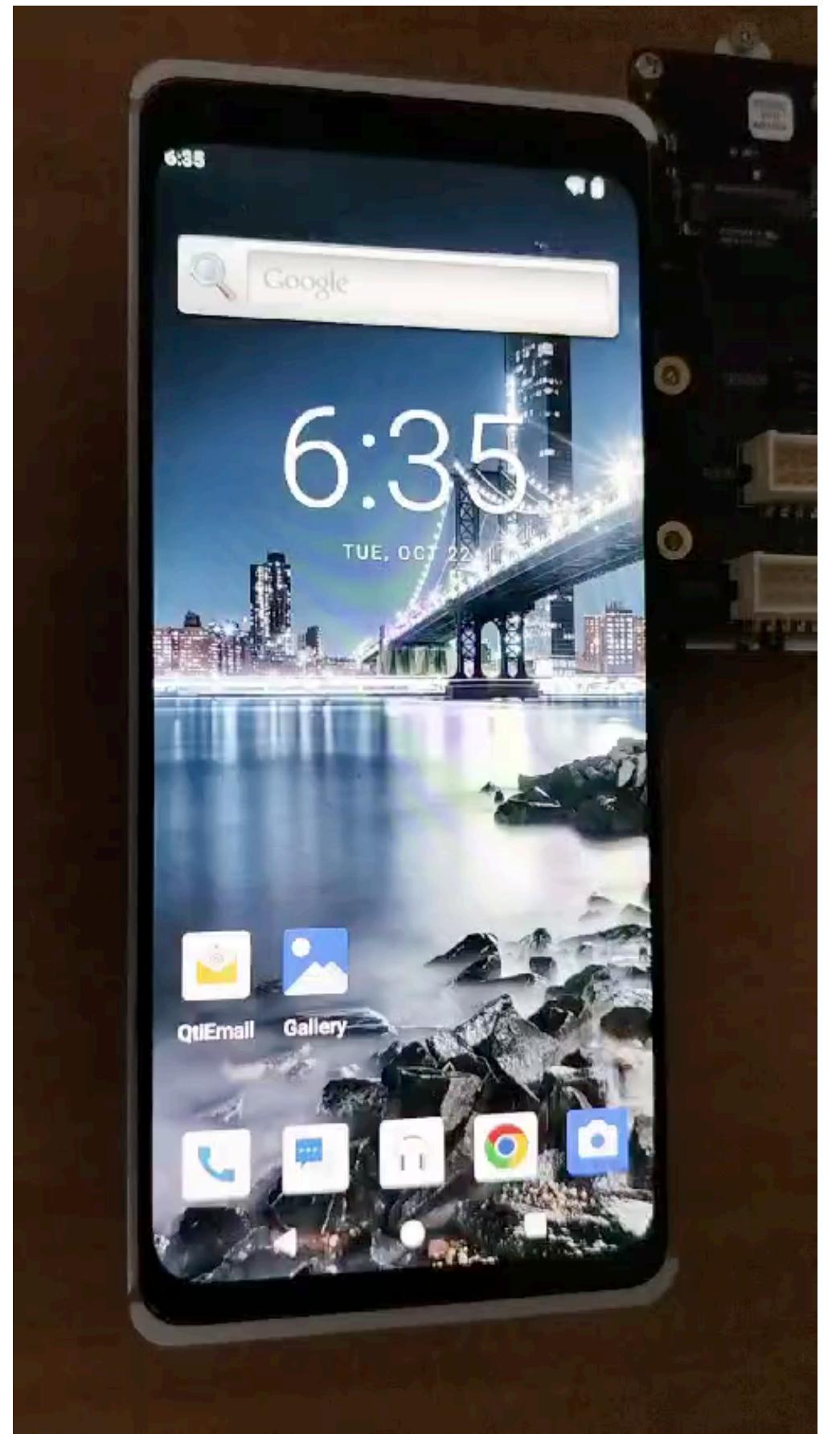
INTERNATIONAL INSTITUTE OF

HYDERABAD

# Into the world of SLMs for SE

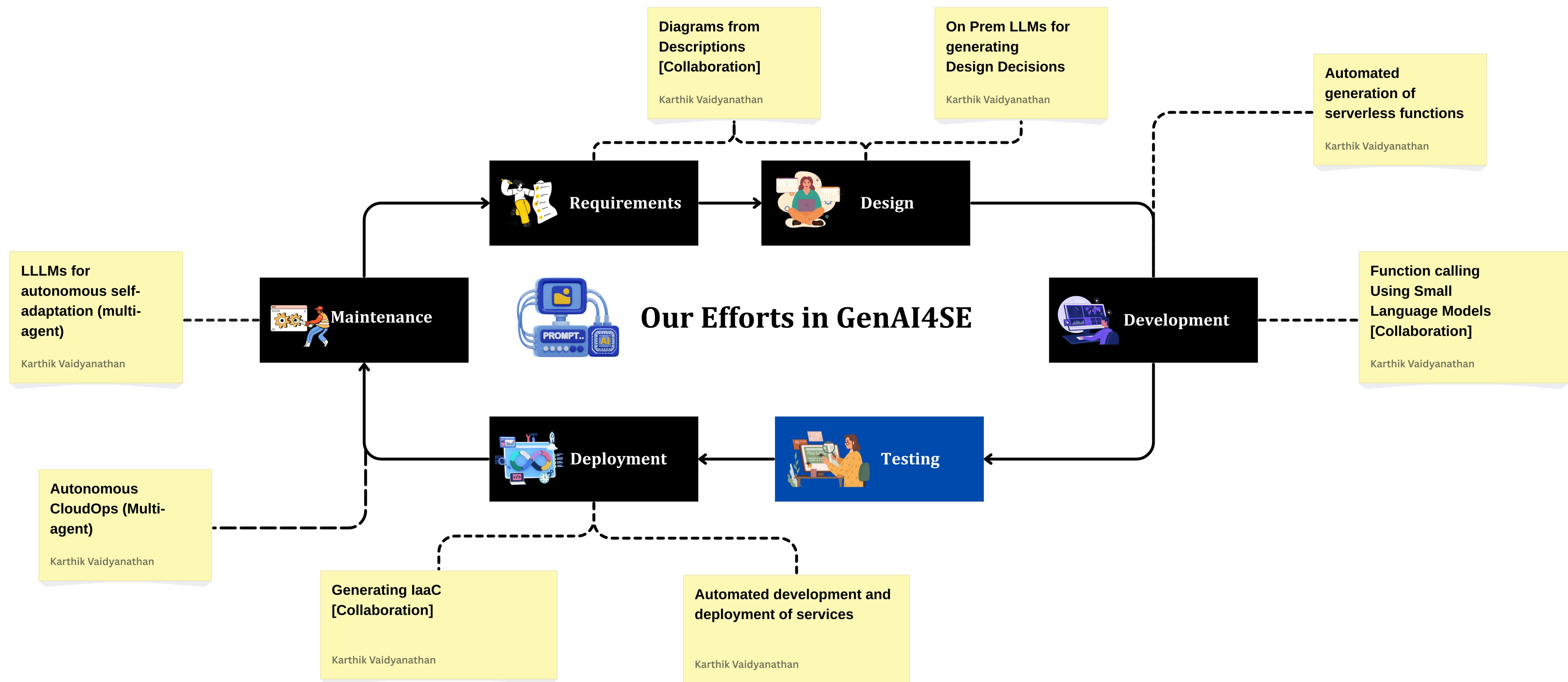
Qualcomm

- LLMs are great but there are also challenges in using proprietary LLMs
- Ongoing research in:
  - SLMs for architects for design decisions
  - SLMs for edge deployment (Qualcomm EdgeAI labs @IIITH) using QIDK, Qualcomm
  - Using SLMs for function calling - Code generation (with Precog, IIITH)
  - .....



Source: <https://www.qualcomm.com/videos/dynamic-ml-model-switching>





# Key Takeaways

***LLMs can be a best friend to a software architect if used wisely!***

- LLM presents a great opportunity for effective SE!
- We need to have an effort to have qualitative data on architecture - **ArchBench is a step!**
- Domain specific LLMs which are smaller shall be the way forward - collection of SLMs!
- Need for better ways to architect/engineer systems around LLMs (SA4LLM-enabled Systems)
- LLMs are not here to replace but to support!
- **Systems thinking is becoming more important!**

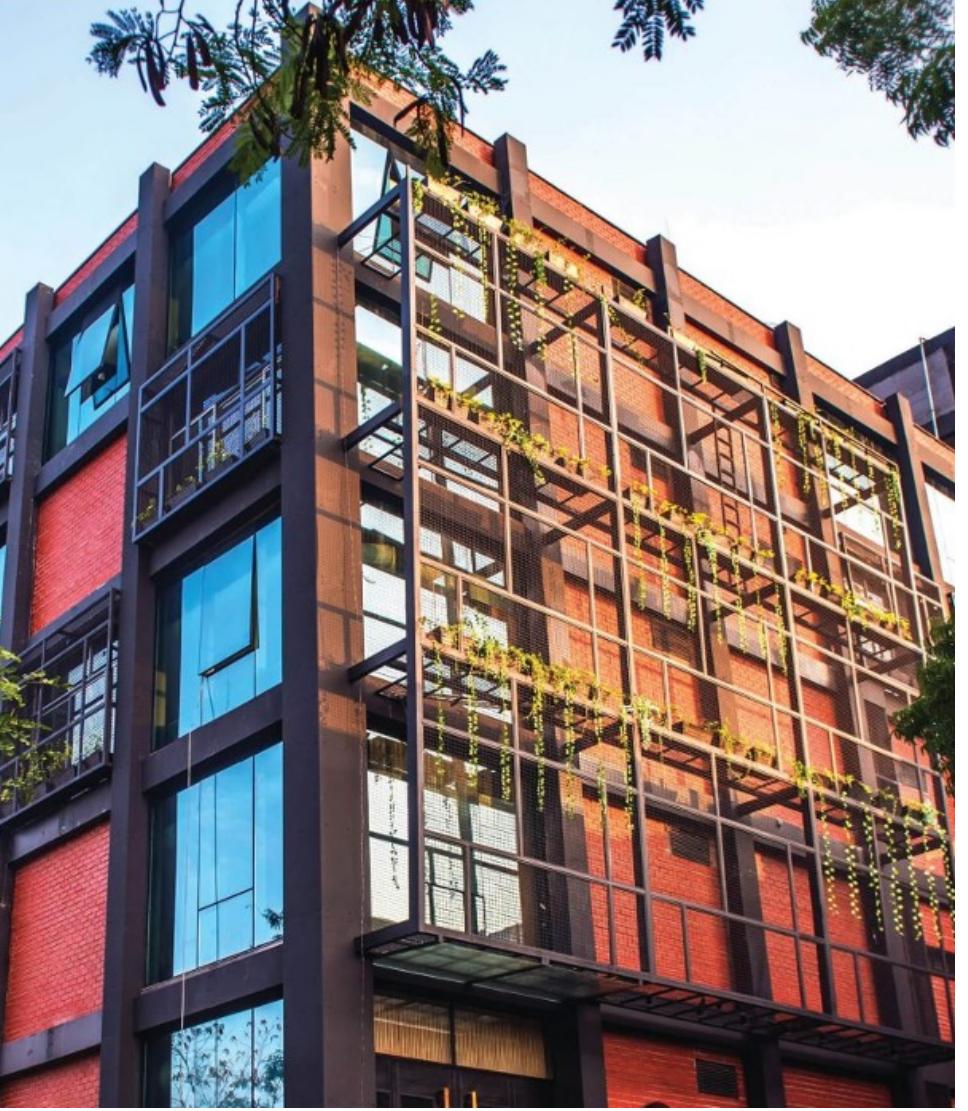
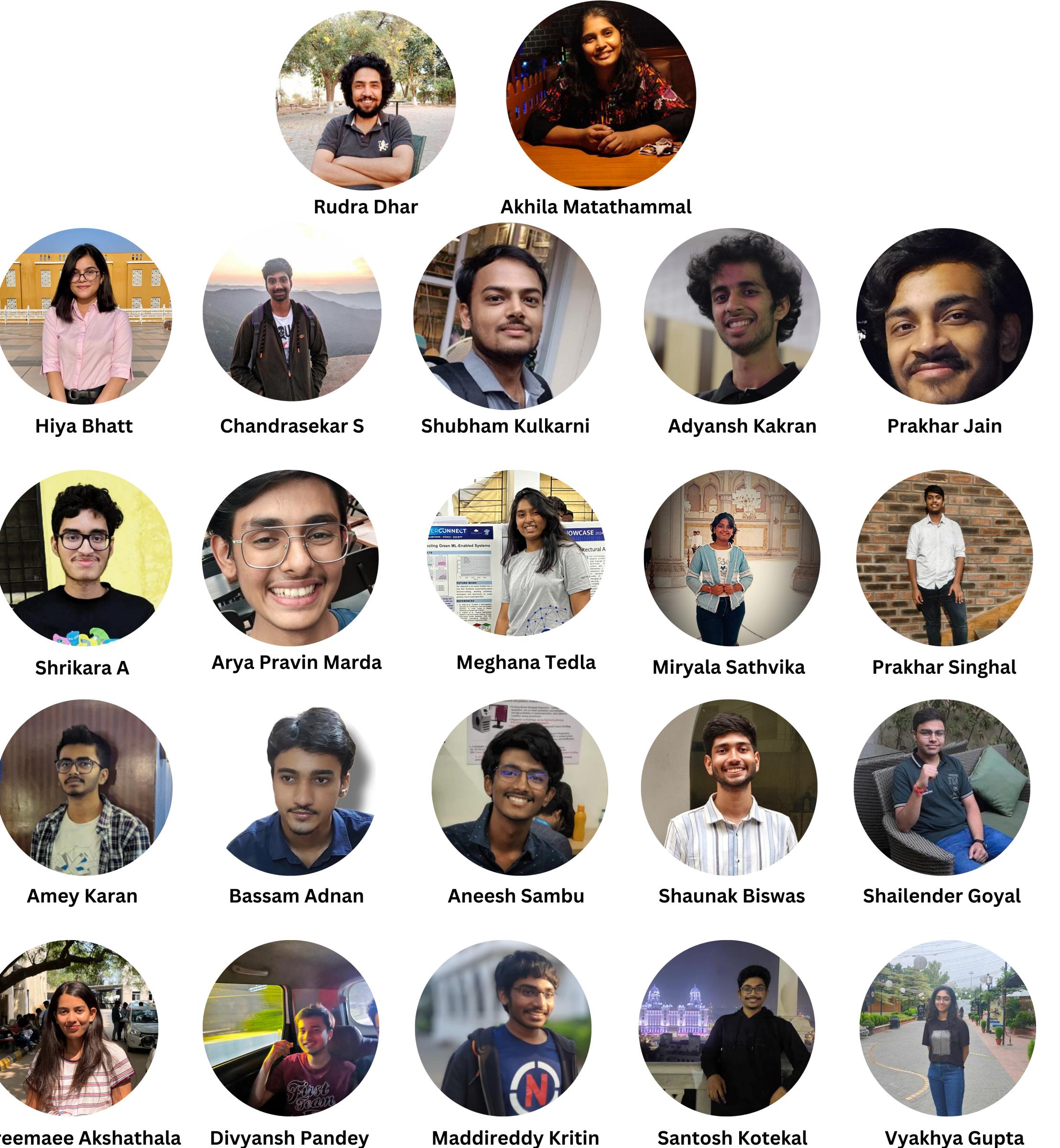


# The Future is here

- **Context is needed** - Capturing organizational aspects needs work, code can help! (Ongoing)
- **LLMs will hallucinate** - No stopping that but we can reduce it - better engineering!
- **Multiple agents collaborating** together to help software engineers
- **Large action models (LAMs)** for self-adaptation, task generation
- Lot of potential for support in **architecture migration**
- **SE Process will also need upgrades!!**



# SA4S@SERC



<https://serc.iiit.ac.in>



Team SA4S

<https://sa4s-serc.github.io>



<https://sa-ml.github.io/saml2025/>

@ ICSA 2025, SAGAI 2025



Thank you

Web: [karthikvaidhyanathan.com](http://karthikvaidhyanathan.com)  
Email: [karthik.vaidhyanathan@iiit.ac.in](mailto:karthik.vaidhyanathan@iiit.ac.in)  
Twitter: @karthi\_ishere



IEEE Software Magazine



SE Radio Podcasts

