



# Deploying multi-agent system in production at scale



The better the question. The better the answer. The better the world works.



Shape the future  
with confidence

Prototypes are easy. Production is, well, hard!

Many organizations embrace agentic AI to automate tasks and boost efficiency. **But does a working multi-agent workflow mean we have a system ready for real-world production?**

# What is multi-agent AI systems?

An AI multi-agent system is a distributed system composed of multiple intelligent agents that can sense, learn, and act autonomously to achieve individual and collective goals.

## Multi agent capabilities:

- Flexibility and scalability
- Robustness and reliability
- Self-organization and coordination
- Real-time operation

### Project manager agent

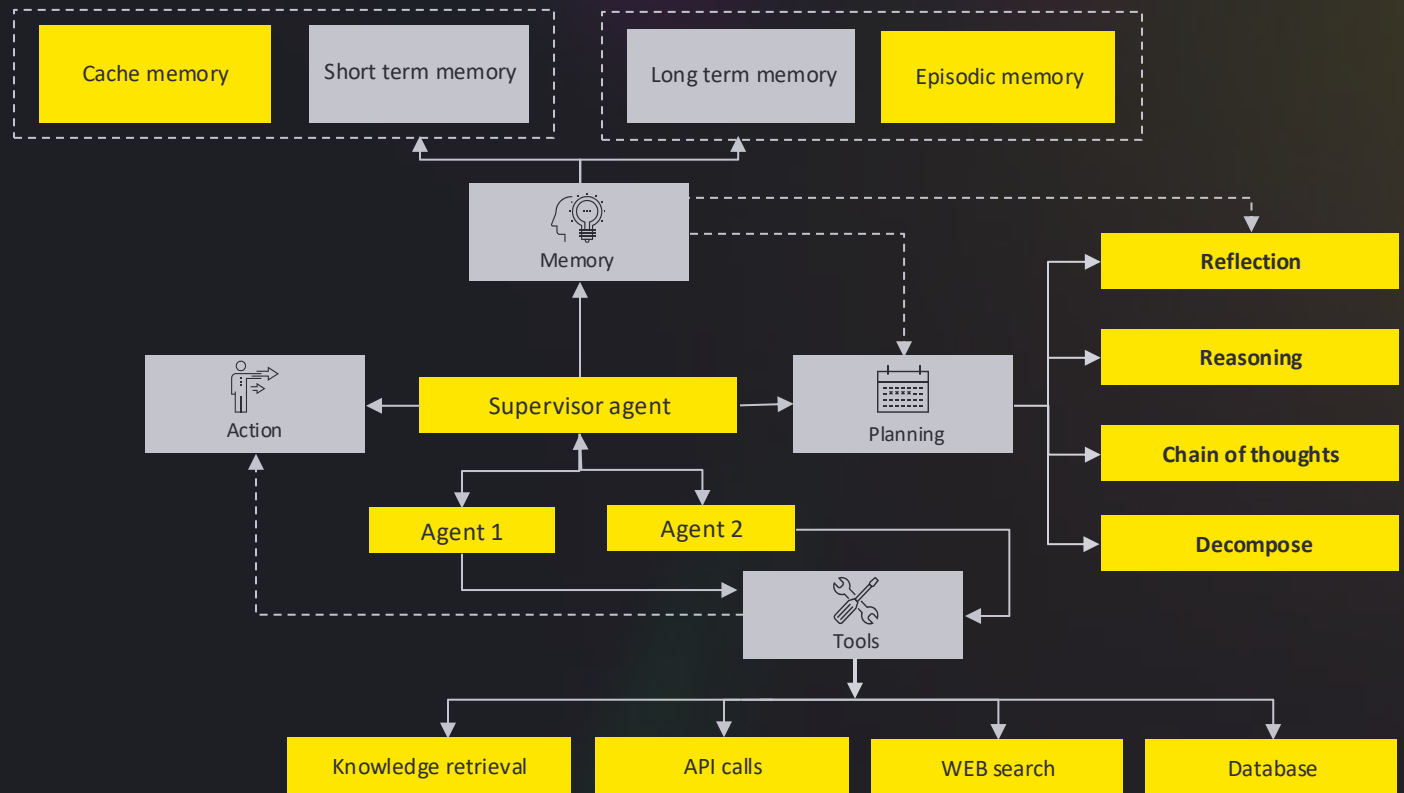
- Task: Oversees project management
- Process: Coordinating efforts among agents and ensuring the process aligns with overall project goals
- Output: Smoothly running project that stays on target

### Content developer agent

- Task: Creates content, writes stuff, and helps make documents
- Process: Creates ideas, writes them down, and polishes them
- Output: Written content and materials for the project

### Market research analyst agent

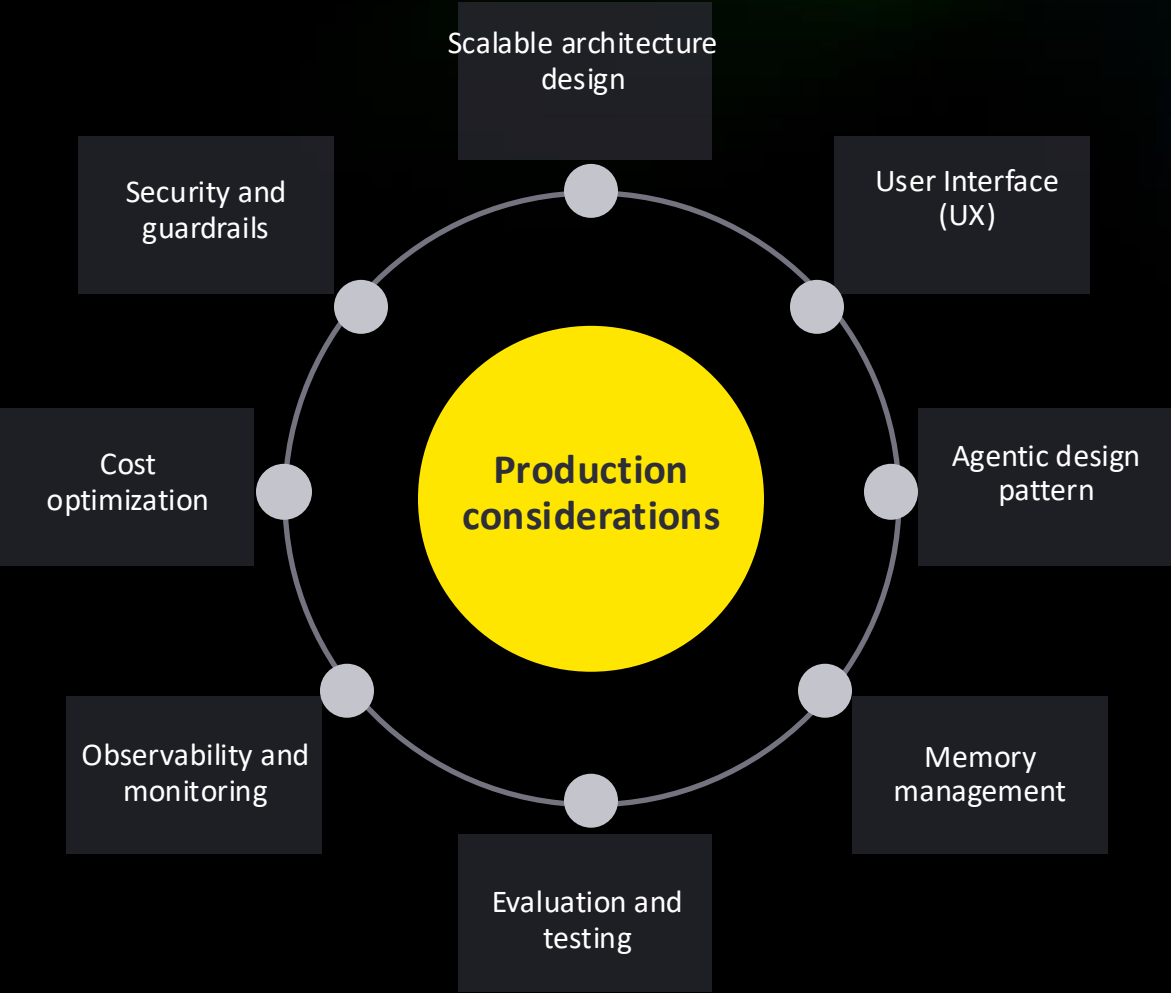
- Task: Looks for market info, studies trends, and gives advice based on what they find
- Process: Collects data, analyses it, and draws conclusions
- Output: Reports and insights about the market



# Overview of production challenges in multi-agent system

Multi agentic workflow challenges	Design challenges	Limit in no of agent and tools	Clear objectives on task and agent definitions. Prefer code logic over agent
		Evaluation challenges	Define robust evaluation pipeline for faster and holistic debugging and proper feedback loop
		Selection of agent pattern	Avoid complex agent network. Break the Agent flow into specialized workflow and assign skill specific tools
	Tool handling challenges	Failure in tools' output	Enforce structured output format. Handle default flow in case of failure
		Vague tools' description and parameters	Ensure well-defined parameters and usage guidelines
		Restrict sensitive or unauthorized access of tools	Define logic for human intervention or blocking user in case of accessing tools
	Agent collaboration challenges	Dependency b/w agent's output	Well defined guidelines with dynamic few shots and structured output schema
		Infinite looping	Prefer prompt logic to justify agent's action or reasoning. Define clear termination conditions
		High cost of running	Implement model tiering based on tasks complexity
		Right agent or tool selection	Specialized prompt and hierarchical design
	Deployment challenges	Guardrails and security	Rule-based and LLM-based validation. Define human in the loop in case of critical tasks. Prefer default behavior over ambiguity
		Scalability	Scalable architecture for state management, agent calling and low latency
		Handling of failure	Clearly defining default behavior based on recommendation for domain expert

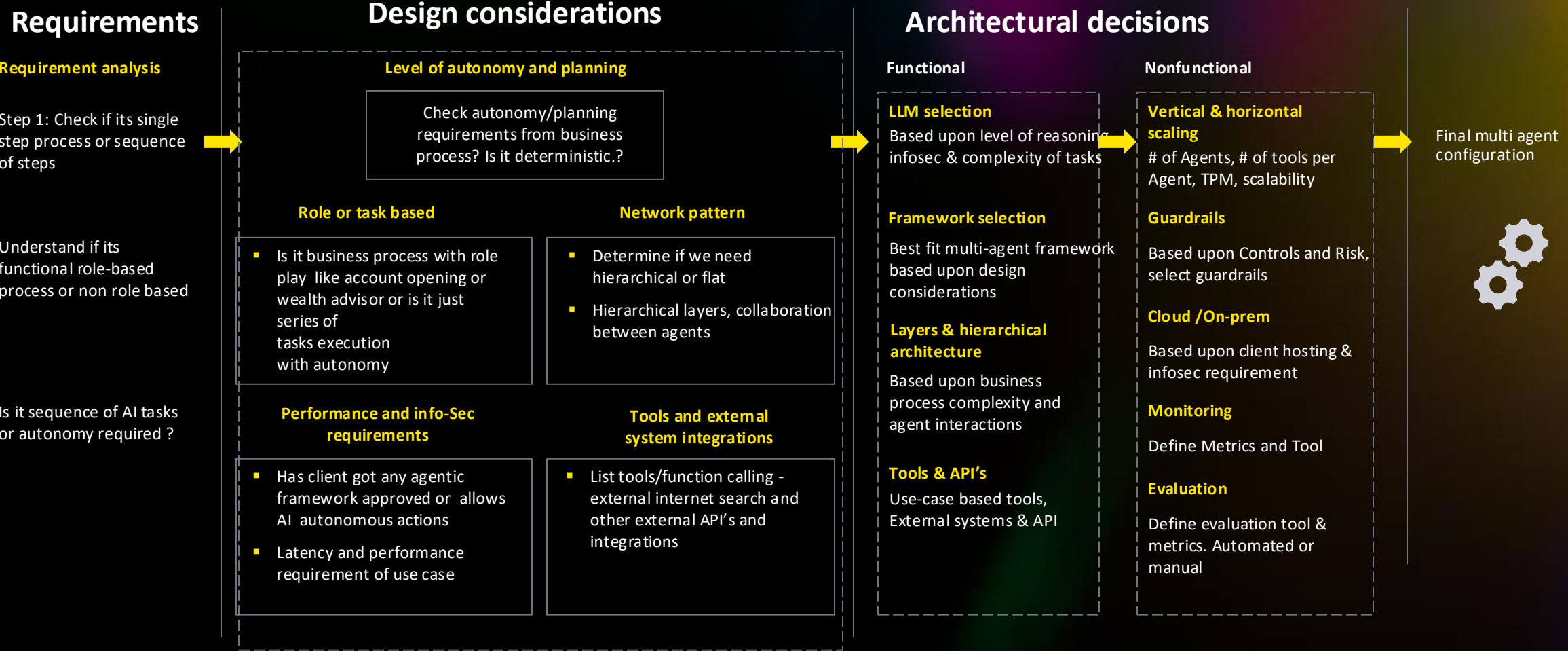
# Key production considerations for multi-agent systems



Security and guardrails	<ul style="list-style-type: none"><li>▪ Prevent unwanted tool/agent calling and prompt injection</li><li>▪ Enforce strict IAM policies to restrict sensitive tools</li></ul>
Cost optimization	<ul style="list-style-type: none"><li>▪ Implement token efficiency strategies – caching repeated queries</li><li>▪ Model tiering routes tasks for cost appropriate backends</li></ul>
Observability and monitoring	<ul style="list-style-type: none"><li>▪ Designing agent and tool specific metrics to assess performance</li><li>▪ Tracking and logging of user feedback to further enhance workflow</li></ul>
Evaluation and testing	<ul style="list-style-type: none"><li>▪ Robust evaluation pipeline to test E2E workflow and agent trajectory</li><li>▪ Metrics to assess tool interactions and failure</li></ul>
Memory management	<ul style="list-style-type: none"><li>▪ Leverage different agentic memories to maintain state in agentic system</li></ul>
Agentic design pattern	<ul style="list-style-type: none"><li>▪ Selection of suitable agentic pattern</li></ul>
Scalable architecture design	<ul style="list-style-type: none"><li>▪ Design infrastructure that scales dynamically and handles fluctuating demands. Horizontal scaling via serverless Kubernetes and vertical scaling by adding multi-region LLM deployment</li></ul>
User interaction	<ul style="list-style-type: none"><li>▪ UX design to allow user to interact with agents in human-in-the-loop scenario</li></ul>

# Solution architecture design consideration

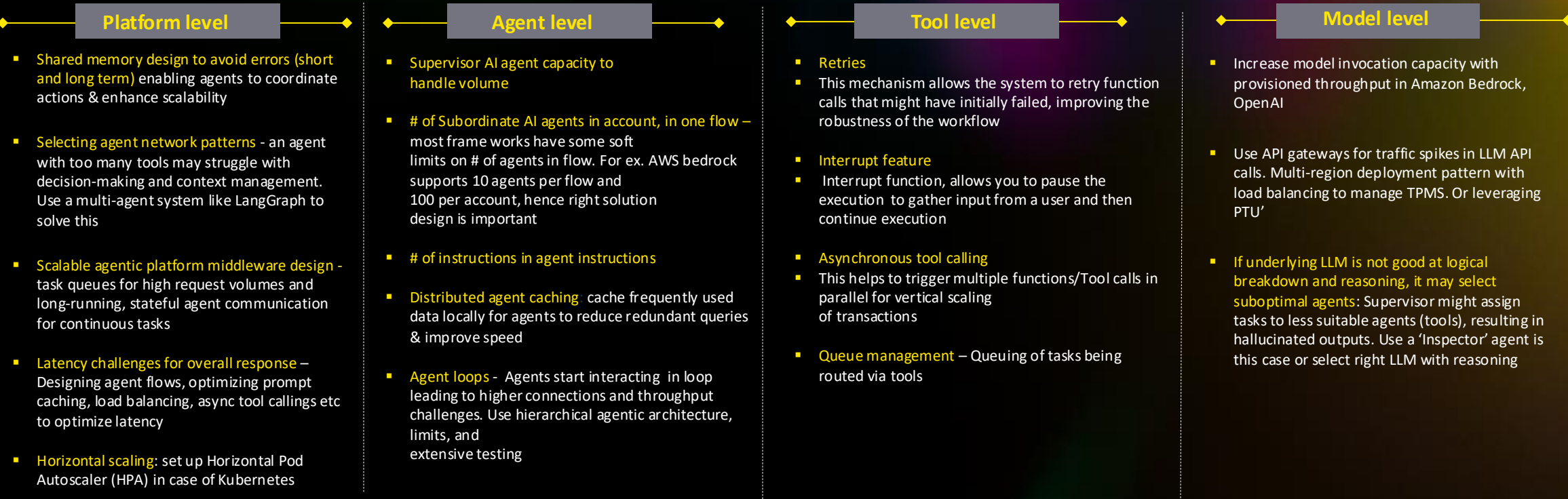
We recommend step by step analysis of business requirements and then mapping it to various design considerations and non-functional requirements to arrive at optimal multi-Agent solution specifications required





# Scalability best practices

Scalability: The ability of AI agents to handle increased workloads, user interactions, and data volumes efficiently without compromising performance or reliability.



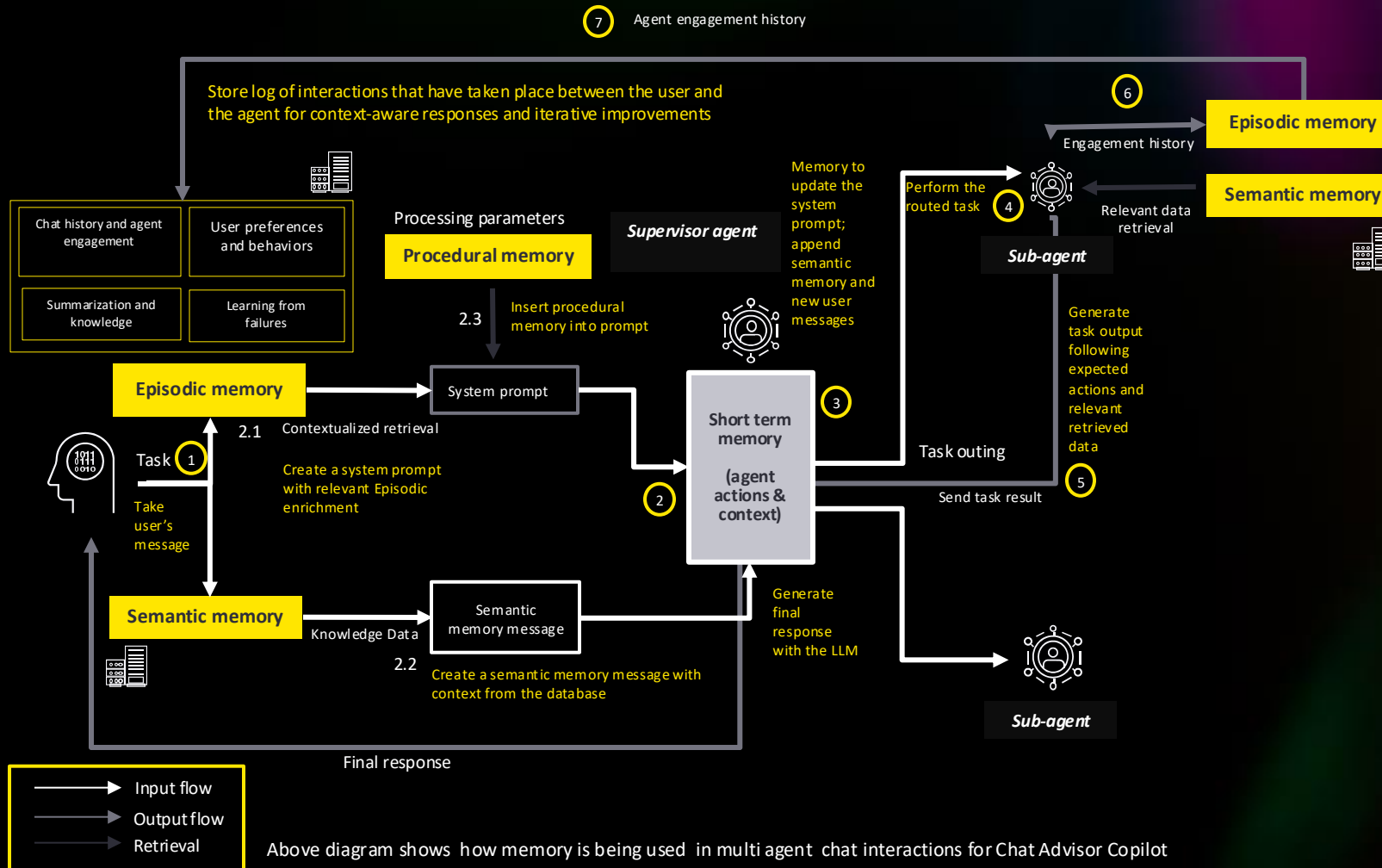
## Challenges in scaling multi framework production deployments

<ul style="list-style-type: none"><li>Cloud does not support agents in all the regions</li><li>Memory of past interactions, state management</li><li>Latency</li></ul>	<b>Limitations at agent level</b> <ul style="list-style-type: none"><li>No of calls each agent can handle</li><li>No of agents per account</li></ul>	<b>Limitations at tool level</b> <ul style="list-style-type: none"><li>Limit on API calls</li><li>Tools failures</li><li>Tools - Synchronous calls</li></ul>	<b>Limitations at model/LLM level</b> <ul style="list-style-type: none"><li>Context window exceeding</li><li>Token limit issue - - No of input and output token processing per minute</li></ul>
--	--	--	---



# Importance of memory management in agentic workflow

## Multi-agent memory management



**Short-Term Memory (STM)**: Temporary storage for immediate tasks like maintaining context of recent conversations optimized for quick access.

In diagram, it is leveraged for storing contextual understanding of conversation across recent chat sessions. (f1 to 30 days)

**Long-Term Memory (LTM)**: Persistent storage for knowledge retention across sessions.

- Episodic: past experiences of users & agents
- Procedural: knowledge of the agent's code,
- Semantic: store general knowledge, facts, and information from various sources

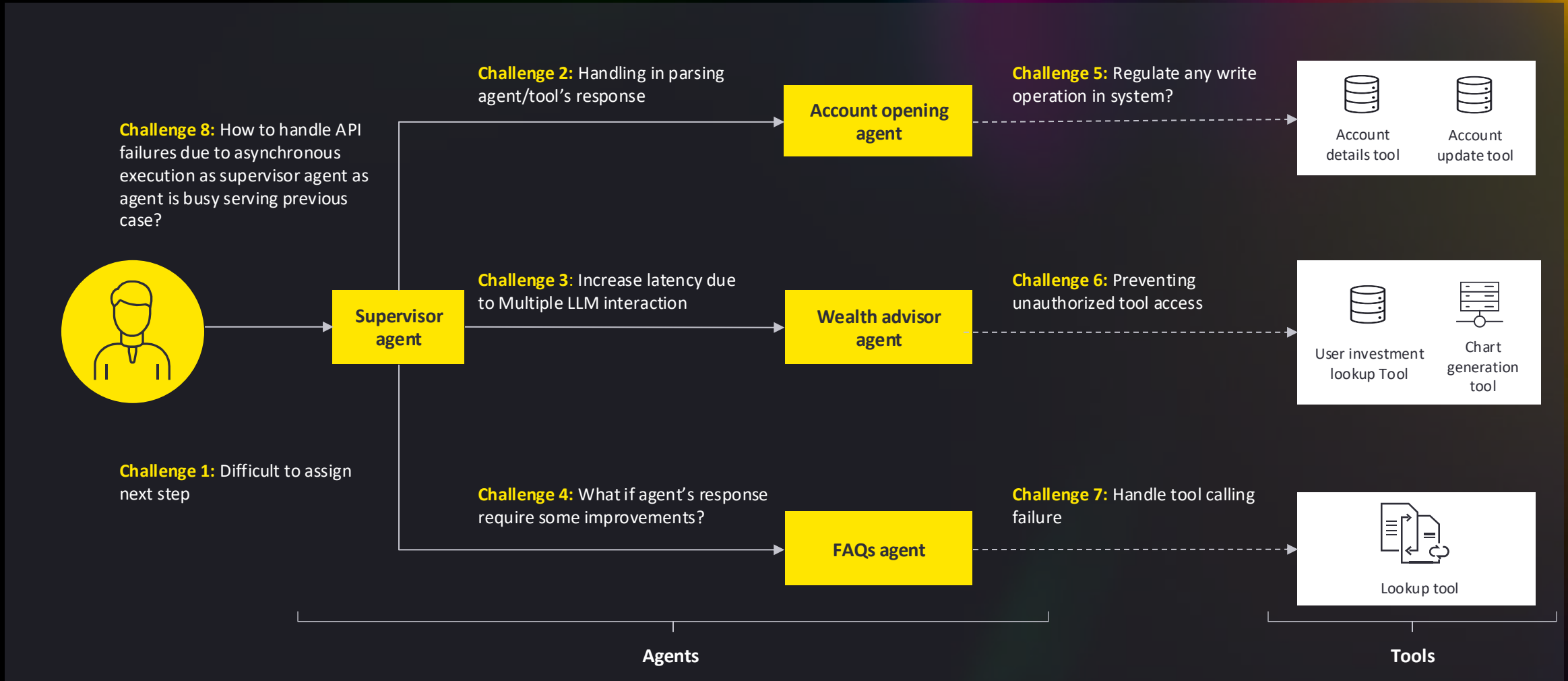
In diagram, it is leveraged to

- Store performance and behaviour of the supervisor agent by maintaining a history of its Reasoning (R), Planning (P), and Tasks (T) etc.
- Stores learnings from failures, user behaviors and preferences, agent history, summarization

Right settings for memory helps in improving accuracy , providing personalized customer experience and improved performance over time.



# Key challenges in deploying multi AI-agent solutions



# Challenges 1 – Agent trajectory and dependence between agent responses (cont..)

## Challenge#1

### Understanding agent's trajectory

In multi-agent system, agent usually assigns direct tasks to other different agents. For example, supervisor agent is responsible for assigning tasks to specialized agents.

## Our approach

- Add prompt logic to reveal an agent's 'thought process' or reasoning to assign next task using CoT
- Implement dynamic shots to make agent aware of context
- Define explicit rules and logic that determine how the supervisor agent assigns tasks

## Challenge#2

### High dependencies between each agent's responses

Multi-agent systems depend heavily on agents passing messages and instructions to one another. But what happens when agents misinterpret those instructions, or worse, don't receive them properly at all?

## Our approach

- Restrict the response into more reliable structures like json-based requests that help to mitigate miscommunication
- Design single supervisor or workflow engine. This controls message flow or task order
- Designed proper agent memory management so it can help reduce ambiguity when making decisions
- Design business criteria that weigh the importance of each agent's input based on the domain/use-case

## Challenge#3

### Increased latency due to multiple LLM interactions

When multiple AI agents collaborate (e.g., a supervisor agent coordinating with account opening and wealth advisor agents), sequential LLM queries introduce latency. This can degrade user experience and slow down decision-making.

## Our approach

- Use asynchronous programming to handle multiple agent interactions concurrently
- Multi-region LLM deployment and API gateway to ensure availability
- Implement short-term memory caching to store responses from LLMs
- Replace large LLMs with lighter, fine-tuned models for non-complex queries

# Multi AI-agent challenges and EY organization's solution approach

## Challenge#4

### Handling cases where agent's response needs improvements

Sometimes an agent may generate incomplete, vague, or incorrect responses. Without quality control, these responses can propagate errors through the system.

## Our approach

- Use an independent verification agent to cross-check responses
- Implement self-evaluation prompts where the agent reflects on its own output
- Provide a human override where experts can modify agent outputs before they reach

## Challenge#7

### Handling unexpected tool responses that could break workflow

In multi-agent system, AI agents rely on external tools (e.g., databases, lookup tools) that may return responses in unexpected formats or fail due to technical issues. If not handled properly, this can break the workflow.

## Our approach

- If a tool fails, provide a default response or an alternative action (e.g., retry, switch tools) or delegate the request to an alternative agent or another tool
- If a request fails due to incorrect parameters, an agent can intelligently modify the query and reattempt

## Challenge#8

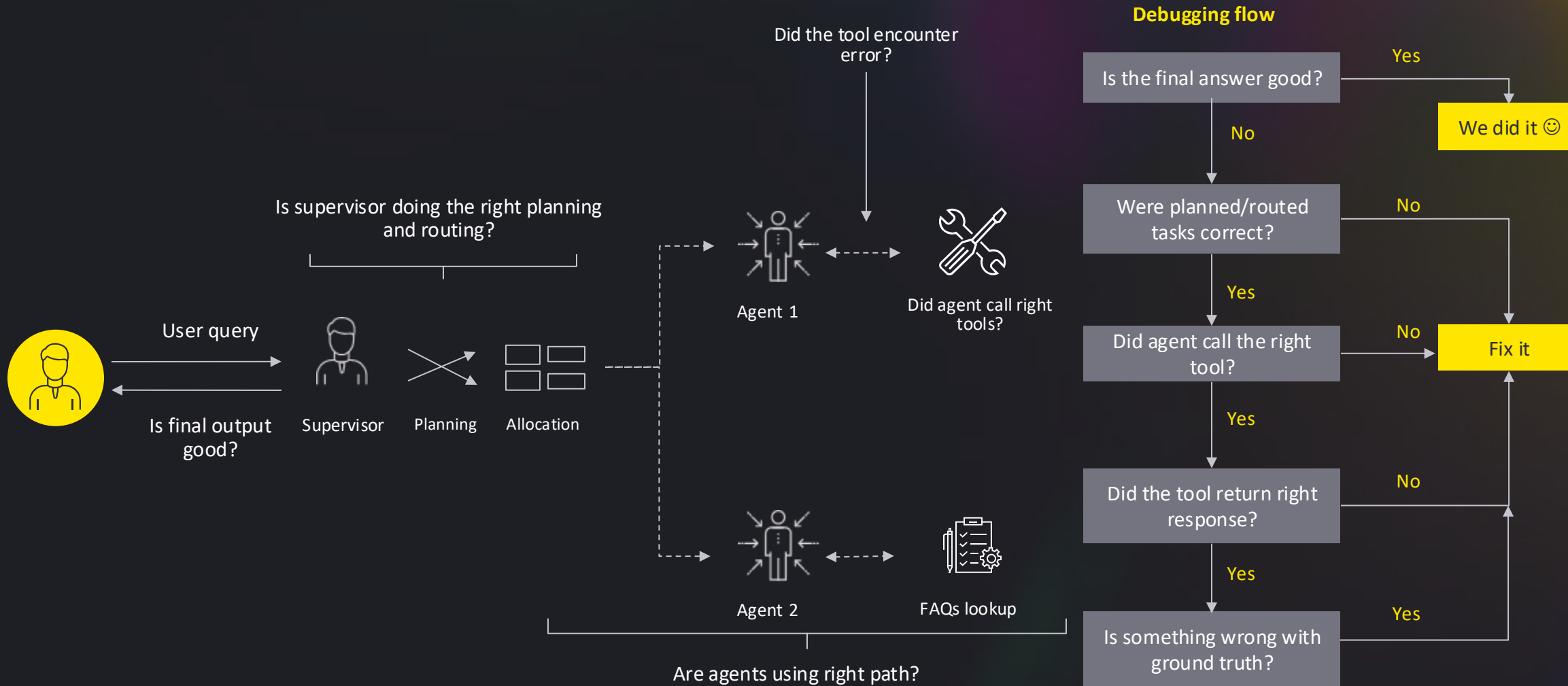
### Handling API failures due to asynchronous execution

In multi-agent system, when a supervisor agent assigns a task, subsequent API calls may fail if the agent is still processing the previous request. This leads to bottlenecks, dropped tasks, and workflow disruptions.

## Our approach

- Instead of waiting for an agent's response synchronously, supervisor agent should use an async task queue to submit requests and continue processing
- Deploy multiple instances of the same agent and distribute tasks among them to avoid congestion
- If an API call fails due to the agent being busy, implement a circuit breaker mechanism to prevent repeated failures

# Multi-agent framework evaluation approach



# Evaluating multi-agent framework

It is important to assess overall and individual performance of flow based on application requirements.

## AI agent evaluation strategies

1

### Define and build golden dataset

To create a comprehensive dataset that enables thorough testing and evaluation of LLM AI agents, ensuring their performance is aligned with real-world applications and challenges.

- Input and expect output (E2E eval)
- Input and expect Trajectory (trajectory eval)
- Input and expect step (step eval)

2

**E2E evaluation:** to assess the overall performance

**Key Metrics:**

Correctness, completeness, faithfulness, relevance

**Trajectory evaluation:** to analyse the sequence of actions by the agent

**Key Metrics:**

Exact match, in order match, number of extra step, number of mismatched entries, trajectory precision, trajectory recall

Step wise evaluation to break down the agent's performance

**Key Metrics:**

accuracy, precision, recall

Tool wise evaluation to break down the tools' performance

**Key Metrics:**

Tool selection, tool error

Evaluation framework (Langgraph, arizeai, aws bedrock agentic eval)

Input (user queries)

Application output/trajectory

Reference output

3

### Experiment results

Evaluate the performance and efficiency of the agentic workflow in a high-demand technical environment.

### Comparative analysis Set up

**Previous run:** Compare with previous run to understand the impact.

**Incremental feature experiment:** Gradual integration of features to isolate impact.

**Golden dataset:** Compare performance with golden dataset



# Agentic evaluation: defining and building evaluation ‘golden’ dataset

An evaluation dataset is a collection of data specifically curated to assess the model's performance. This golden dataset should clearly define the branch user queries, the relevance of the policy/circular to the query and the expected response.

## Considerations for building the golden dataset for Evaluation

Initiate incrementally	LLM-assist data generation	Expand dataset progressively
<ul style="list-style-type: none"><li>▪ Build a manageable scope or subset of data to start</li><li>▪ Include false positive and example of jail-break, prompt injection to evaluate solution’s response</li><li>▪ Start with a smaller dataset to pilot the process and refine methodologies before scaling up</li></ul>	<ul style="list-style-type: none"><li>▪ Utilize LLM capabilities to analyze policy document to extract key information</li><li>▪ This LLM-assist approach could speed up compilation of evaluation set</li><li>▪ Perform human-based feedback mechanism to validate the quality of generated data with respect to use-case requirements</li></ul>	<ul style="list-style-type: none"><li>▪ Continuously evaluate and refine the dataset based on feedback and insights gained during initial phases and discussion with branches</li><li>▪ Gradually incorporate additional data sources and variables to enrich the dataset's depth and quality</li><li>▪ <b>Track user feedback</b> to append edge cases</li></ul>

## Types of queries

Simple query	Reasoning based query	Multi-hop based query
Straightforward query which is direct and does not involve complex reasoning or inference.	Complex inquiries that require logical deduction, inference, or contextual understanding to generate the desired response.	Involves harnessing multiple pieces of information or knowledge sources to answer a complex question

# Agentic evaluation: challenges and their mitigation strategy

## Challenges with auto-evaluation

Utilizing LLM-as-judge is powerful and can be helpful to automate the evaluation but we have faced challenges around:

- 1 Most of the evaluation frameworks such as ragas, deepeval, provide metrics with numeric score (out of 10). We observed, the range value has variation in evaluation output
- 2 These frameworks provides some out-of-the-box metrics (faithfulness, relevancy etc.), but was difficult to infer for use-cases
- 3 Due to non-deterministic nature, multiple runs of LLM-as-judge yields different score/results. Because of which evaluation results are not consistent
- 4 These evaluation framework has generic prompts and criteria to perform evaluation

			answer_correctness					
1806_Correct/Wrong/Incomplete		max	min	q25	median	mean	q75	max
Correct	5	1.00	0.21	0.39	0.57	0.55	0.68	0.98
Incomplete	8	1.00	0.18	0.41	0.52	0.51	0.61	0.85
Wrong	6	1.00	0.19	0.25	0.43	0.42	0.50	0.80

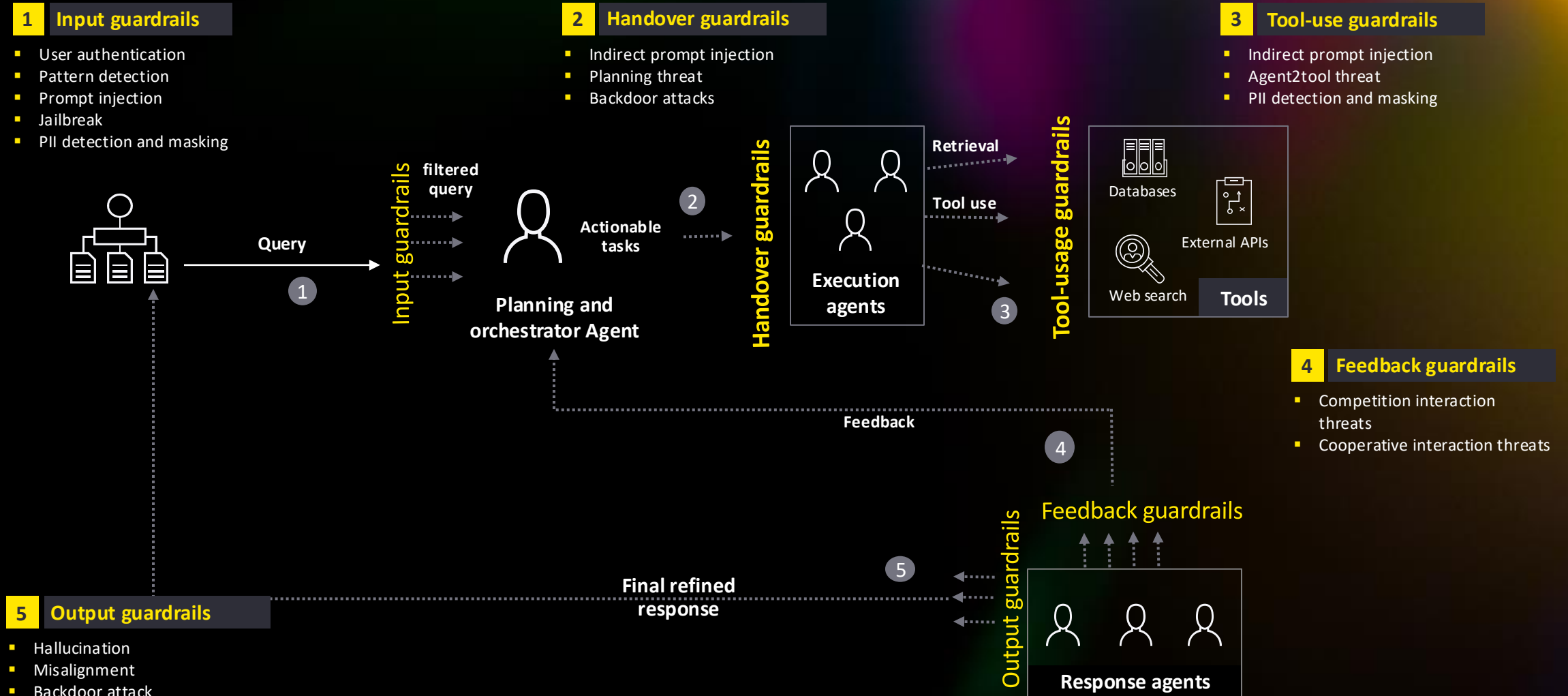


## Our unique approach/frameworks to solve key challenges

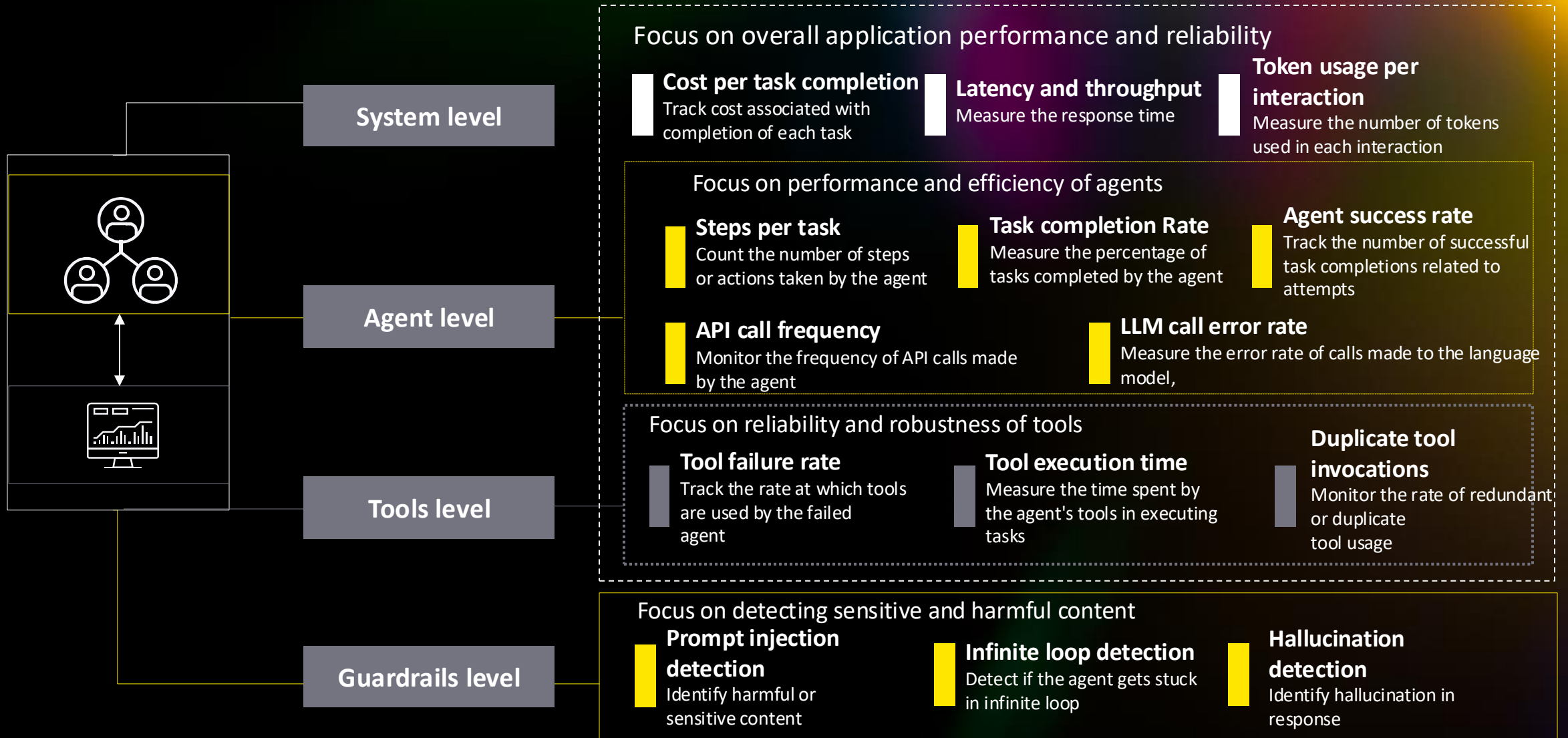
Came up with strategies to mitigate issues and improve the performance of auto-evaluation

- 1 **Modified criteria for evaluation:** instead of LLM performing numeric evaluation, we modified prompts to provide binary score based on comparison between Ground Truth data and LLM Output. As mentioned in couple of papers, we also devised criteria to compare two consecutive LLM responses to score
- 2 **Modified evaluation metrics:** instead of using typical ragas' metrics, we defined set of metrics i.e., completeness, contradiction etc. to judge LLM response on various aspects
- 3 **Self-consistency (CoT@k) approach** to mitigate the variation in LLM evaluation. We performed iteration of evaluation and opted for most frequent response
- 4 **Design domain and use-case specific prompt with few-short steps** to allow LLM to improve the evaluation response quality and provide better reasoning to assign specific labels

# Guardrails in multi-agentic system



# Key metrics for agentic AI monitoring



## EY | Building a better working world

EY is building a better working world by creating new value for clients, people, society and the planet, while building trust in capital markets.

Enabled by data, AI and advanced technology, EY teams help clients shape the future with confidence and develop answers for the most pressing issues of today and tomorrow.

EY teams work across a full spectrum of services in assurance, consulting, tax, strategy and transactions. Fueled by sector insights, a globally connected, multi-disciplinary network and diverse ecosystem partners, EY teams can provide services in more than 150 countries and territories.

All in to shape the future with confidence.

EY refers to the global organization, and may refer to one or more, of the member firms of Ernst & Young Global Limited, each of which is a separate legal entity. Ernst & Young Global Limited, a UK company limited by guarantee, does not provide services to clients. Information about how EY collects and uses personal data and a description of the rights individuals have under data protection legislation are available via [ey.com/privacy](https://ey.com/privacy). EY member firms do not practice law where prohibited by local laws. For more information about our organization, please visit [ey.com](https://ey.com).

Ernst & Young LLP is one of the Indian client serving member firms of EYGM Limited. For more information about our organization, please visit [www.ey.com/en\\_in](https://www.ey.com/en_in).

Ernst & Young LLP is a Limited Liability Partnership, registered under the Limited Liability Partnership Act, 2008 in India, having its registered office at Ground Floor, Plot No. 67, Institutional Area, Sector-44, Gurugram, HARYANA, 122003, India.

© 2025 Ernst & Young LLP. Published in India.  
All Rights Reserved.

ED None

GDS Creative Support (GDS CS): CRS\_GDS BMC\_153240819

This publication contains information in summary form and is therefore intended for general guidance only. It is not intended to be a substitute for detailed research or the exercise of professional judgment. Neither EYGM Limited nor any other member of the global Ernst & Young organization can accept any responsibility for loss occasioned to any person acting or refraining from action as a result of any material in this publication. On any specific matter, reference should be made to the appropriate advisor.

[ey.com](https://ey.com)