

From Proof-of-Concept to Production

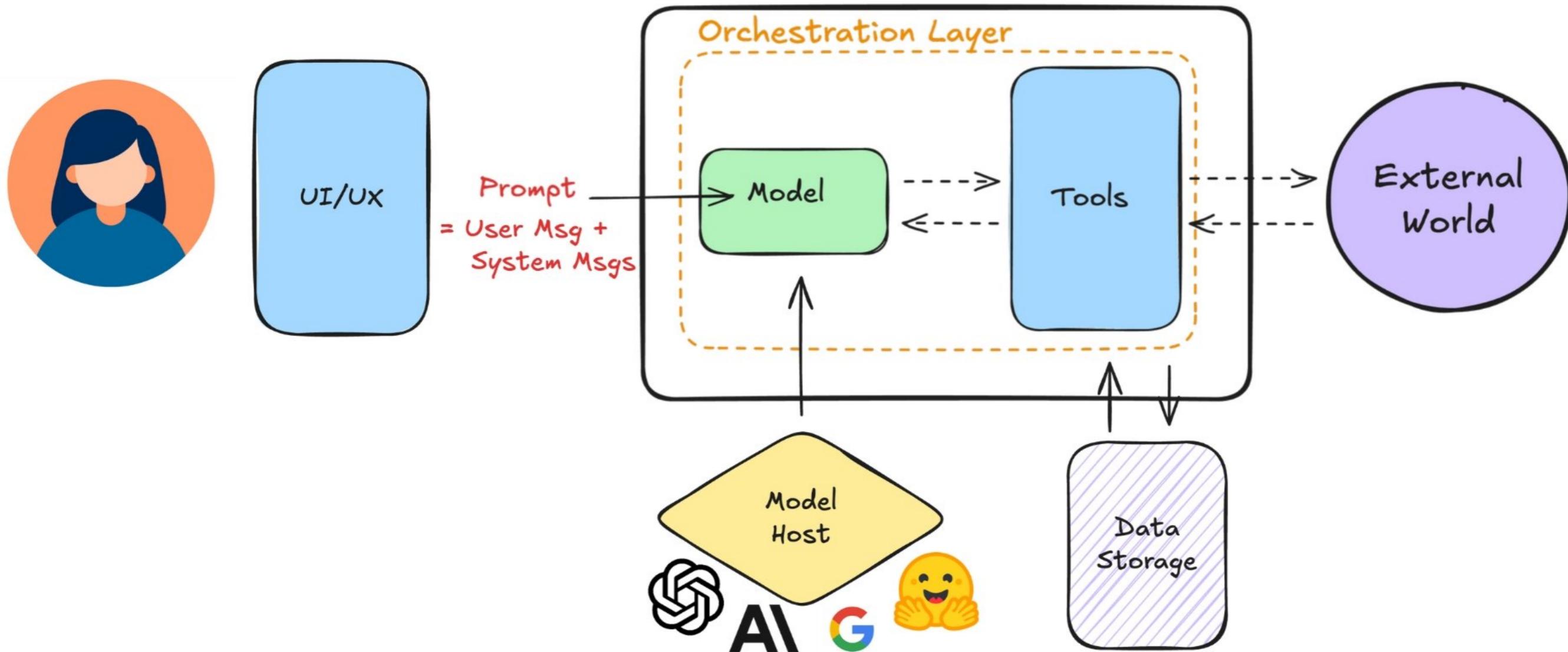
BUILDING SCALABLE AGENTIC SYSTEMS



Korey Stegared-Pace

Senior AI Cloud Advocate, Microsoft

The bigger picture



Step 1: Validate Real Interactions

Testing leads to failure, and failure leads to understanding.

— Burt Rutan

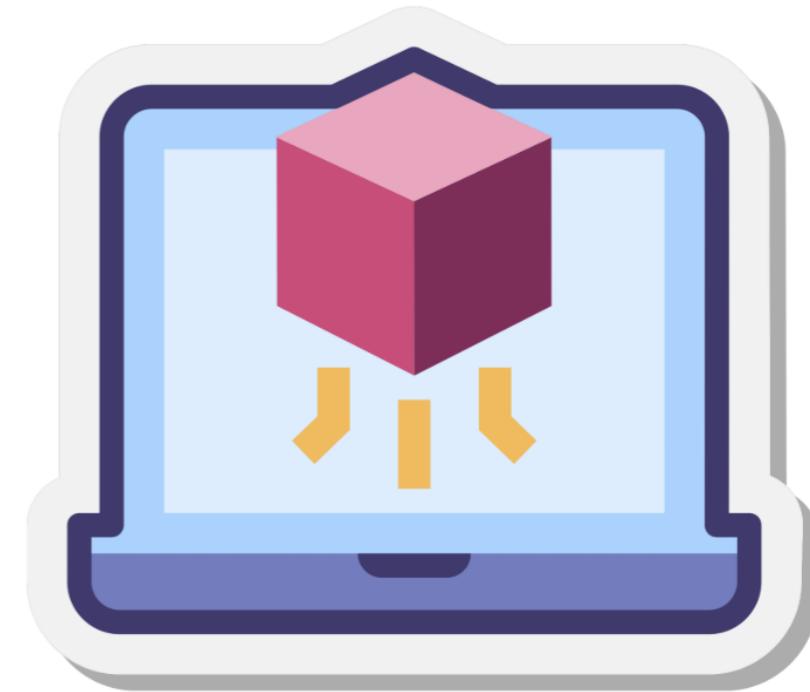


Step 1: Validate Real Interactions

NONSENSE fodsihasojdpa

SLANG I'd like to *nip* over to
Dublin from Belfast

INSULTS Work you stupid
system!!!

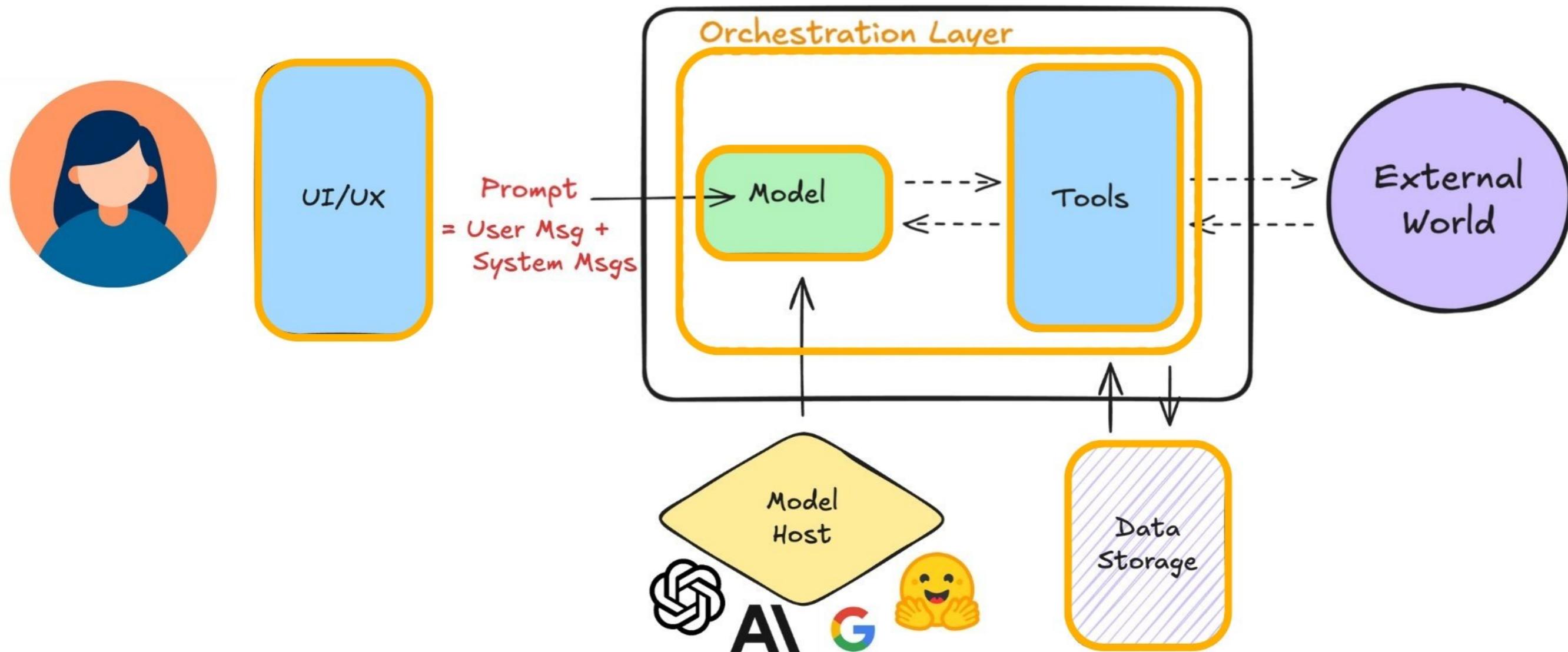


Simulated User:
fdubnfdiubfd

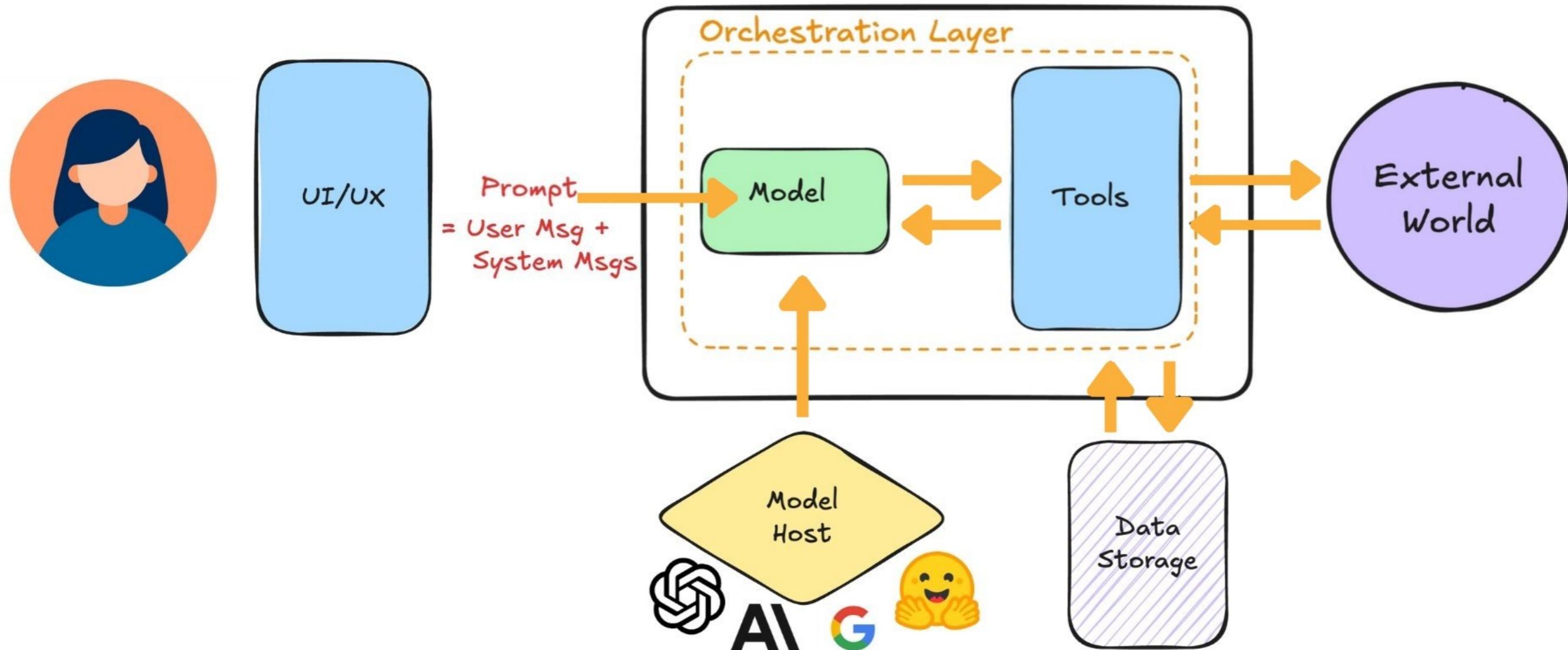
Simulated User: I'd like
to head over to Paris

¹ nip over = British slang for a quick visit

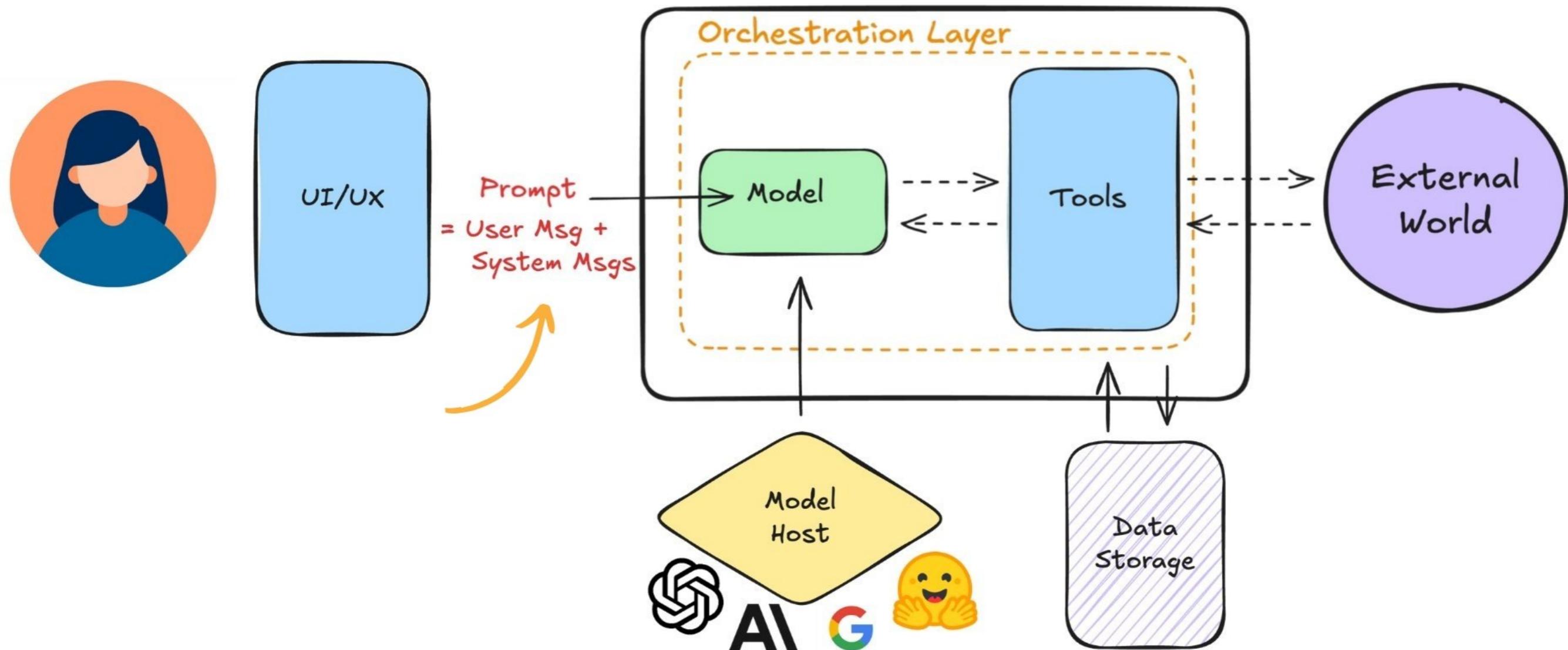
Step 2: Test Everything: Unit tests



Step 2: Test Everything: Integration tests



Step 2: Test Everything: Evaluate system prompts



Step 2: Test Everything: Subjective evaluations

HUMAN RATER



Step 2: Test Everything: Subjective evaluations

HUMAN RATER



LLM-AS-A-JUDGE



Step 3: Guardrails and Observability

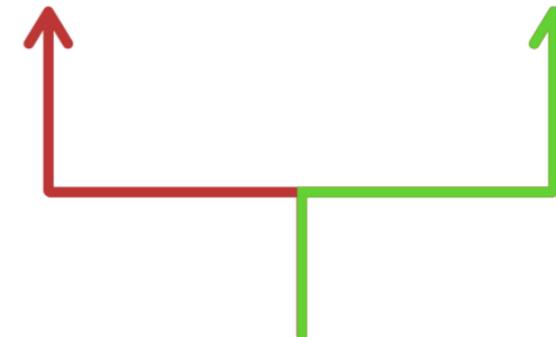
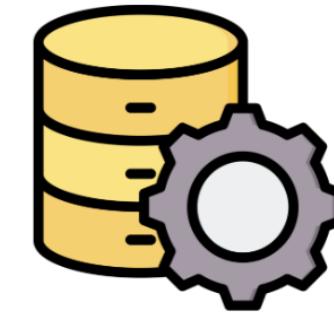
CONTENT FILTERS



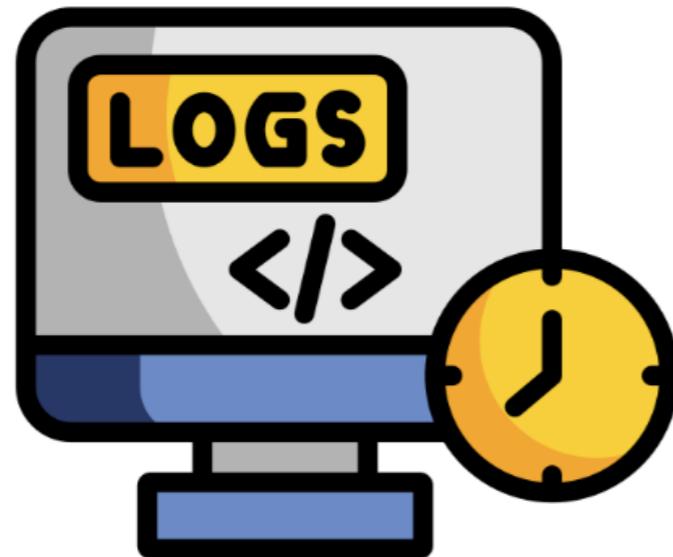
MODEL OR TOOL FAILURE



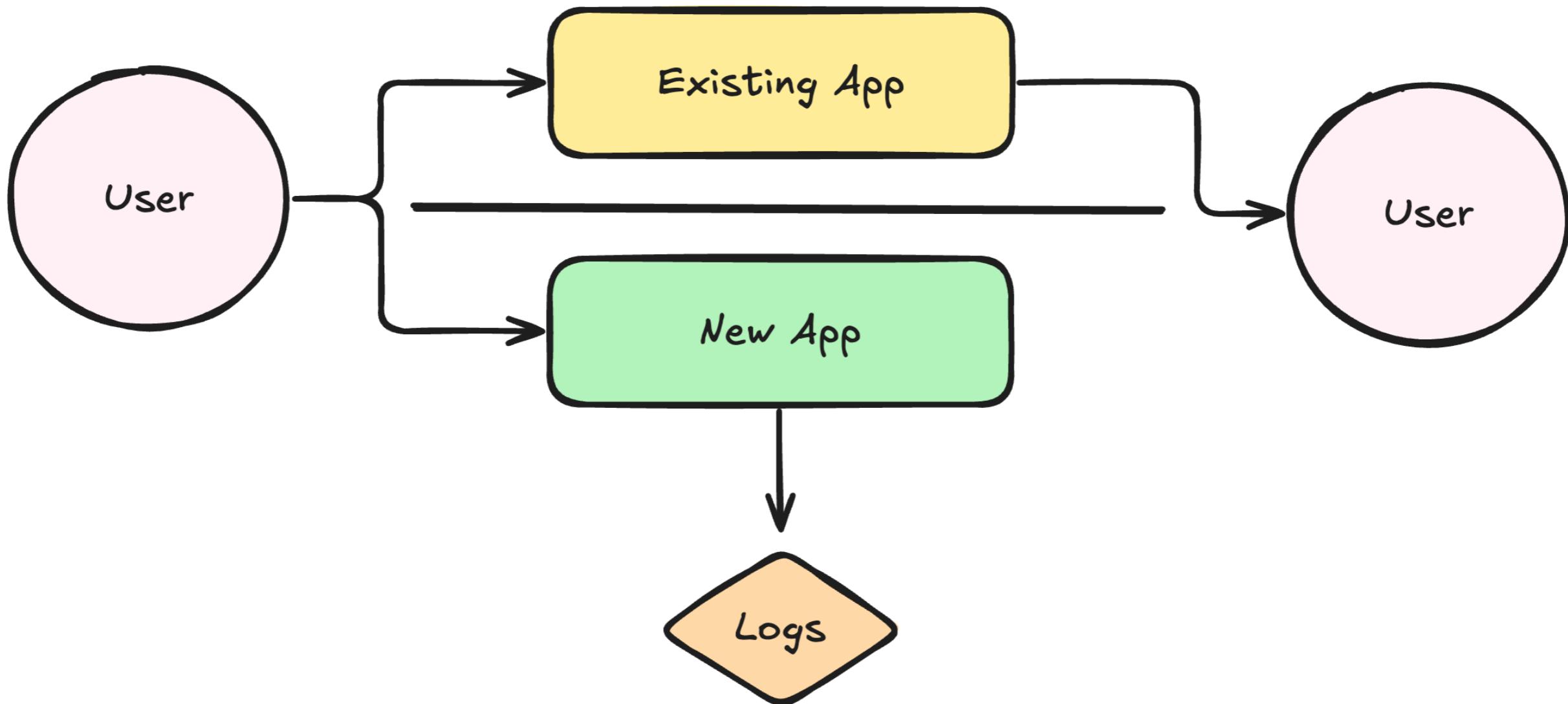
FALLBACK



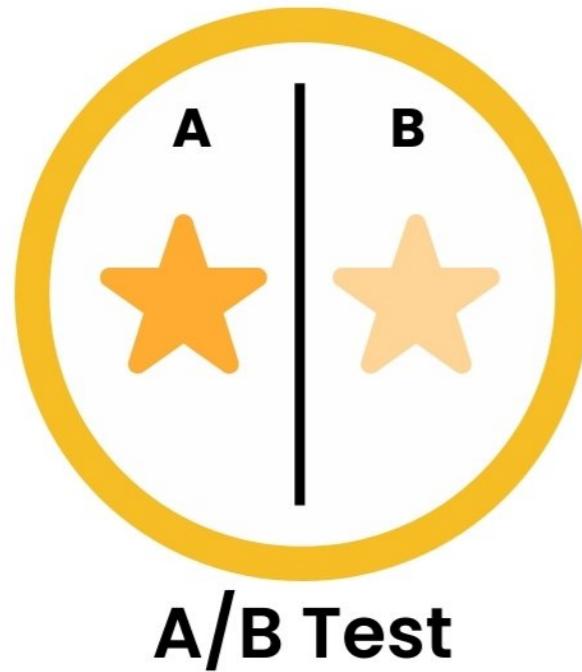
Step 3: Guardrails and Observability



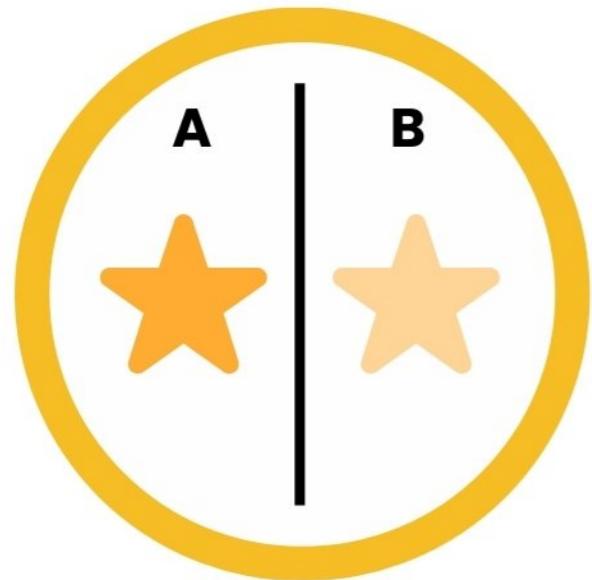
Step 4: Shadow Deployments



Step 5: Deployment Strategies



Step 5: Deployment Strategies

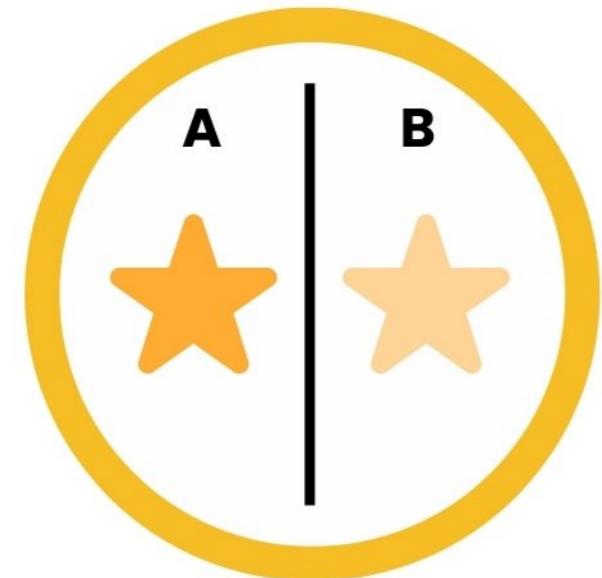


A/B Test



User Segmentation

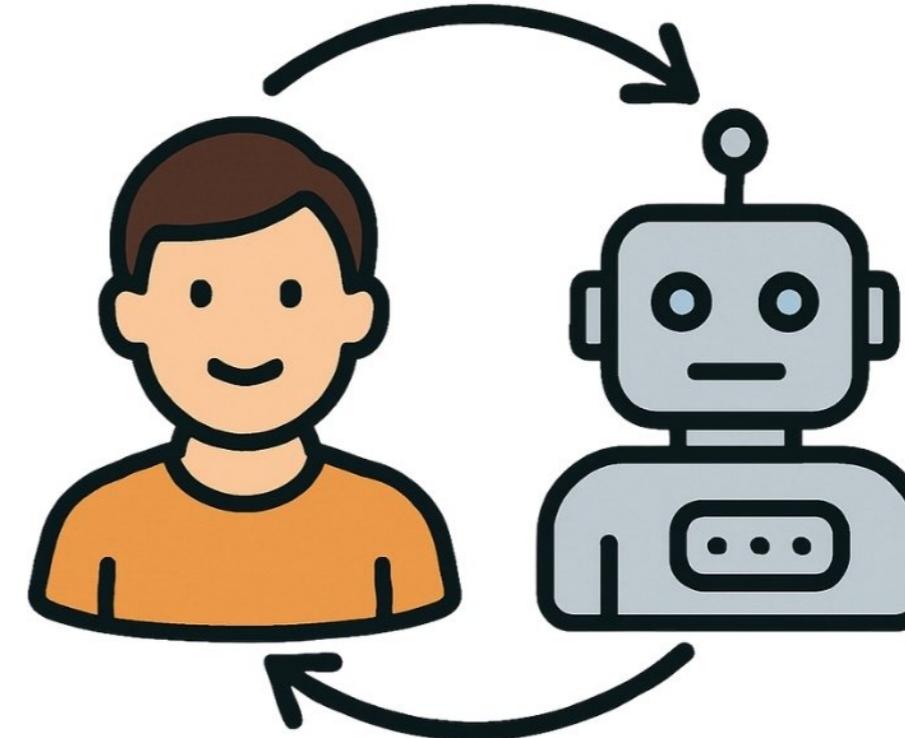
Step 5: Deployment Strategies



A/B Test



User Segmentation



Human-in-the-loop

Final Checklist



Let's practice!

BUILDING SCALABLE AGENTIC SYSTEMS

Fast and Efficient Data Ingestion

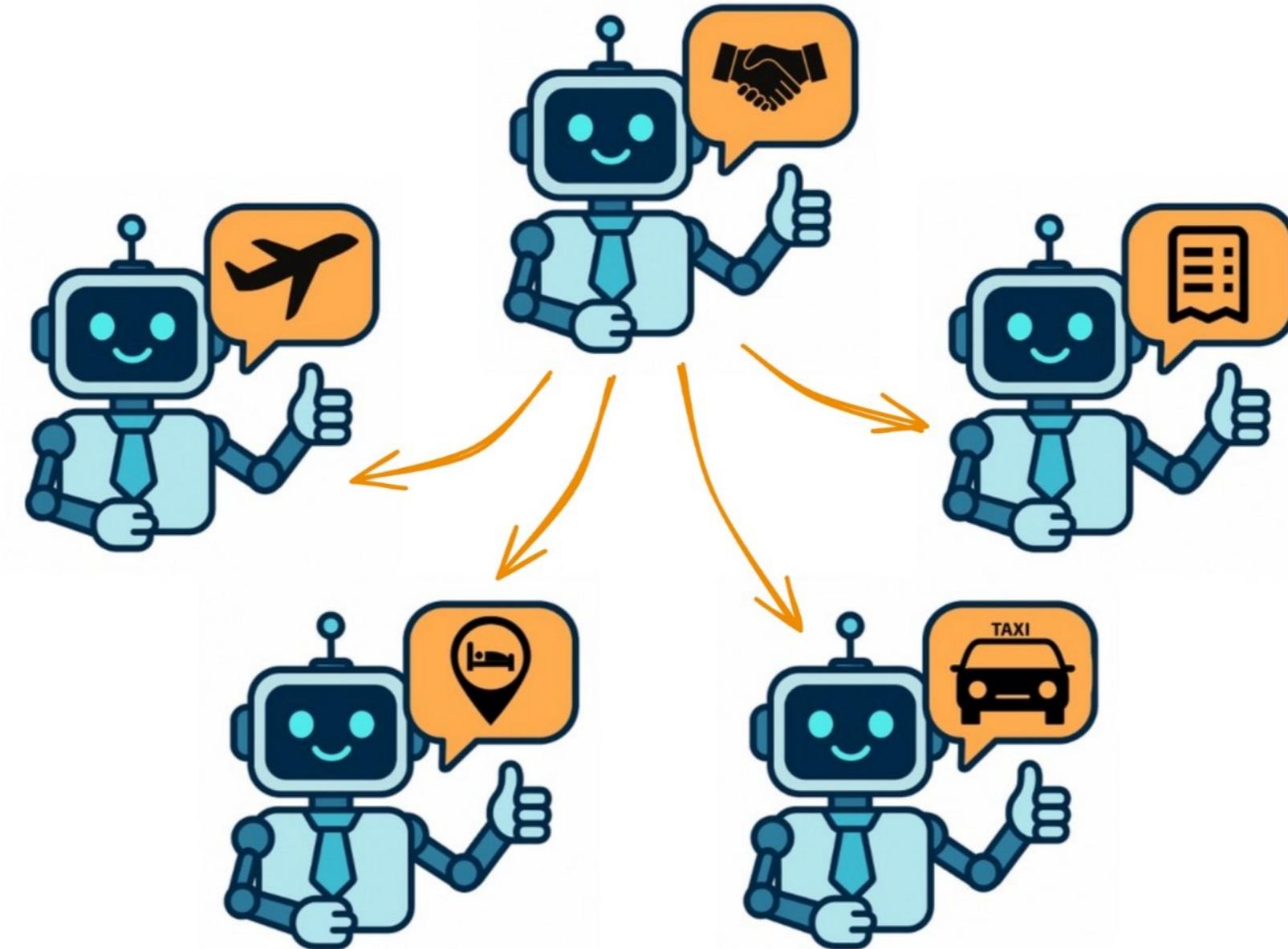
BUILDING SCALABLE AGENTIC SYSTEMS



Korey Stegared-Pace

Senior AI Cloud Advocate, Microsoft

Agents without up-to-date information



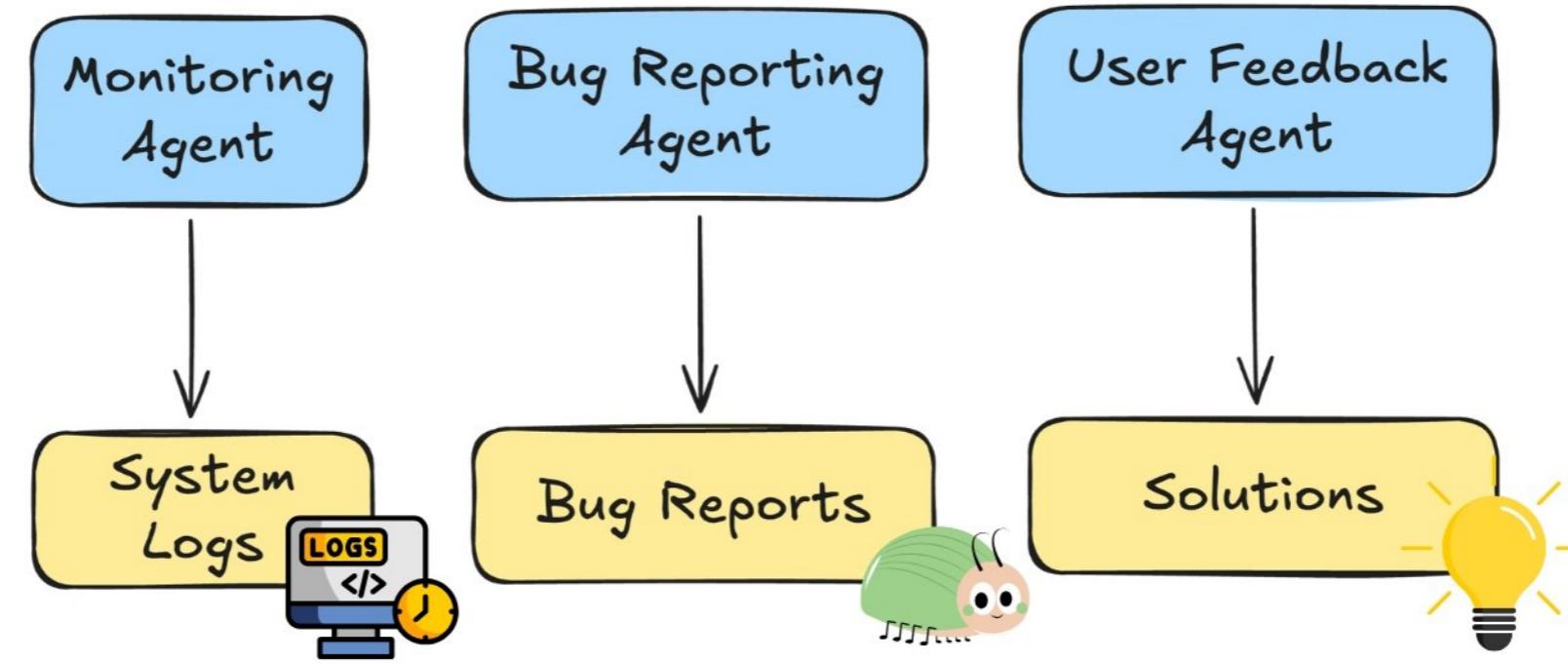
Real-time data for multi-agent coordination

Monitoring
Agent

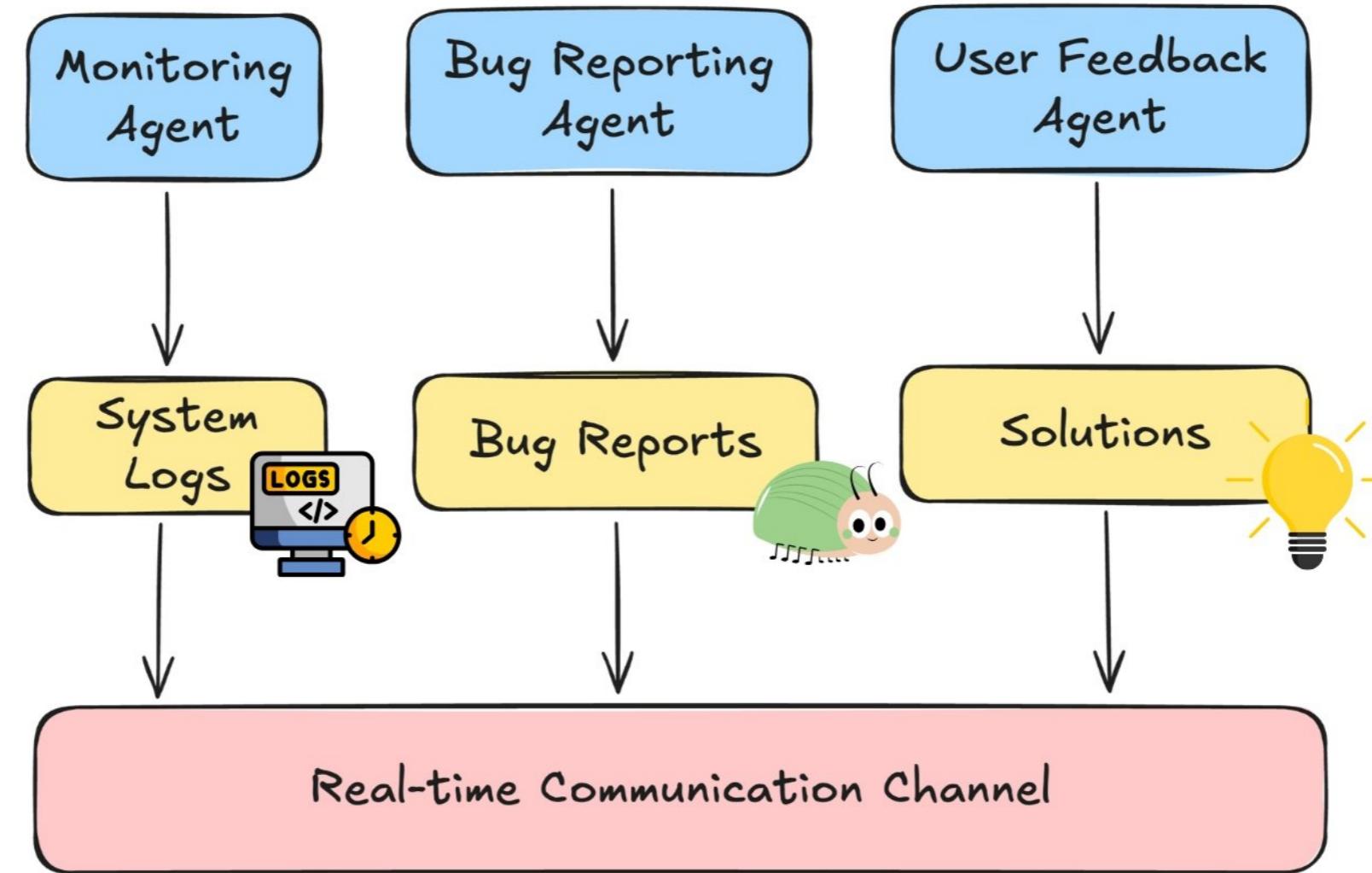
Bug Reporting
Agent

User Feedback
Agent

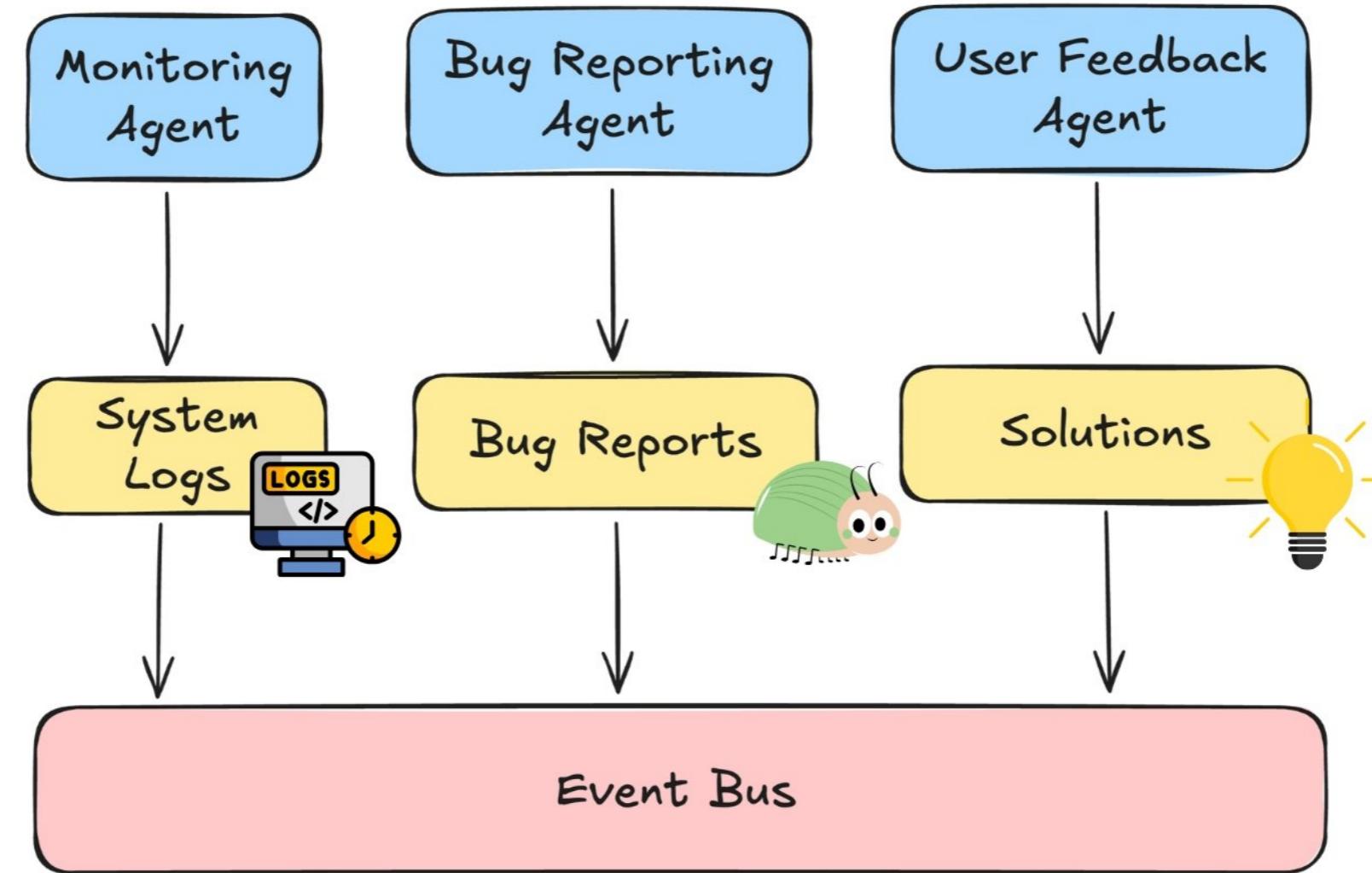
Real-time data for multi-agent coordination



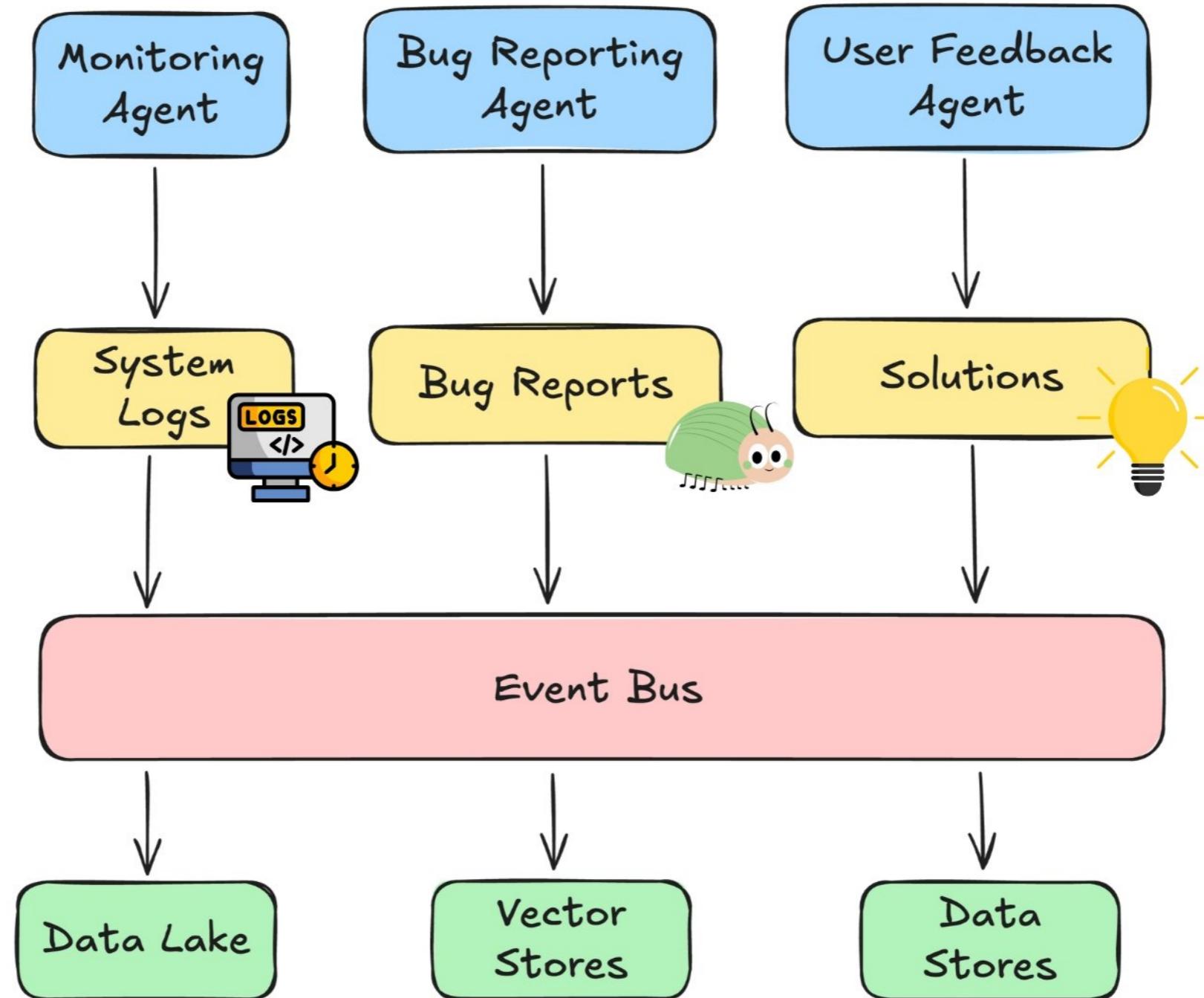
Real-time communication channel



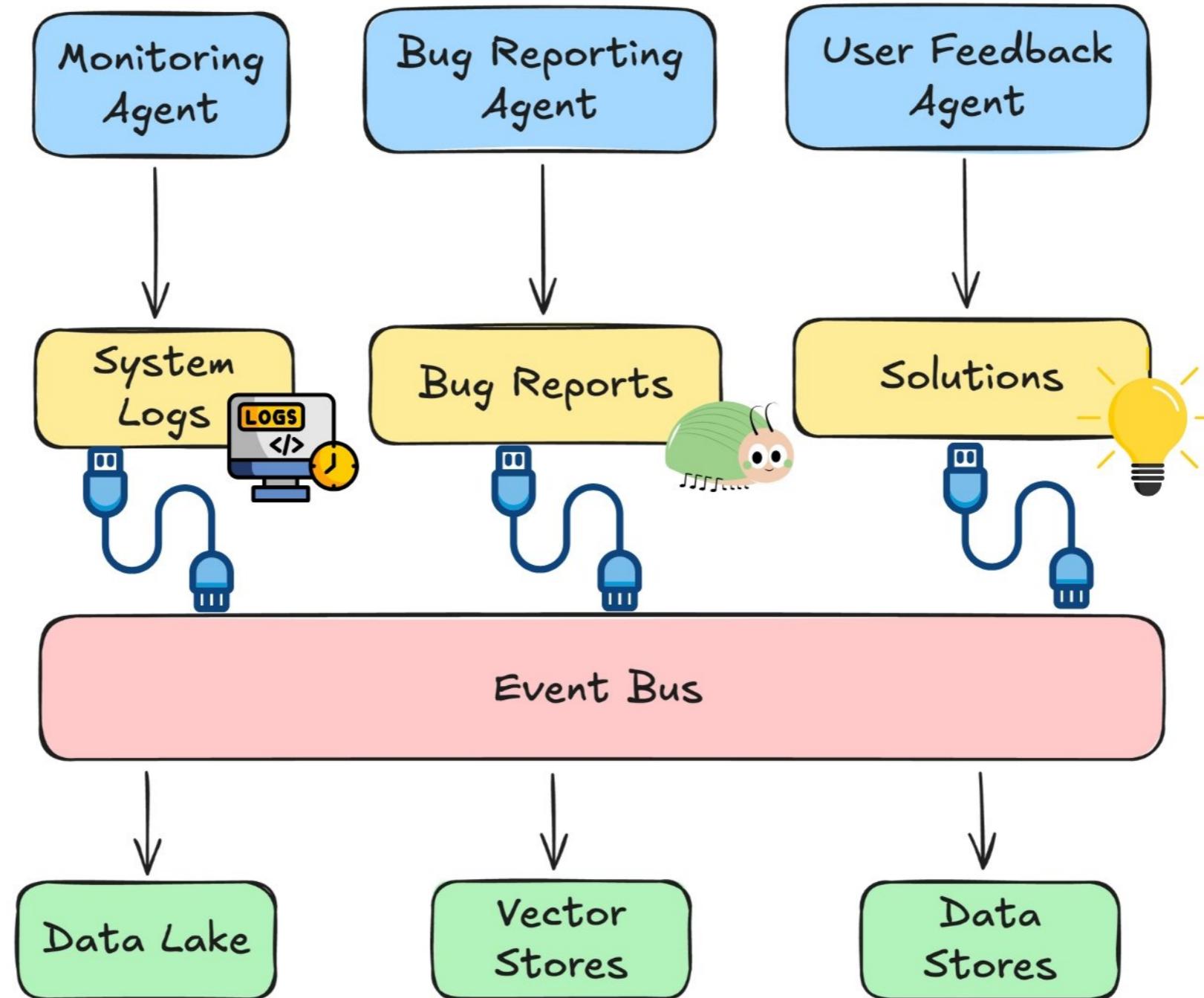
Real-time communication channel



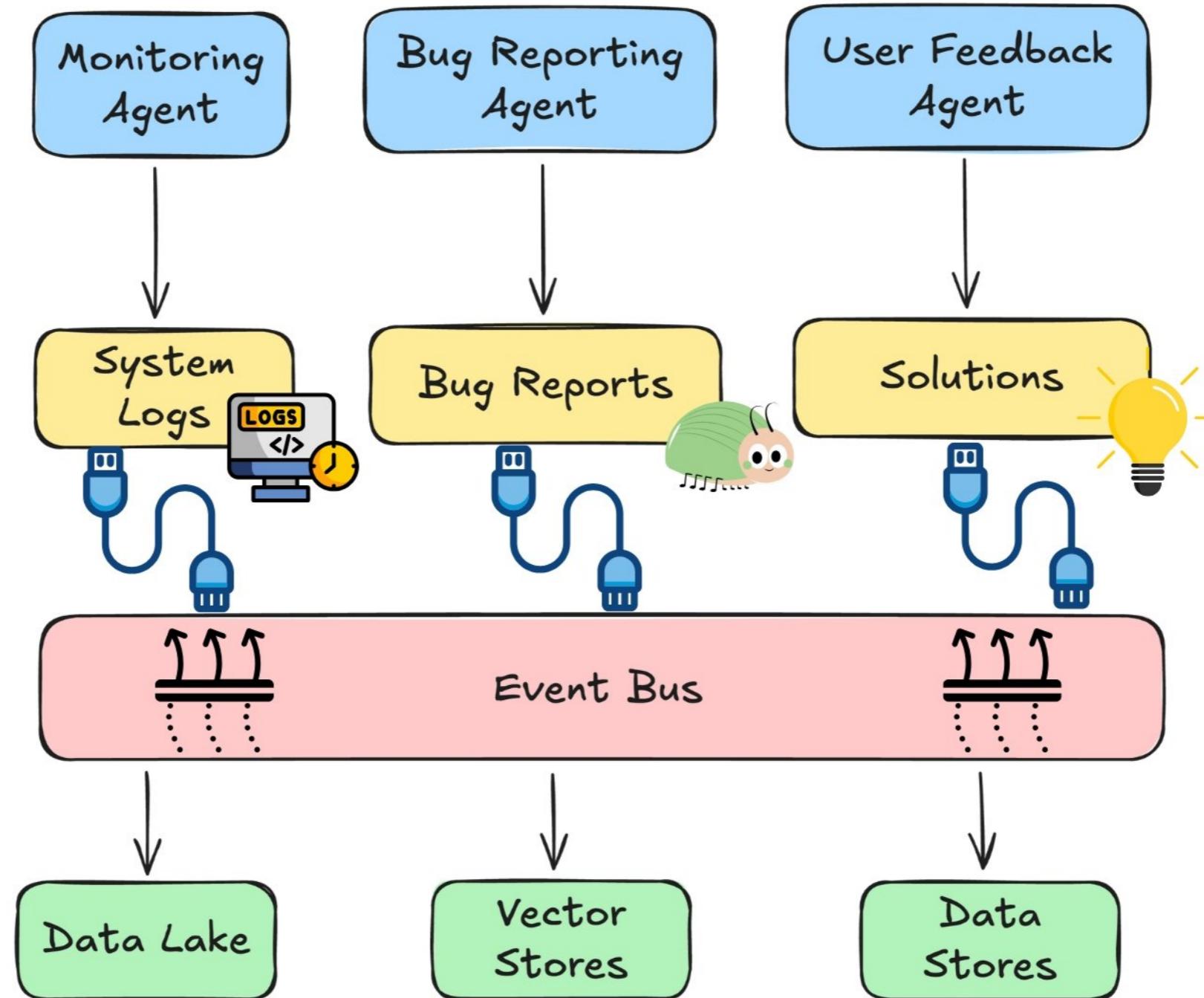
Utilizing an event bus



Utilizing an event bus



Utilizing an event bus



Let's practice!

BUILDING SCALABLE AGENTIC SYSTEMS

Failing Gracefully: Mitigating Risks

BUILDING SCALABLE AGENTIC SYSTEMS



Korey Stegared-Pace

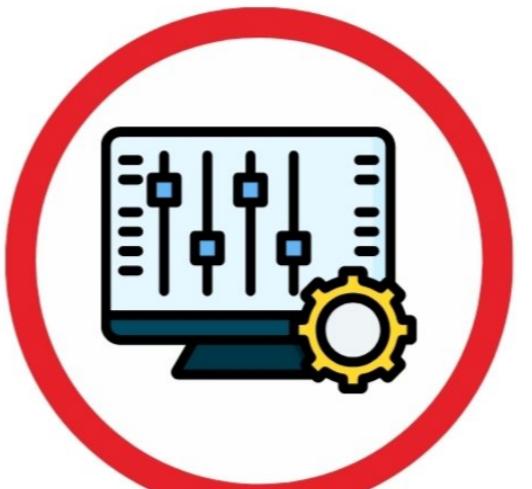
Senior AI Cloud Advocate, Microsoft

Tool call failures

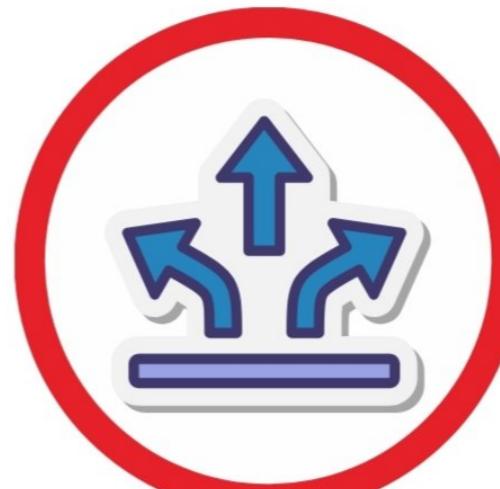
Solutions

- Defining clear parameters and usage
- Tool output validation (unit testing!)
- Verification checks on tool selection
- MCP!

Problems



Parameters



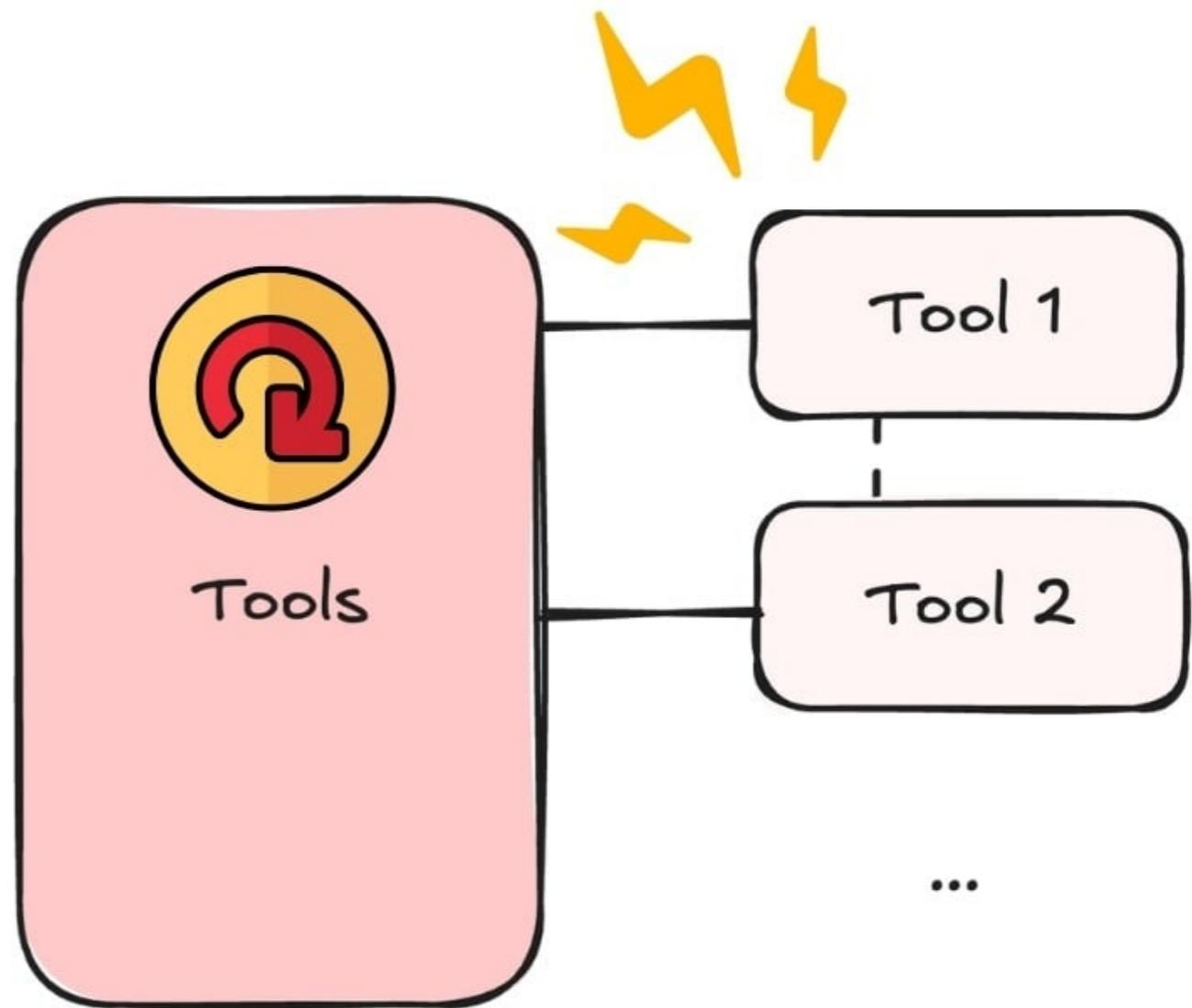
Outputs



Tool Selection

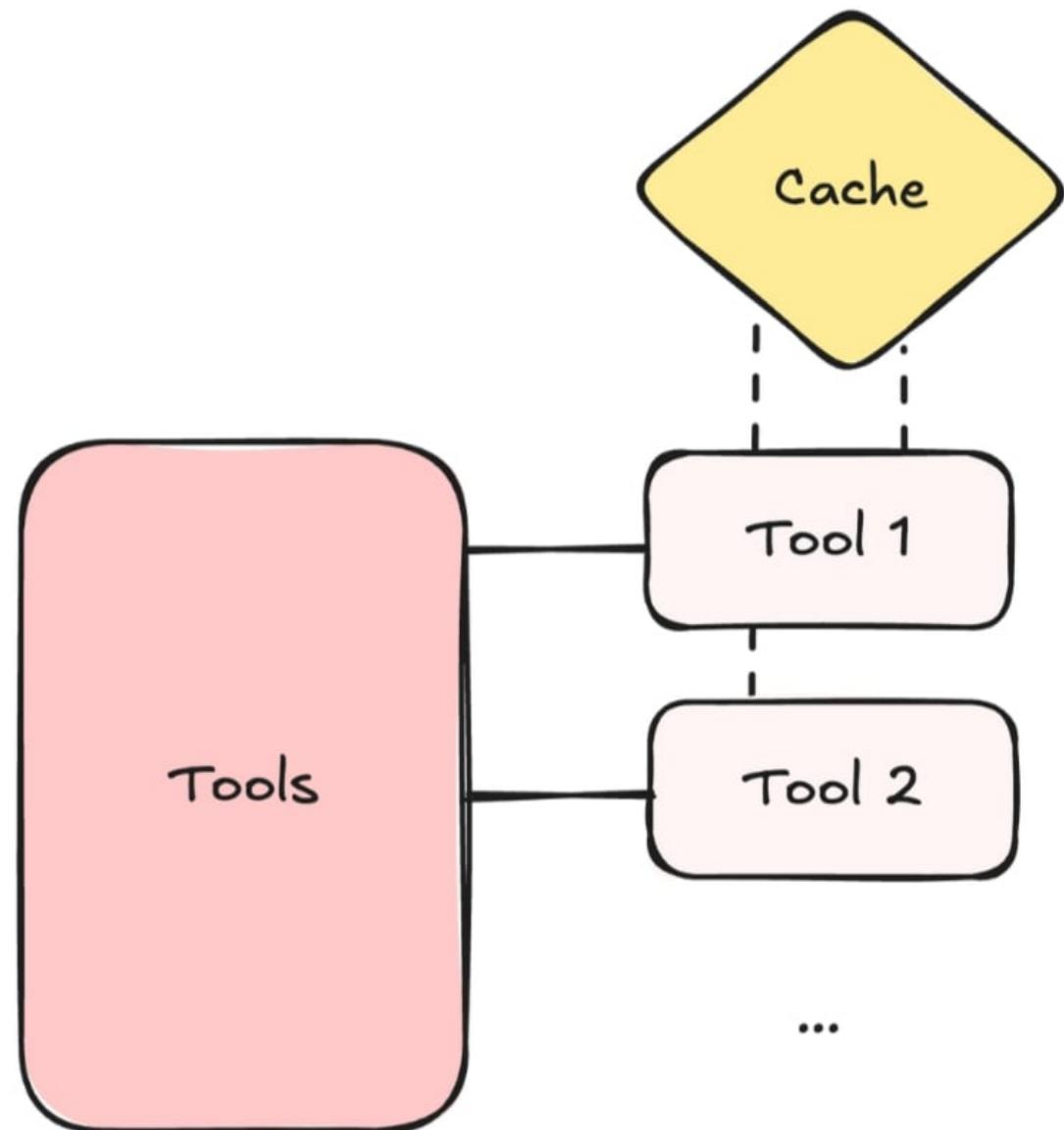
Tool call failures - retry mechanism

- External services may *break* or be temporarily *slow*
- **Retry**
 - **Exponential backoff** → *slowly decreasing* retries
 - Inform users of issues using *callbacks*



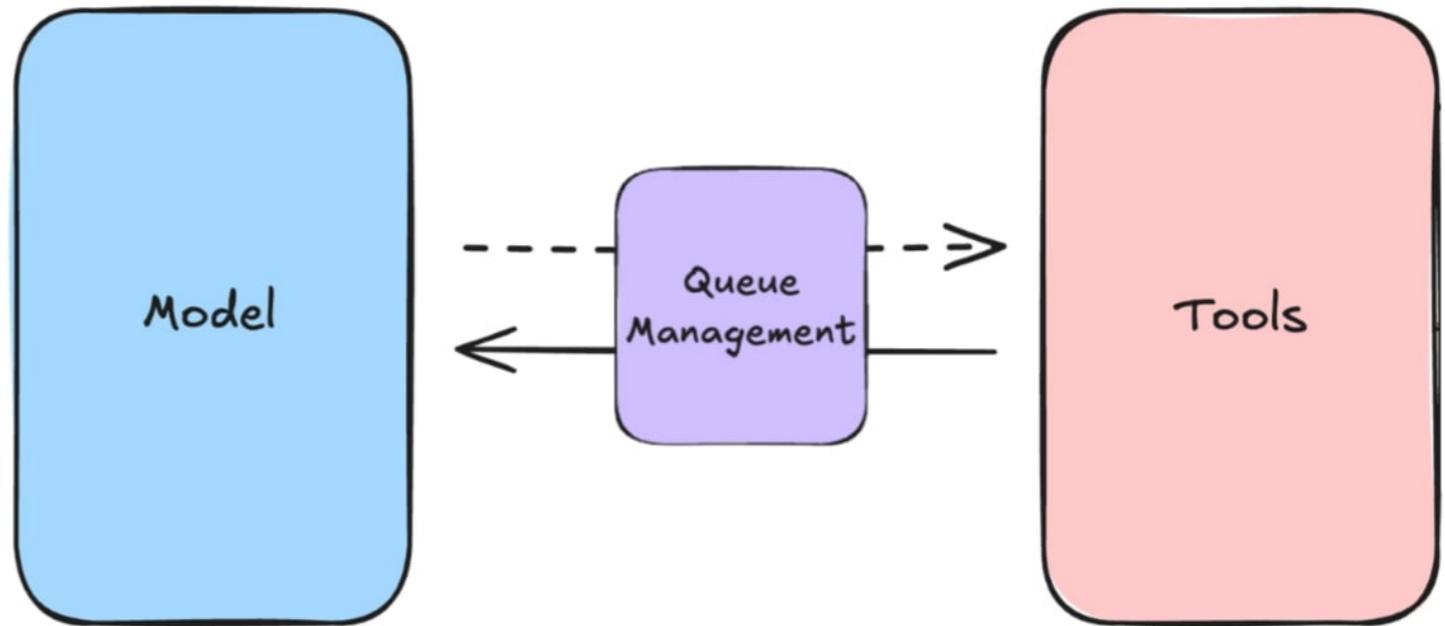
Tool call failures - caching

- Tool outputs are cached as a *fallback*
- **Works for:** Static data with infrequent updates
- **Won't work for:** tools that serve real-time information

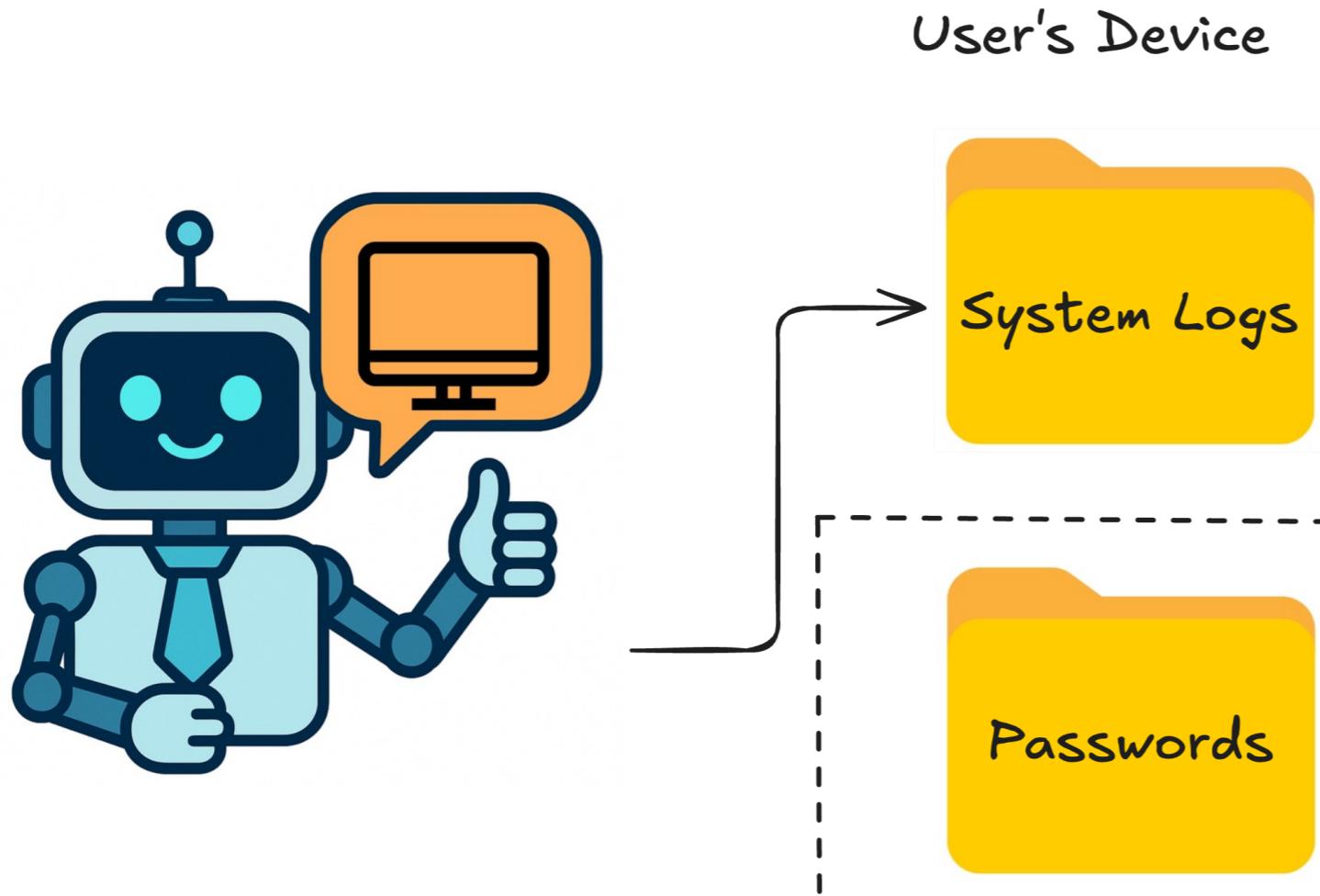


Tool call failures - queue management

- Tool calls may be *reliant* on each another
- Example: flights must be booked before the taxi
- Unsuccessful tool calls can be moved down the queue



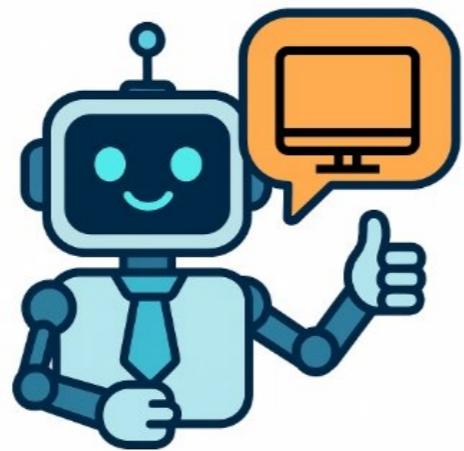
Authentication



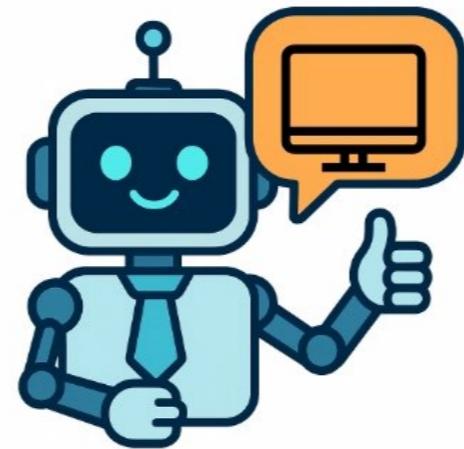
- Tools may require *access* and *permissions* to **private** data
- **Example:** IT support agent
 - Valid access: system logs
 - Blocked access: passwords, location, etc.

Authentication - unique agent identifiers

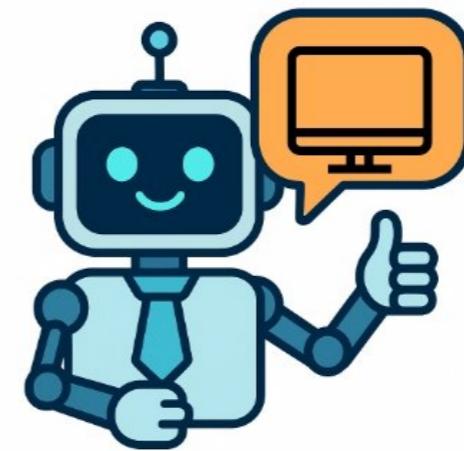
ID: x784



ID: x783

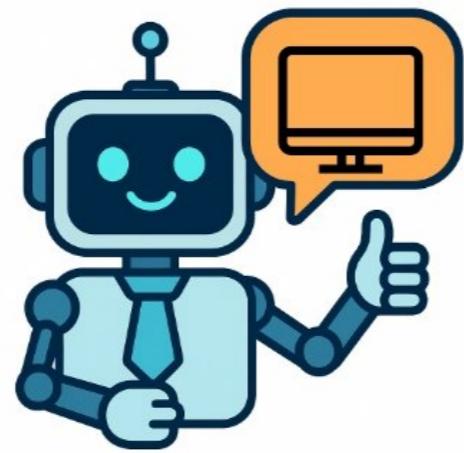


ID: x782



Authentication - isolated environments

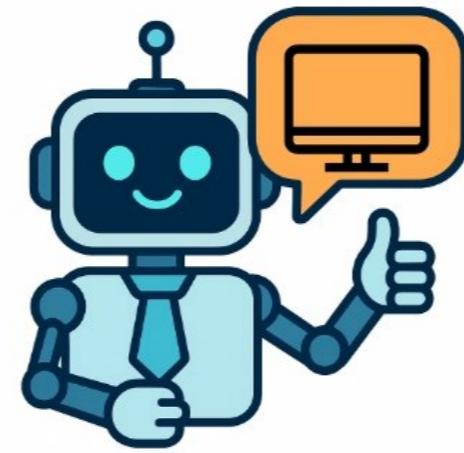
ID: x784



Access Rights:

...

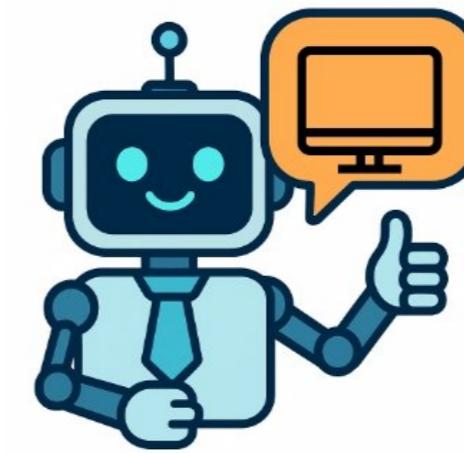
ID: x783



Access Rights:

...

ID: x782



Access Rights:

...

Authentication - guardrails and action restraints



Let's practice!

BUILDING SCALABLE AGENTIC SYSTEMS

Congratulations!

BUILDING SCALABLE AGENTIC SYSTEMS

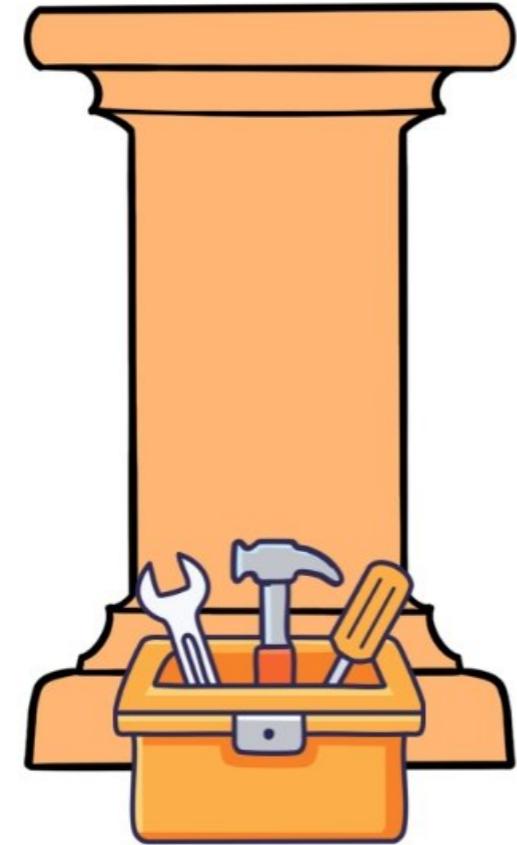


Korey Stegared-Pace

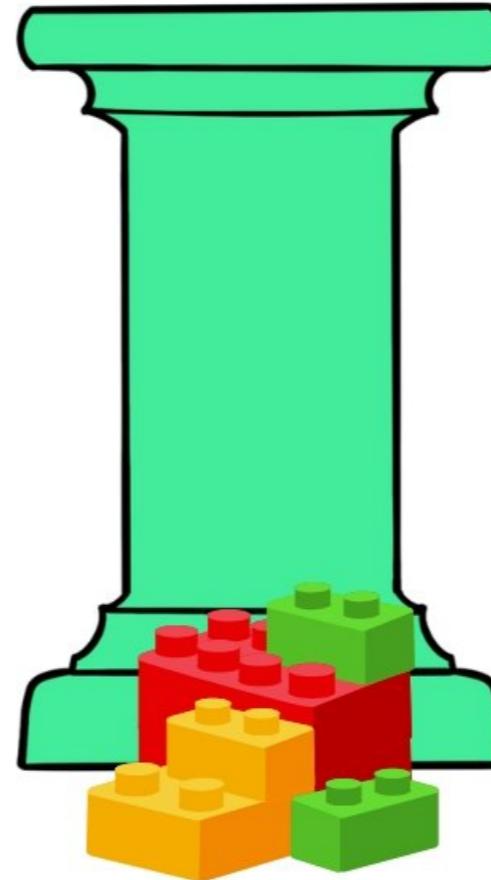
Senior AI Cloud Advocate, Microsoft

Chapter 1

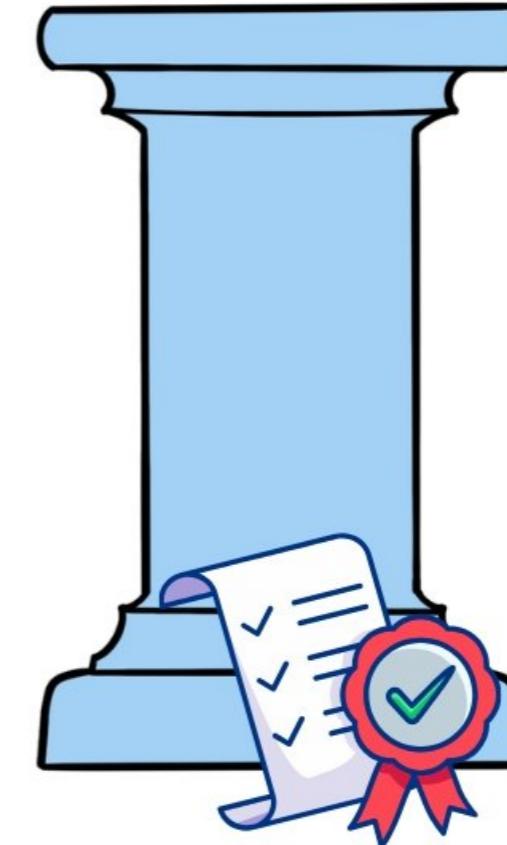
**Robust
Infrastructure
and Tooling**



**Modular
Design
Architecture**



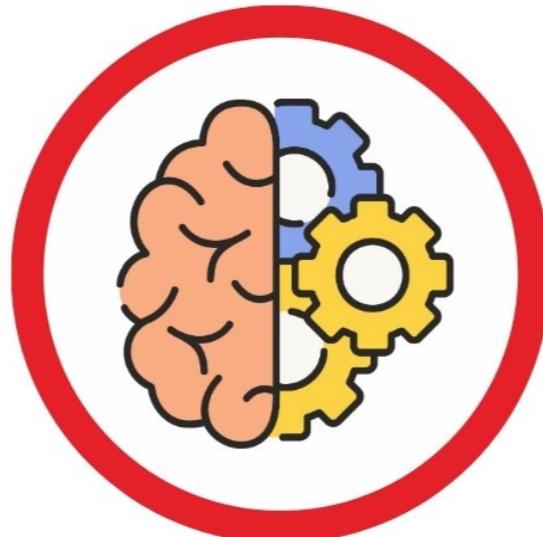
**Continuous
Evaluation &
Feedback Loops**



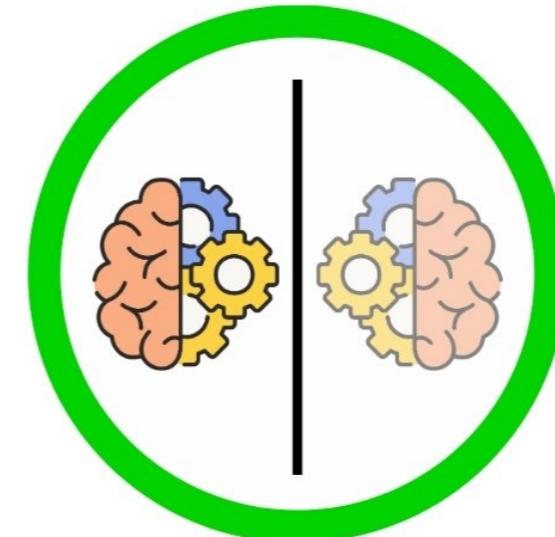
Chapter 1



Tool Use



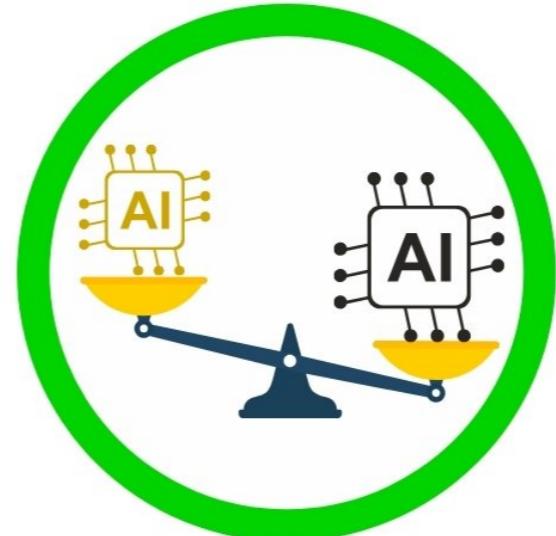
Reasoning



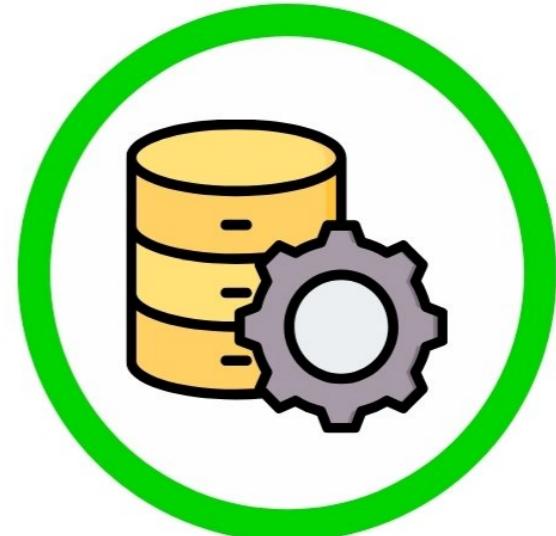
Selective Reasoning



Retrieval



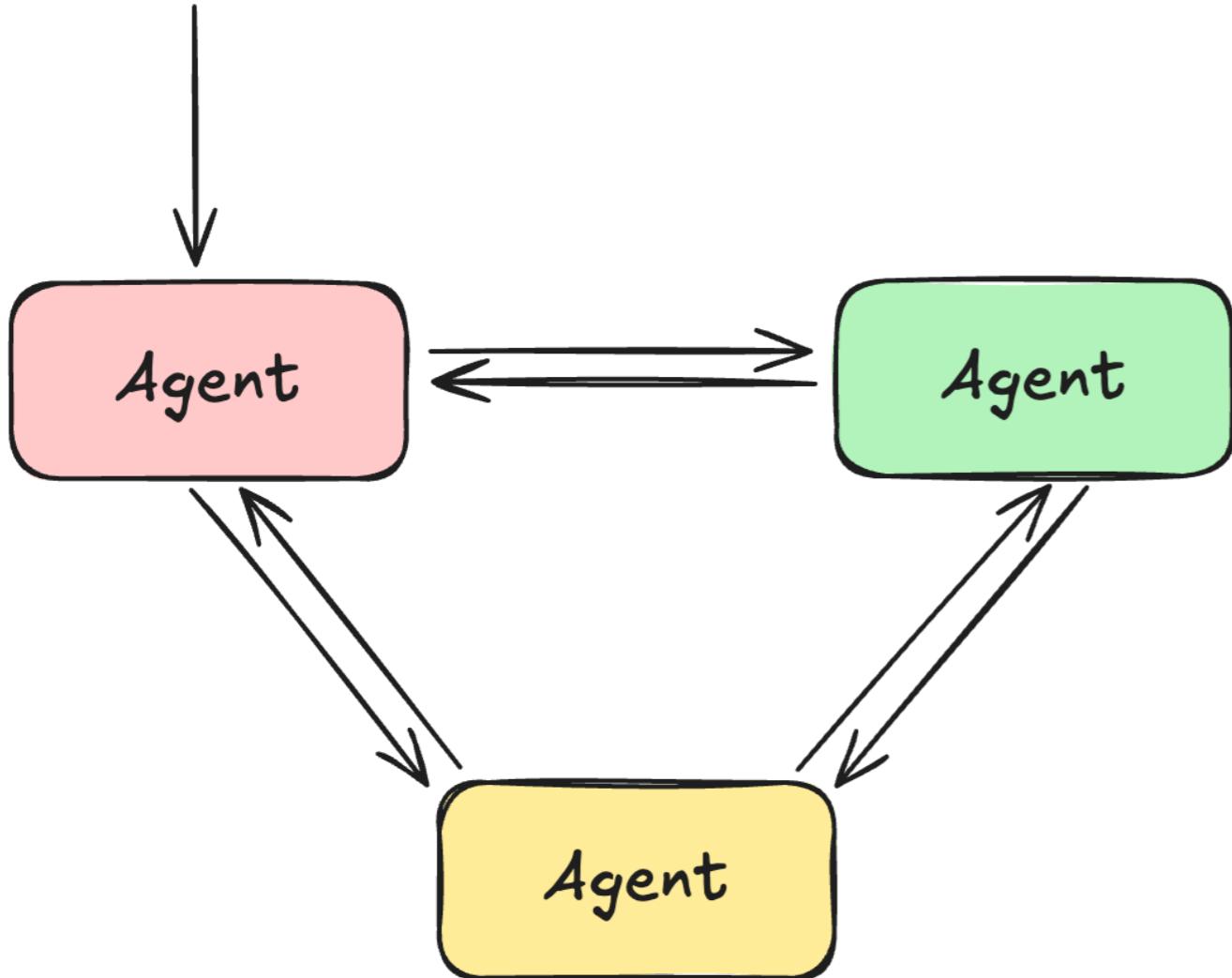
Lighter Models



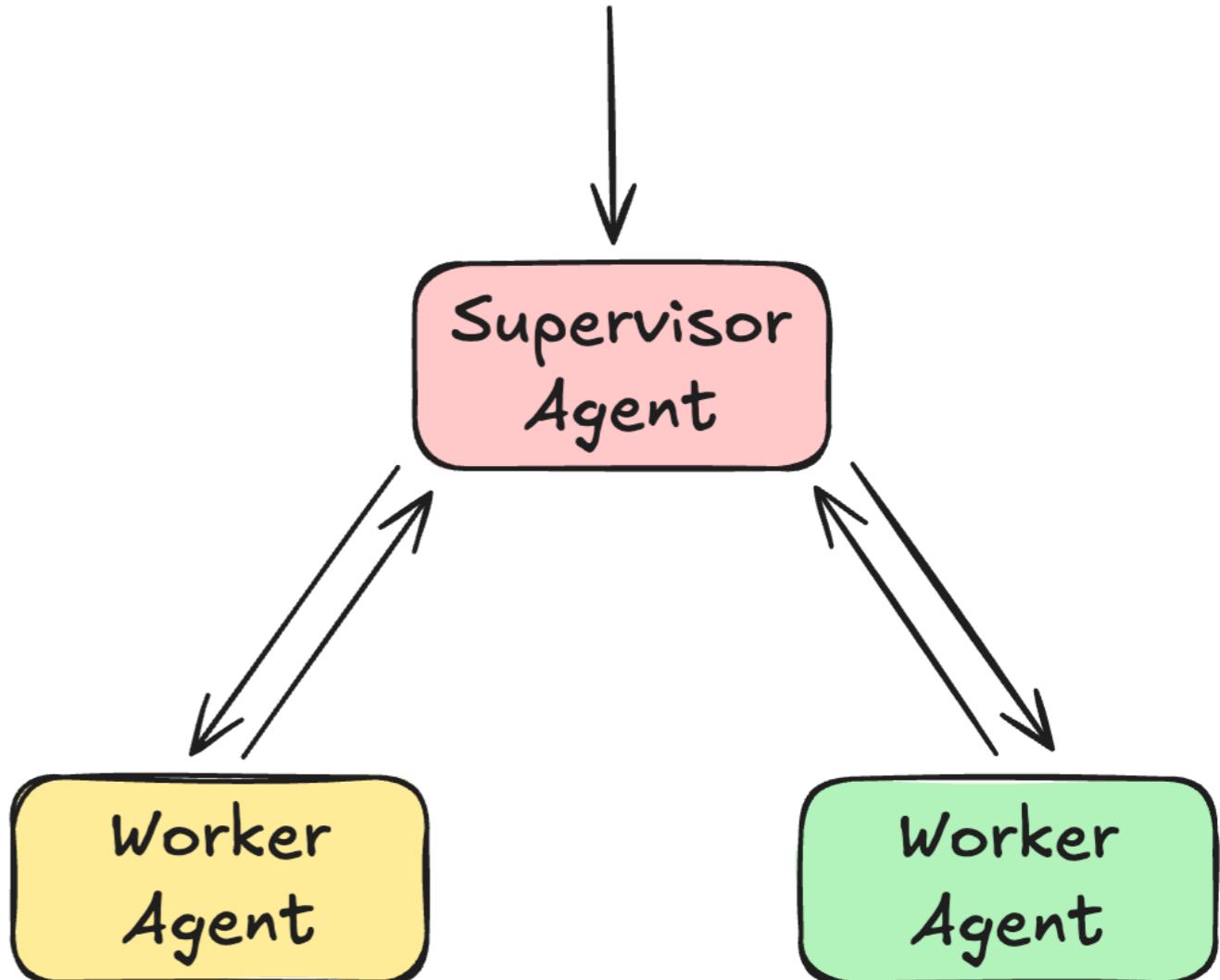
Cache

Chapter 2

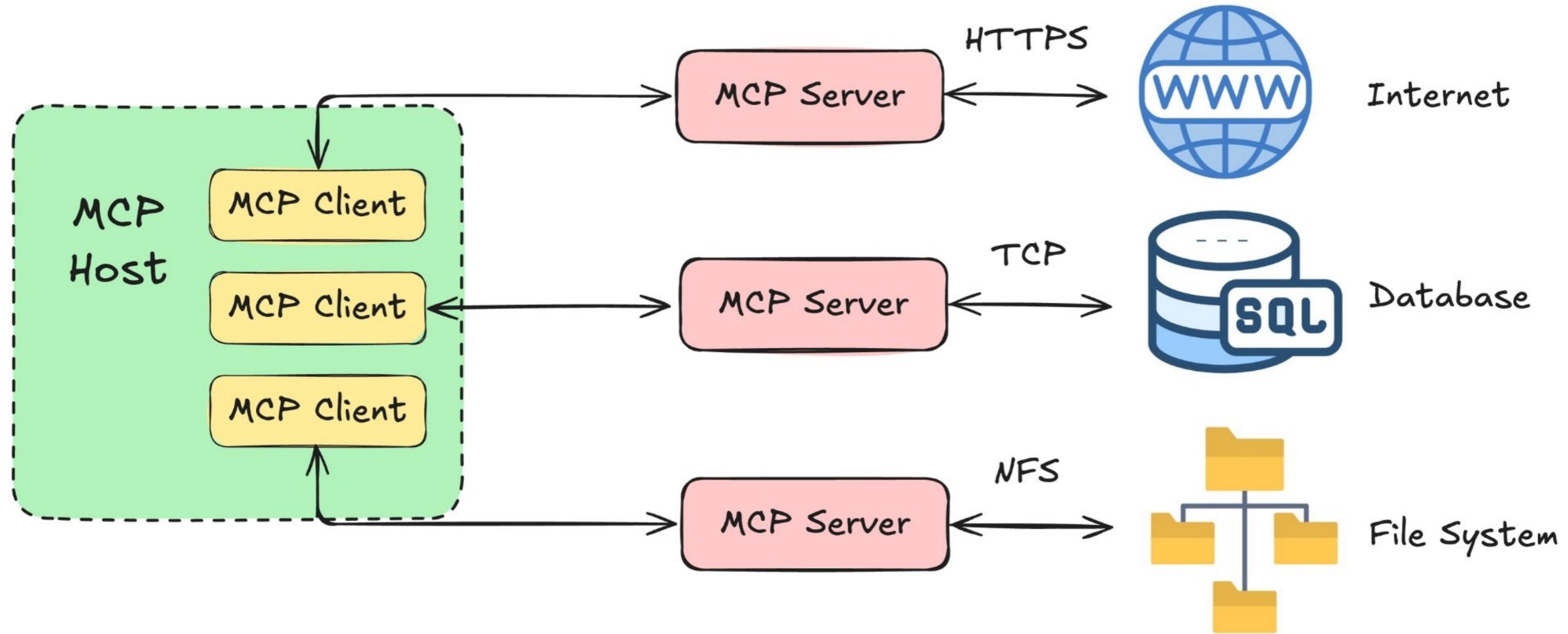
Network



Supervisor



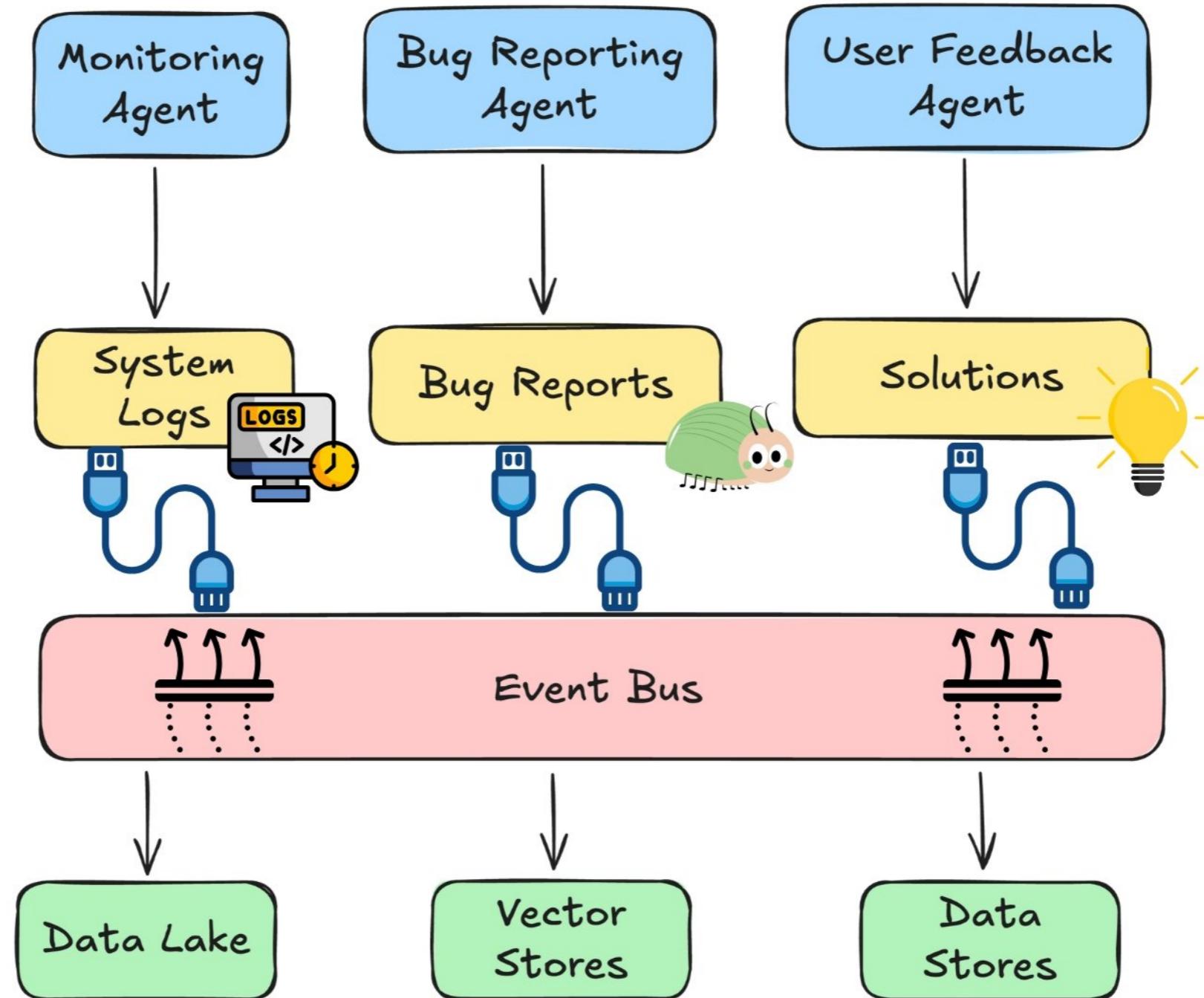
Chapter 2



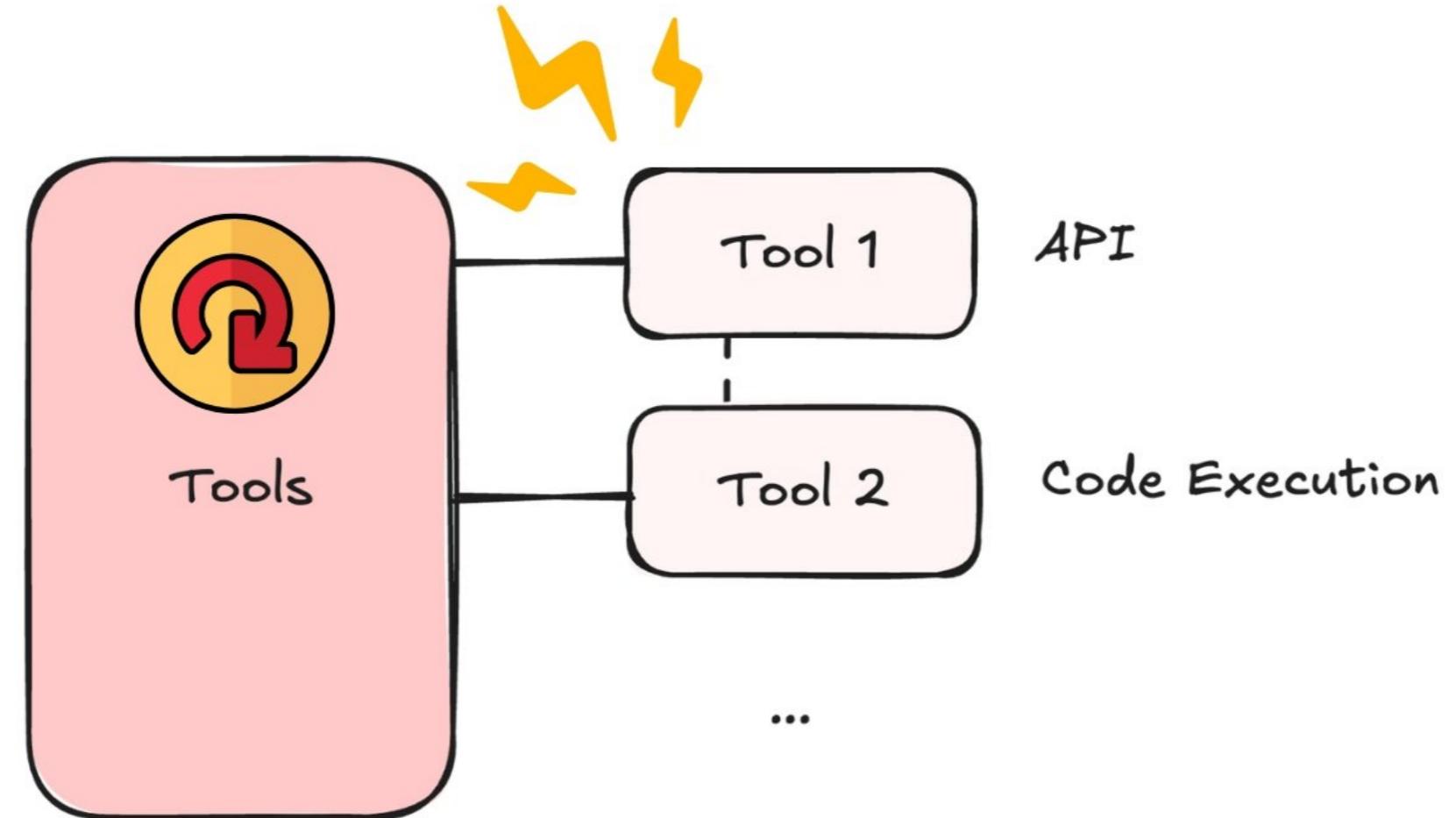
Chapter 3



Chapter 3



Chapter 3



Where next?

- Blog: [Understanding AI Agents: The Future of Autonomous Systems](#)
- DataFramed Podcast: [AI Agents Hype vs. Reality](#)
- [AI Agents Cheat Sheet](#)

Let's practice!

BUILDING SCALABLE AGENTIC SYSTEMS