# AI is here. And it's changing everything.
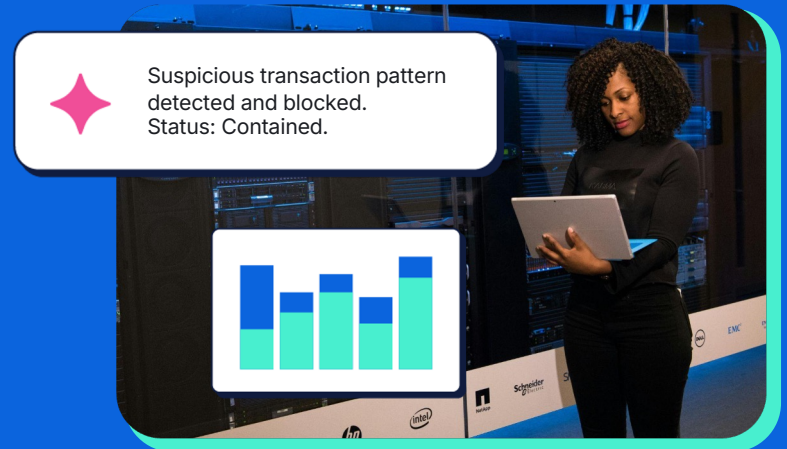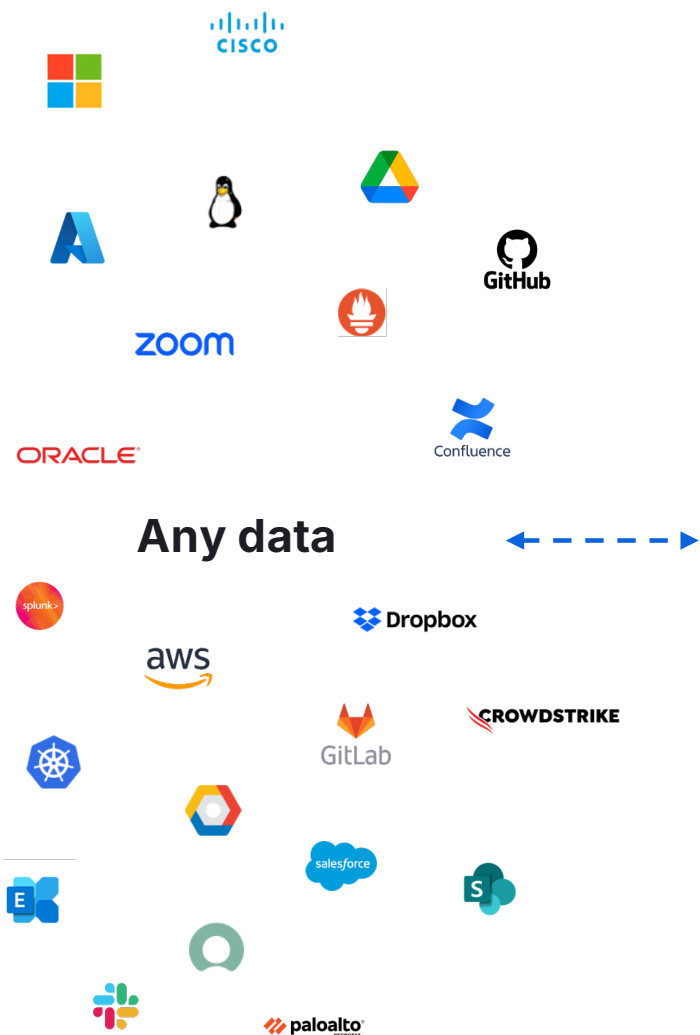
In the workplace, in our personal lives, and more ...



**People** driving the process with AI as assistants & copilots



Suspicious transaction pattern detected and blocked.
Status: Contained.

**AI agents** driving the process with people reviewing

elastic

**Any data**

Success depends on providing AI **high-quality, relevant** access to **your data and powerful tools**

**AI experiences**

elastic

# Search relevance and context has never been **more important**



Widespread
**Search Powered Applications**

Current
**RAG**

Emerging
**Agents**

Human judged results

One relevant answer

Iterative, looping retrievals and LLM calls to accomplish task.

Importance of Relevance

elastic

# Improving Agent outcomes through Context Engineering

**Possible context to give the model (Potentially Bns of tokens)**

Sys Prompt    Docs    Memory

Tools    Database    Chat Hist

**Agent**

**Curation**

**Large Language Model call**

**Curated Context**     **Task**

**Limited to 100k -> 10M tokens**

AI experiences depend on a limited LLM context window.

Dynamically delivering the right data for the LLM's task is called **Context Engineering**.

elastic

# Agentic AI

# Agentic AI

- **Agentic AI** refers to AI systems that can act autonomously to achieve goals, making decisions and taking actions without direct human intervention
- An **agent** is usually a software that makes decisions and takes actions often by calling tools or APIs in a loop, based on its current understanding of the world and its objective
- Key features of an agent:
  - **Goal-oriented**
  - **Autonomous**
  - **Interactive**
  - **Iterative**: reason → act → observe → repeat
  - **Memory** (optional)

elastic

# Agent workflow example

- **Goal:** Book a flight to Paris under $500.
  - **Plan (reason)**: "Check available flights."
  - **Act**: Calls search_flights(to="Paris", max_price=500)
  - **Observe**: Gets a list of flights.
  - **Reason**: "There's a flight on Tuesday under budget. Book it!"
  - **Act**: Calls book_flight(flight_id)
  - **Finish**: Reports the result back to the user.
- The agent monitors its own process, deciding what to do next based on each result

# Limitations

- **Execution time:** not really fast since it requires many steps
- **Expensive**: agents typically use about 4x more tokens than chat interactions, multi-agent use about 15x more tokens than chat
- **Complexity**: multi-agent are good with parallel tasks but not so good to manage many dependencies between agents*
- **Unpredictability**: agent interactions can lead to unexpected outcomes. Needs of Guardrail systems.
- **Evaluation**: it's difficult to measure the collective success or alignment of multiple agents

* How we built our multi-agent research system, Anthropic
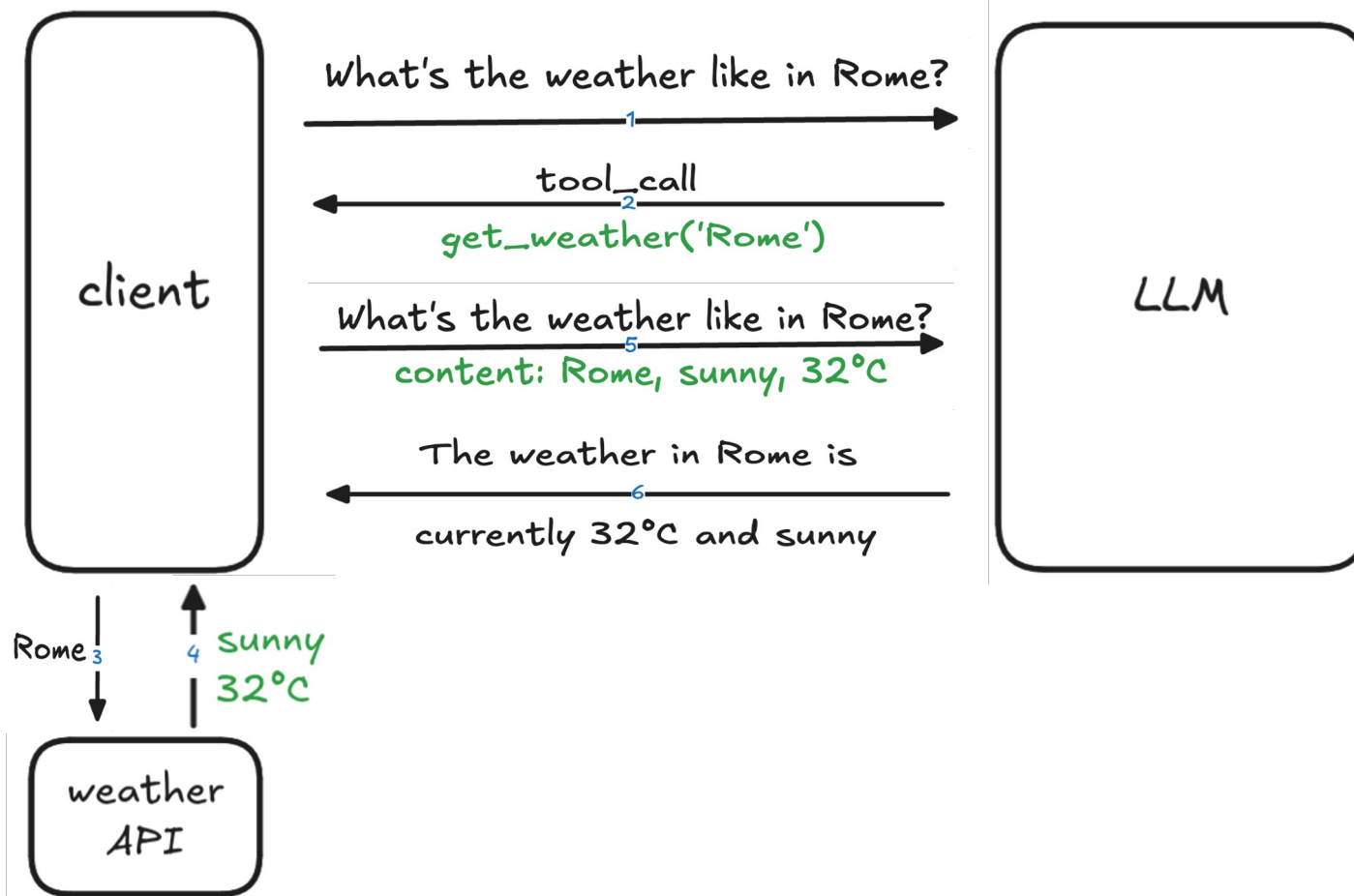
elastic

# Tool calling

# History of tool calling

- **Tool calling** (or function calling) is an emerging property in LLM
- Relevant papers that investigated the topic:
  - Nakano et al., [WebGPT: Browser-assisted question-answering with human feedback](#), OpenAI, 2022
  - Timo Schick et al., [Toolformer: Language Models Can Teach Themselves to Use Tools](#), Meta AI Research, 2023
  - [Function calling and other API updates, OpenAI](#), June 13, 2023

elastic

# Tool calling (or Function calling)

- Tool calling is the ability of LLM to recognize the need to execute external functions (tools) as part of its reasoning process
- The LLM recognizes when it needs of additional information or actions and request the usage of tools (preparing the generation a function call in JSON function)
- The client is responsible for executing the function call (not the LLM) and this step is usually monitored by a human

# A diagram of Tool calling

# Tool calling in OpenAI

POST  https://api.openai.com/v1/chat/completions

```
{
    "model": "gpt-4.1",
    "messages": [
        {
            "role": "user",
            "content": "What is the weather like in Rome today?"
        }
    ],
    "tools": [
        {
            "type": "function",
            "function": {
                "name": "get_weather",
                "description": "Get current temperature for a given location.",
                "parameters": {
                    "type": "object",
                    "properties": {
                        "location": {
                            "type": "string",
                            "description": "City and country e.g. Rome, Italy"
                        }
                    },
                    "required": [
                        "location"
                    ],
                    "additionalProperties": false
                },
                "strict": true
            }
        }
    ]
}
```

## Response

```
[{
    "id": "call_12345xyz",
    "type": "function",
    "function": {
        "name": "get_weather",
        "arguments": "{\"location\":\"Rome, Italy\"}"
    }
}]
```
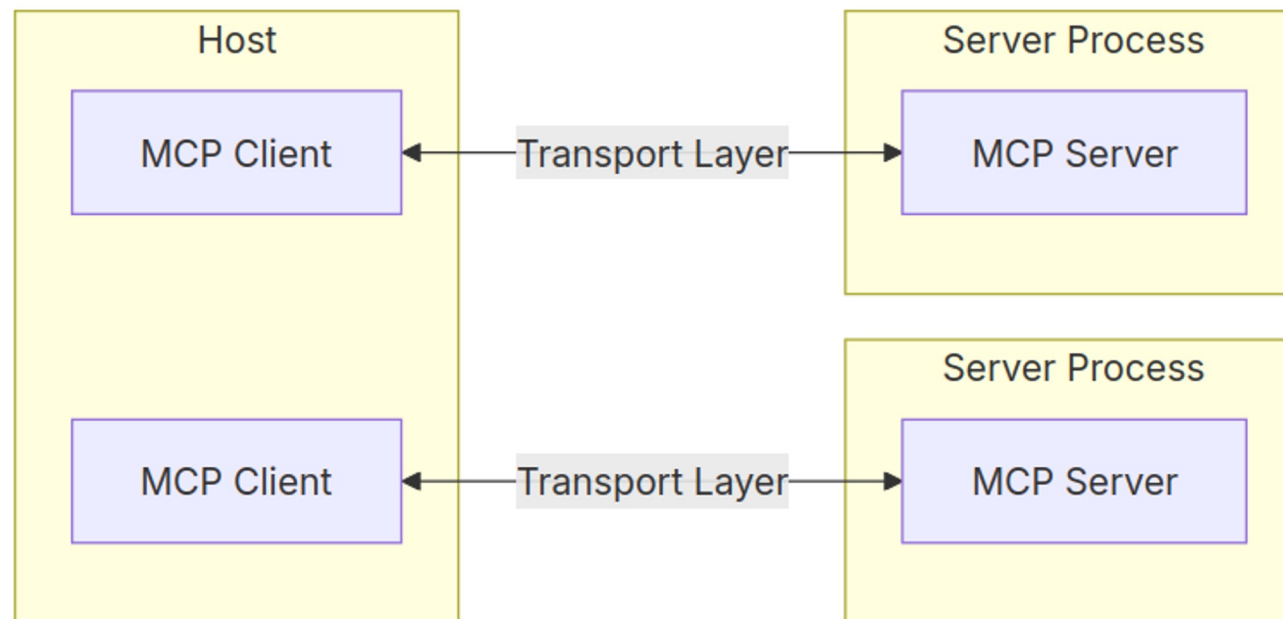
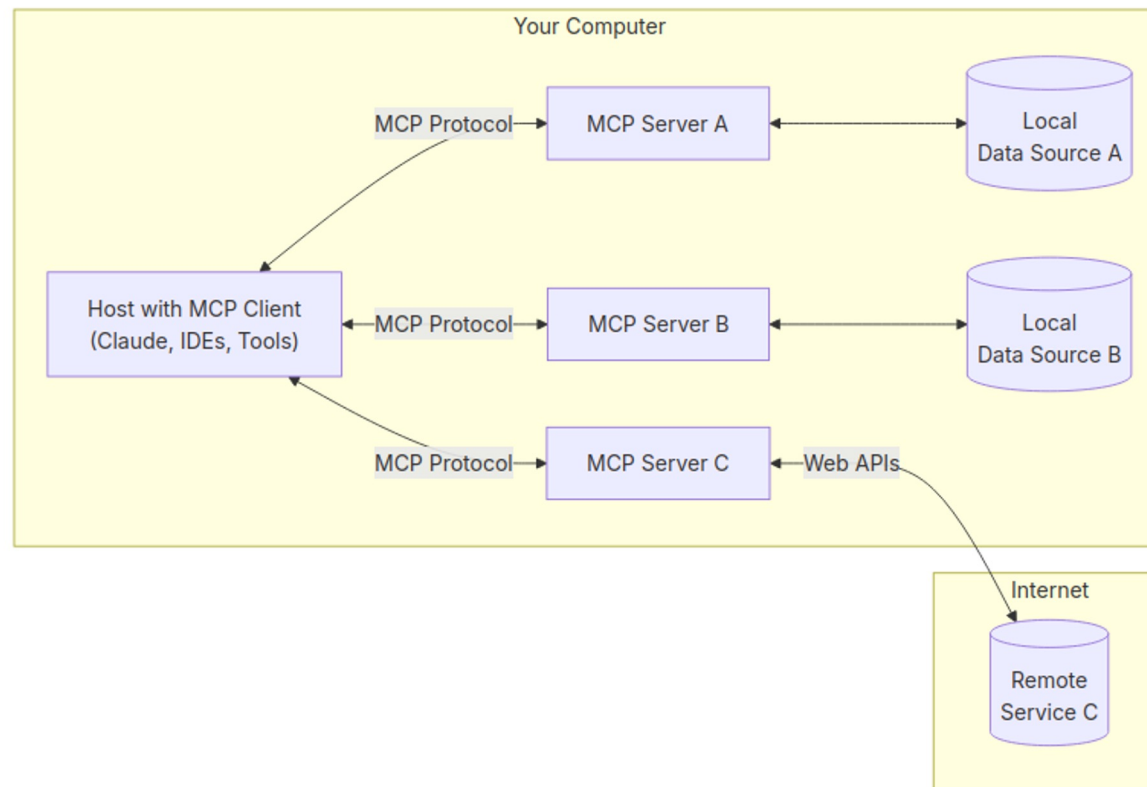# Model Context Protocol

# Model Context Protocol

- Model Context Protocol (MCP) is an open protocol that standardizes how applications provide context to LLMs
- MCP provides a standardized way to connect AI models to different data sources and tools
- Client/server architecture

elastic

# Core architecture

- MCP follows client-server architecture
- Transport layer: Stdio, Streamable HTTP
- All transport use JSON-RPC 2.0 to exchange messages



elastic

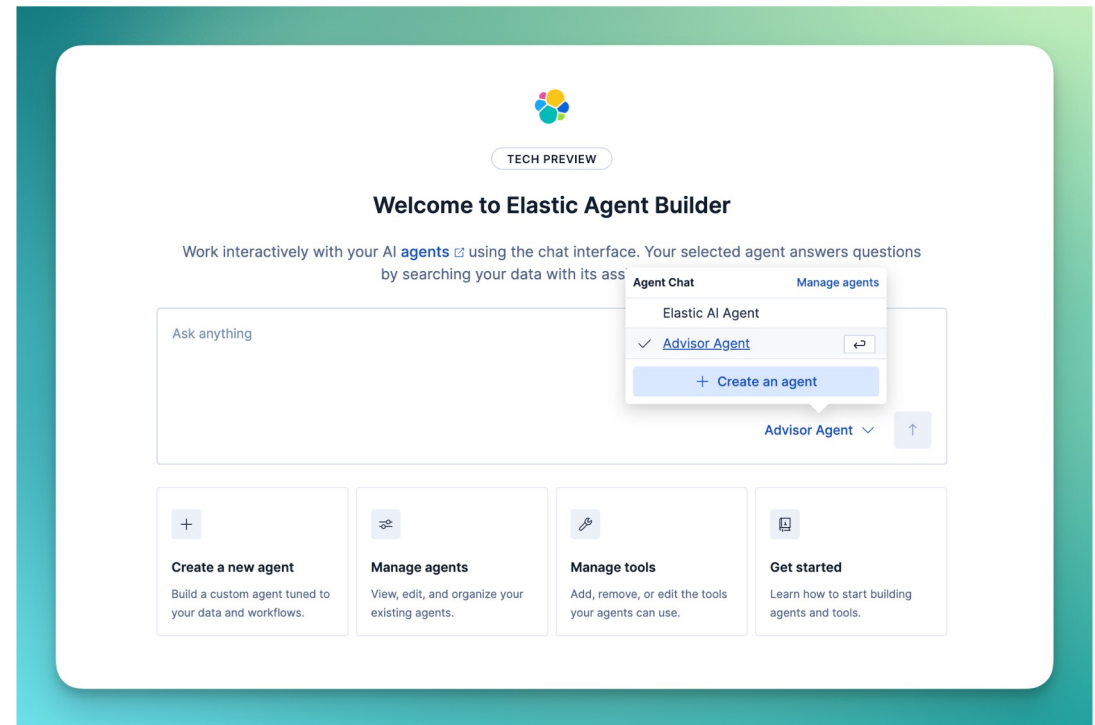# Example: multiple MCP servers

# Meet Elastic Agent Builder

Develop custom AI agents and chat experiences in minutes with Elasticsearch leading relevance and new agent-building capabilities

➔ **Natively chat with any data in Elastic** using a built-in agent

➔ **Build custom AI agents** that achieve higher accuracy, relevance, and efficiency based on the power of hybrid search

➔ **Create powerful tools:** Give your agent new powers. Build custom tools with the full power of ES|QL for fine-grained control over data, relevance, and security.

➔ **Expose** your data by hosting **MCP and A2A** compatible interfaces to your AI ecosystem.

# Conversational chat can help you …

Improve technical support resolution times with Agents to aid investigations

Improve sales reporting and analysis with Agents to summarize and deliver insights

Understand and create policies (e.g. HR, IT, Legal) based on existing and usage (e.g. past searches/conversations)

Investigate purchase and product trends to make better recommendations

Create intelligent product or content recommendations based on catalog and past searches, interactions

Any other internal knowledge worker use case...

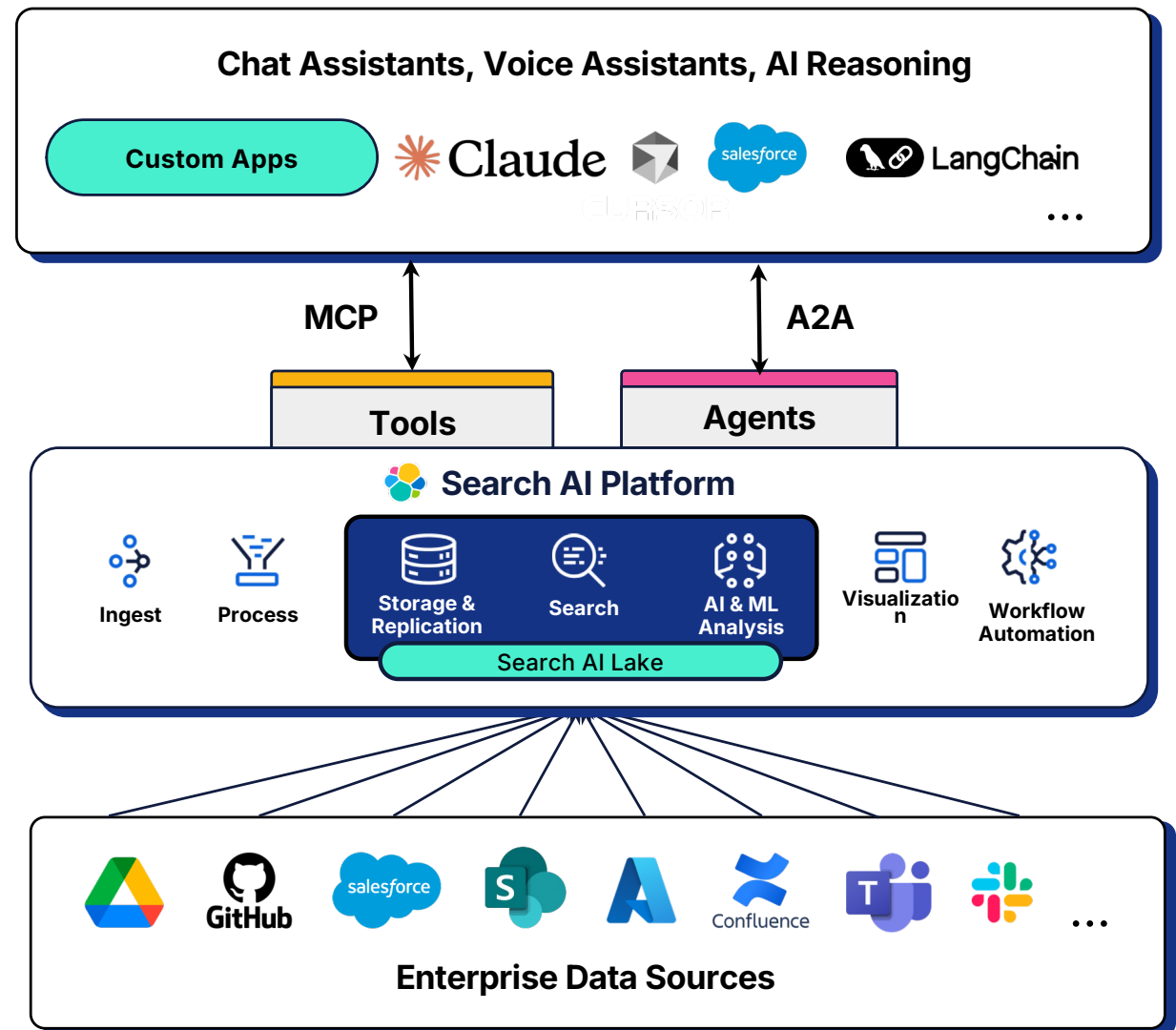Ensure compliance across unstructured/structured data sources (e.g. KYC)

# One unified platform to build custom AI agents fast

**AI Native Experiences**

**Chat Assistants, Voice Assistants, AI Reasoning**

Custom Apps · Claude · CURSOR · salesforce · LangChain · ...

**Powered by...** Tools & Agents

MCP ↕ | A2A ↕

**Tools** | **Agents**

**Enabled by...** The Platform

**Search AI Platform**

Ingest · Process · Storage & Replication · Search · AI & ML Analysis · Visualization · Workflow Automation

Search AI Lake

**Built on...** Enterprise Data

Google Drive · GitHub · salesforce · SharePoint · Azure · Confluence · Teams · Slack · ...

**Enterprise Data Sources**

elastic

# Quickly build custom agents that utilize all your data powered by Elasticsearch context

**Build on the best Vector Database technology**

---

*Speed, Scale, Efficiency, Relevance*

**Data Relevance built to match the intent of users**

---

*Hybrid Search, Machine Learning, Small Language Models*

**Unify Agent Data layer with Enterprise quality**

---

*Data Security, Multi-Cloud*

elastic
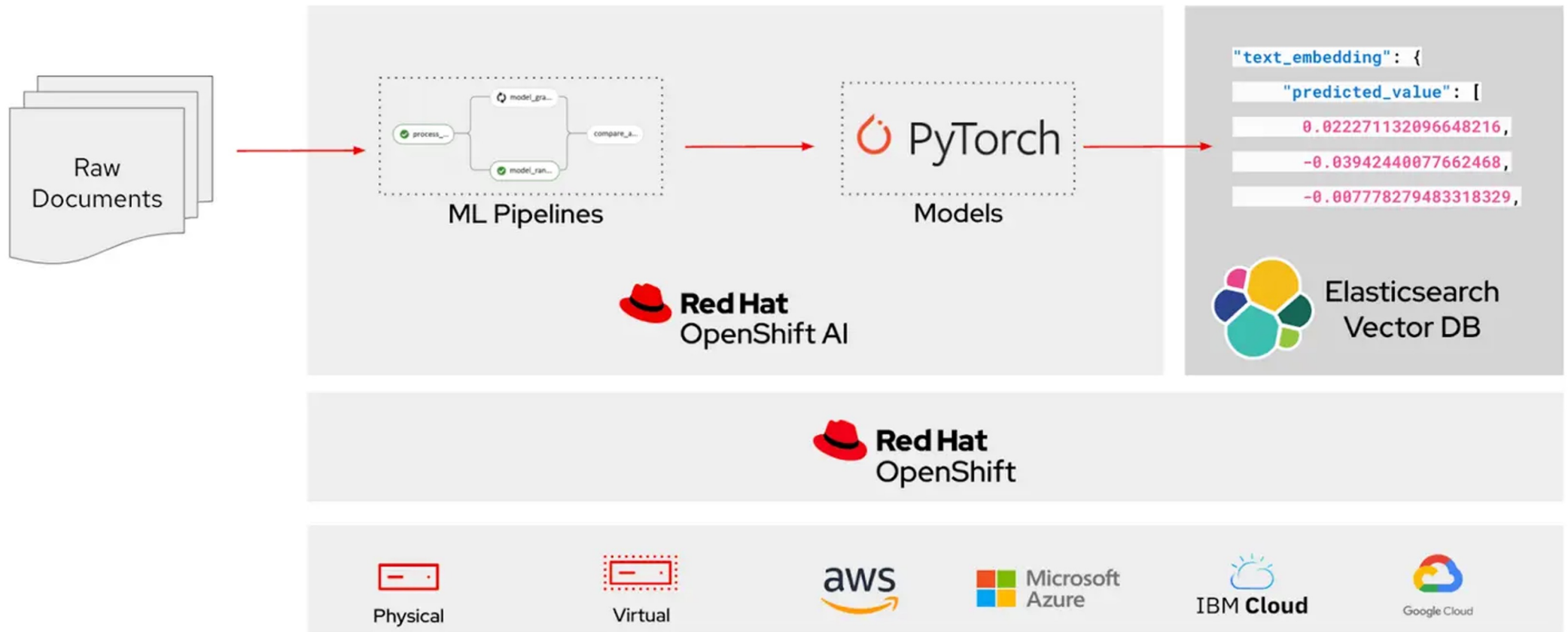
# Elasticsearch (ECK) Operator ✓ Certified
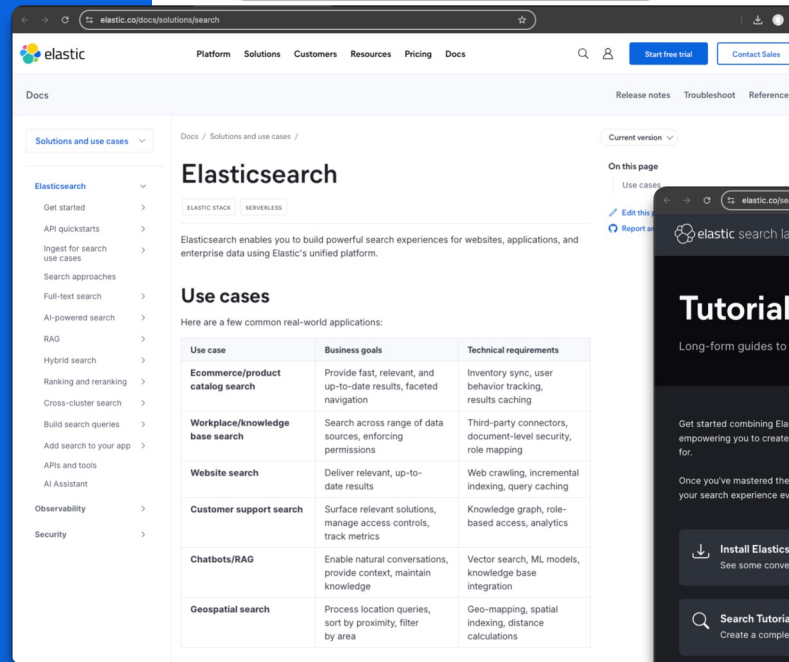
**ECK into Red Hat Catalog**

## Elastic & Red Hat

- Easily manage the Elastic Stack on Openshift
- Secure by Default
- Advanced Topology Support
- Snapshot scheduling and keystore support
- Cross-cluster search and cross-cluster replication
- Store local, search Global

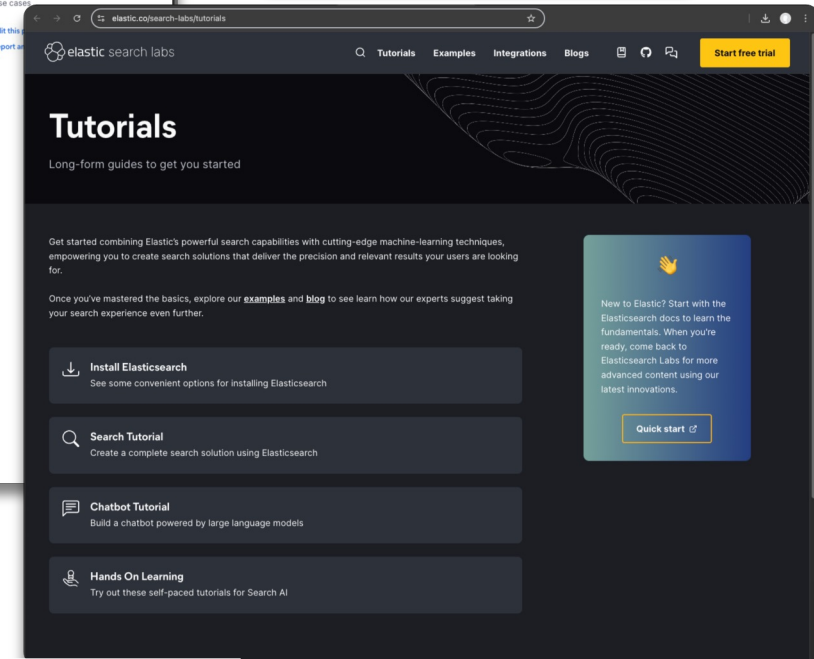# Red Hat Openshift AI & Elasticsearch Vector DB

# Want to
## learn more?

## Elastic Agent Builder - Tools, Agents, and MCP

In this hands-on course, learn how to create custom AI tools from your own business logic, build specialized agents to accurately chat with your data, and integrate them into external applications via the built-in MCP server.

elastic

**Red Hat Summit**

## Connect

linkedin.com/company/red-hat

facebook.com/redhatinc

youtube.com/user/RedHatVideos

twitter.com/RedHat