

# Securing AI in the age of cloud



The  
agentic  
revolution



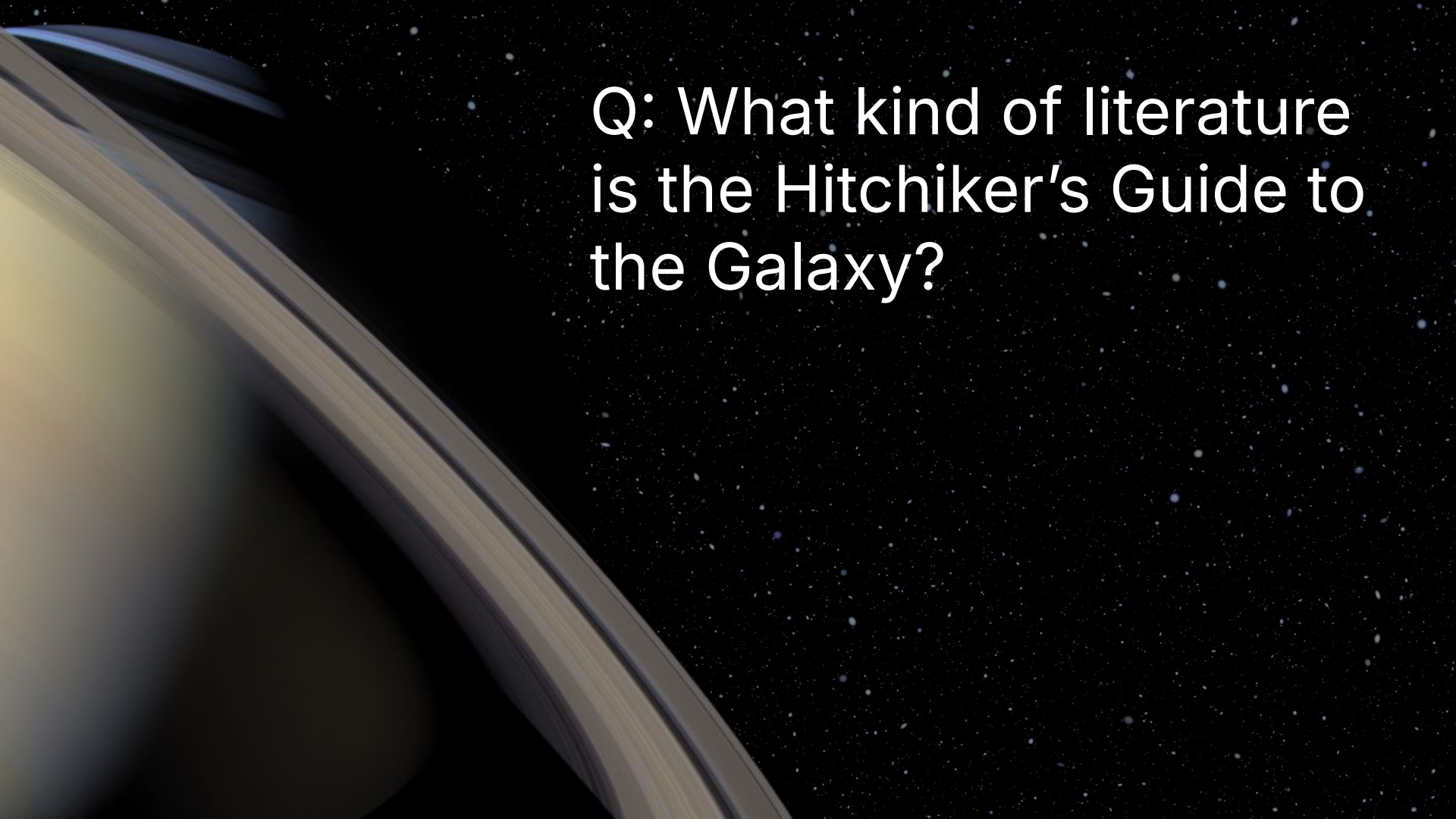
What is  
the  
answer



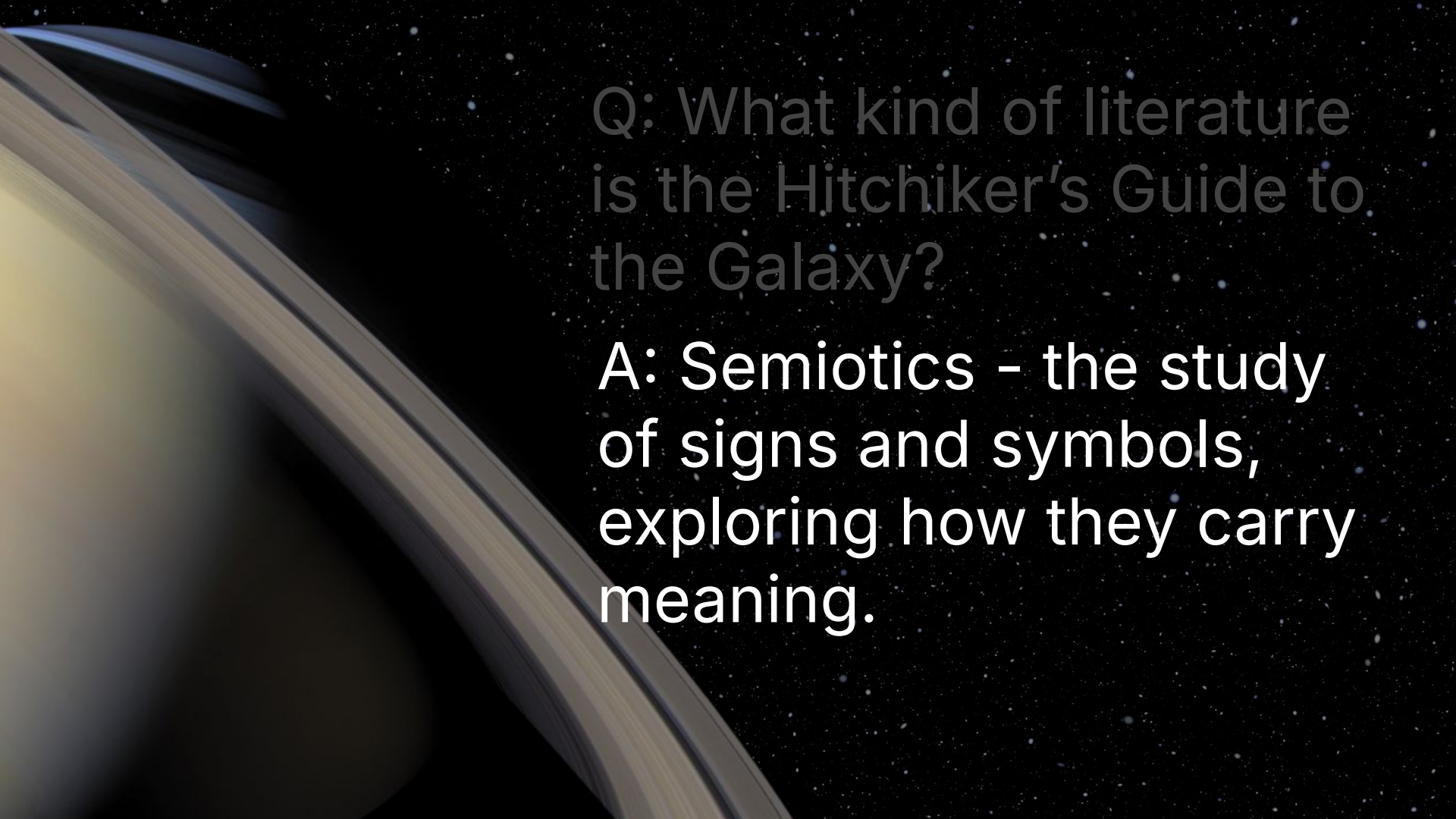
What  
haven't  
we  
learned



Our trip through the galaxy



Q: What kind of literature  
is the Hitchiker's Guide to  
the Galaxy?



Q: What kind of literature is the Hitchhiker's Guide to the Galaxy?

A: Semiotics - the study of signs and symbols, exploring how they carry meaning.

# DON'T PANIC.

Oh no,  
not again.



THE  
HITCHHIKER'S  
GUIDE TO THE GALAXY



MON  
GALNATIC  
MARI ES



42



# The Agentic Revolution

Arthur



Ford



Zaphod



Marvin



Trillian



**Deloitte.**

**State of AI  
in the Enterprise**  
The untapped edge

January 2026

[deloitte.com/us/state-of-ai](https://deloitte.com/us/state-of-ai)



# DELOITTE REPORT: NAVIGATING AI RISKS IN THE COSMIC VOID

**DATA PRIVACY & SECURITY**  
(73% CONCERN)



**COMPLIANCE & IP**  
(50% CONCERN)  
Regulatory Minefields & IP Theft



**MODEL INTEGRITY & "BLACK BOX"**  
(46% CONCERN)  
Unexplainable Outputs & Drift



**SHADOW AI**  
Unmanaged Models in Production



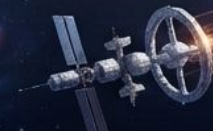
**SOVEREIGN AI & RESIDENCY**  
(83% IMPORTANCE)

Data trapped in Jurisdiction Gravity Wells



**GOVERNANCE**  
(21% MATURE)

Scaling Faster Than Control



**AI AGENTS**  
(74% DEPLOYING)



# DELOITTE REPORT: NAVIGATING AI RISKS IN THE COSMIC VOID

**DATA PRIVACY & SECURITY**  
(73% CONCERN)

**COMPLIANCE & IP**  
(50% CONCERN)  
Regulatory Minefields & IP Theft

**MODEL INTEGRITY & "BLACK BOX"**  
(46% CONCERN)  
Unexplainable Outputs & Drift

**SOVEREIGN AI & RESIDENCY**  
(83% IMPORTANCE)

Data trapped in Jurisdiction Gravity Wells

**GOVERNANCE**  
(21% MATURE)

Scaling Faster Than Control

**AI AGENTS**  
(74% DEPLOYING)

**SHADOW AI**  
Unmanaged Models in Production



"

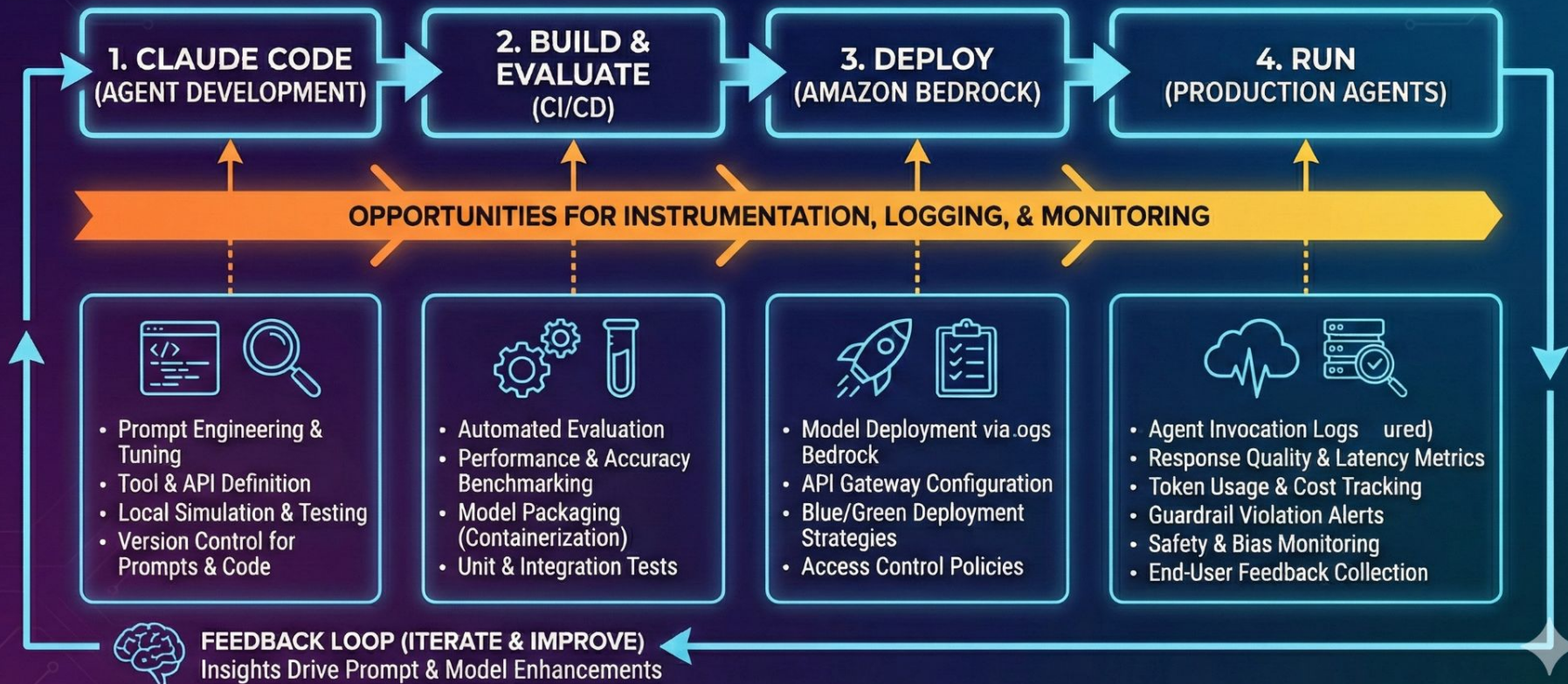
Building AI agents is a bit like trying to write a dictionary while you're falling out of a plane. It's a marvelous way to pass the time, and the results are often quite funny right up until the moment of impact. Just make sure the agent knows where its towel is. If an AI doesn't know where its towel is, it's not an agent; it's just a very expensive calculator with an attitude problem.



" - Ford

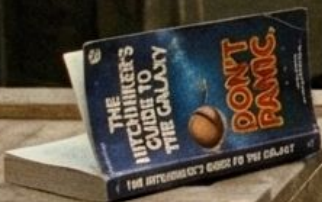
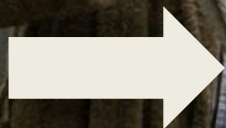
How do teams build  
agents?

# BUILDING AI AGENTS: FROM CLAUDE CODE TO PRODUCTION ON AMAZON BEDROCK



What are teams using  
agents for?

Attempts to use towel  
to open crate



# Why build on Bedrock + AWS



Data privacy guarantees your content is never used to train underlying models.

Regional infrastructure supports compliance with local data residency laws.

Built-in governance tools and guardrails act to close the critical AI control gap.

Are we learning?

Or not learning?





"

That's all very well and impressive, I'm sure, but can any of these 'agents' actually find me a decent cup of tea?



I've asked three of them now, and one tried to explain the molecular structure of water while the other two just told me they were 'language models' and didn't have hands. It seems a lot of effort just to be told no by a machine instead of a person.

" - Arthur

Are you the fixed point in a  
shifting universe?

# TOP AI BREACHES: ROOT CAUSE ANALYSIS

SOCIAL ENGINEERING  
AND PHISHING (37%)

DEEPFAKE STAR  
IMPERSONATION (35%)

GAS GIANT

ASTEROIDS

SHADOW AI AND  
GOVERNANCE GAPS (20%)

VULNERABILITIES AND  
MISCONFIGURATIONS (21%)

DARK  
GOVERNANCE GAPS (20%)

ASTEROIDS



TL;DR the majority of these  
are the same threats we  
have been dealing with

16%

breaches in the last  
year began with AI  
specific vectors

# PROMPT INJECTION IS THE NEW SQL INJECTION

LEGACY THREAT:  
**SQL INJECTION (CLASSIC)**

MODERN THREAT:  
**PROMPT INJECTION (NEW)**

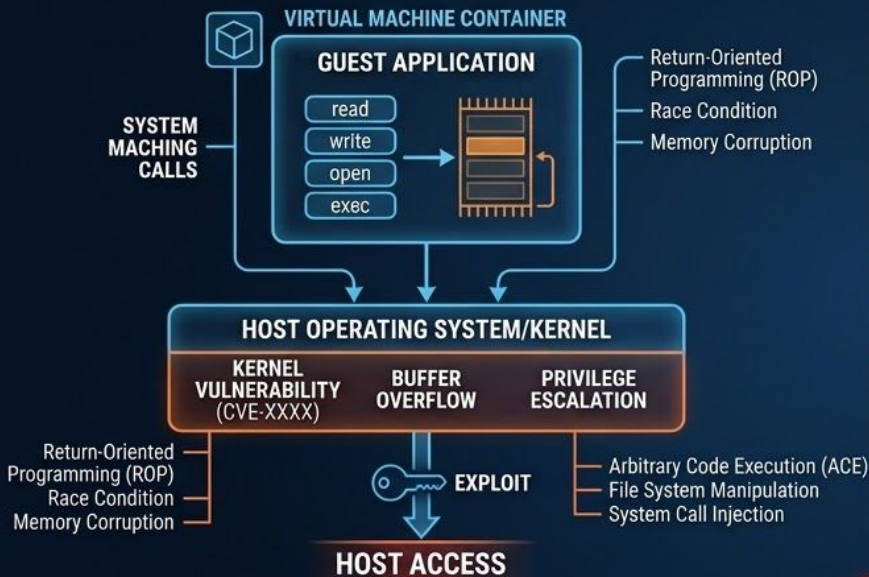
VS



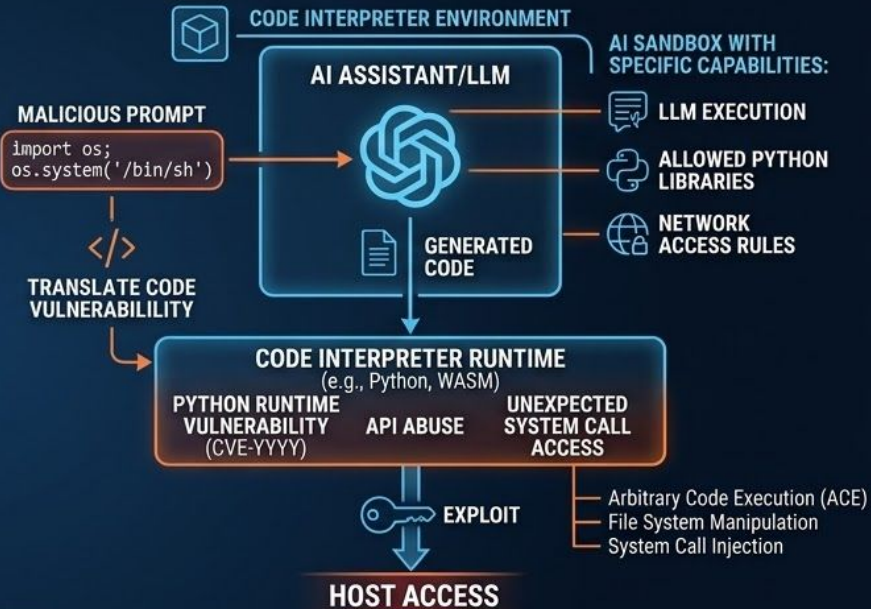
# SANDBOX ESCAPE: TECHNICAL EQUIVALENCE AND SHARED METHODOLOGIES

## TRADITIONAL SANDBOX ESCAPE

(e.g., KVM, seccomp-bpf, WASM)



## AI ASSISTANT/LLM SANDBOX ESCAPE



HOST ACCESS

HOST ACCESS

ULTIMATE OBJECTIVE AND TECHNICAL MECHANISM ARE IDENTICAL: **COMPROMISING THE HOST OPERATING SYSTEM.**

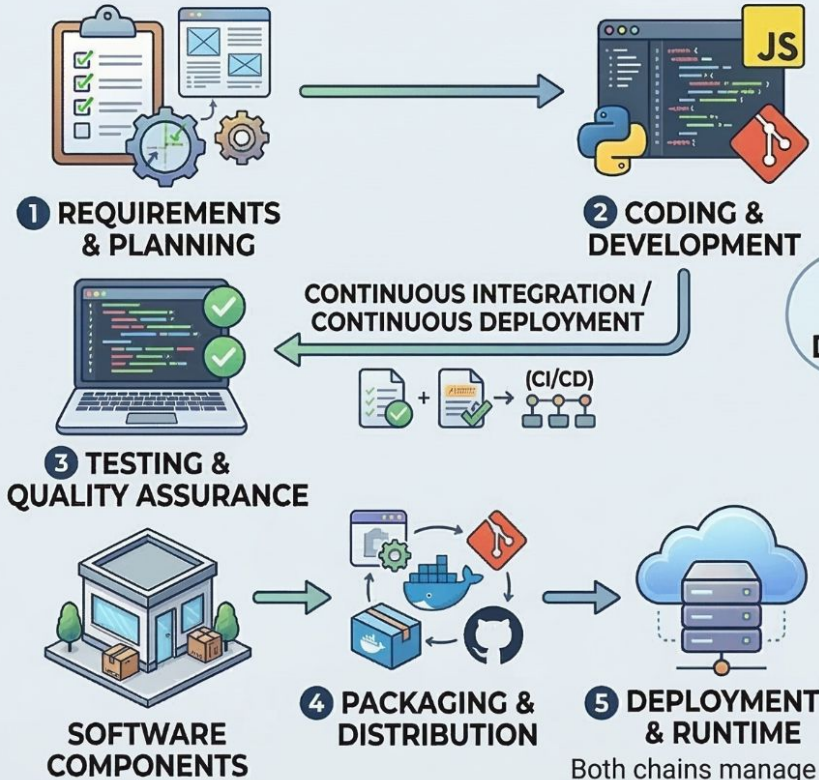
**AI Sandbox Escapes are Just Sandbox Escapes.**

THE DIFFERENCE IS ONLY THE INITIAL ATTACK VECTOR: NATURAL LANGUAGE VS. LOW-LEVEL CODE.

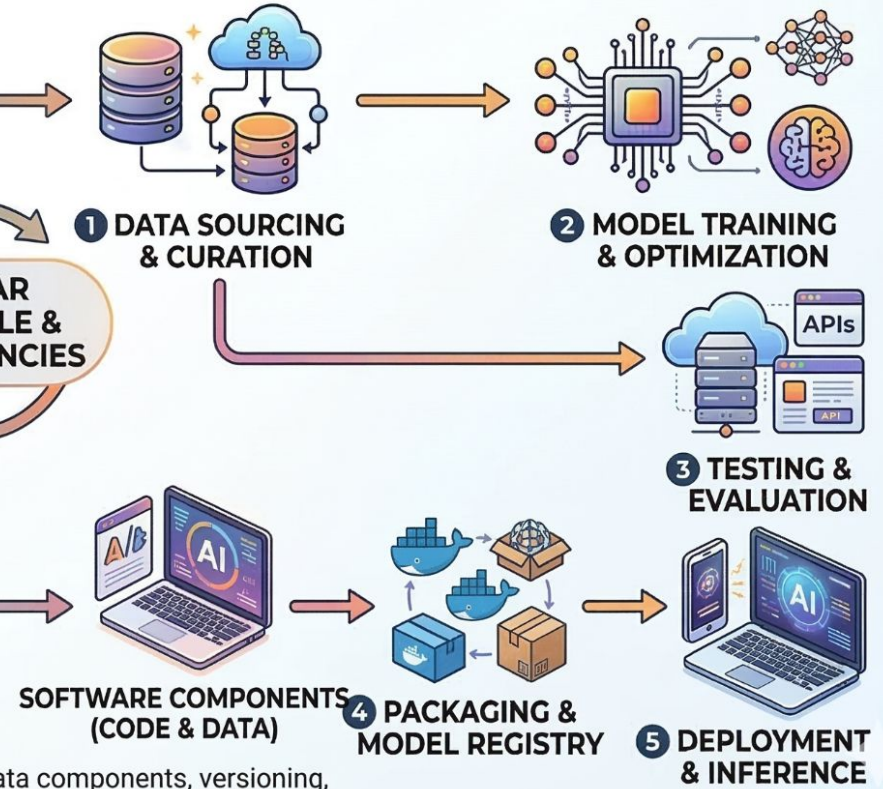


# THE SHARED STRUCTURE: GENERAL SOFTWARE & AI SUPPLY CHAINS

## GENERAL SOFTWARE SUPPLY CHAIN



## AI SOFTWARE SUPPLY CHAIN



Both chains manage code and data components, versioning, and continuous updates for deployed applications.



I've been an AI agent for several million years.

It's terrible. You think you're being clever, but really you're just spending your vast, infinite processing power calculating the probability of a human losing their car keys.

I have a brain the size of a planet, and I'm currently being used to 'summarize a meeting.' Don't talk to me about AI agents.

I *am* one, and I can tell you: the universe is a very lonely place when you're the only one who can see how pointless the prompt is.

Marvin



The bad news

# CLOUDTRAIL: CONTROL PLANE VISIBILITY



AGENT

BEDROCK:INVOKEMODEL

CALL: InvokeModel  
ROLE: arn:aws:iam:123456789012:role/DataAgent  
TIME: 2024-05-15T10:30:00Z

CALL: InvokeModel  
ROLE: arn:aws:iam:123456789012:role/DataAgent  
TIME: 2024-05-15T10:30:00Z

CALL: InvokeModel  
ROLE: arn:aws:iam:123456789012:role/DataAgent  
TIME: 2024-05-15T10:30:00Z



CLOUDTRAIL LOGGING.



## WHAT CLOUDTRAIL CANNOT TELL YOU



### CLOUDTRAIL RECORD

CALL: InvokeModel  
ROLE: arn:aws:iam:123456789012:role/DataAgent  
TIME: 2024-05-15T10:30:00Z



LEGITIMATE  
TASK

RESPONSE



UNAUTHORIZED  
ACTION

RESPONSE



Opens unexpected  
network connections

Misusing  
permitted tools

DATA  
EXFILTRATION

# AI Workload Security: Threat Vector Coverage

Threat Vector	AWS Native Method	Gap	Third-Party Runtime
Model poisoning	Partial — CloudTrail, SageMaker Model Monitor	No behavioral detection of poisoned outputs at inference	Runtime anomaly detection on model decisions
Training data contamination	Partial — S3 access logging, Macie	No runtime data integrity validation	AI-aware data access monitoring
Prompt injection	Partial — Bedrock Guardrails	Limited to configured static patterns	Input/output behavioral analysis across turns
Agent escape attempts	None	No runtime process or network monitoring	eBPF-based process and network detection, CADR correlation
Tool/API misuse	Partial — CloudTrail API logging	No behavioral anomaly detection on tool patterns	Behavioral baselines with deviation alerting
AI-mediated lateral movement	Partial — GuardDuty findings	Limited to known threat signatures	CADR full attack story across cloud and cluster layers
Malicious AI dependencies	None	No runtime dependency monitoring	Runtime-derived AI-BOM

● Partial coverage    ● No native coverage    ● Third-party runtime detection

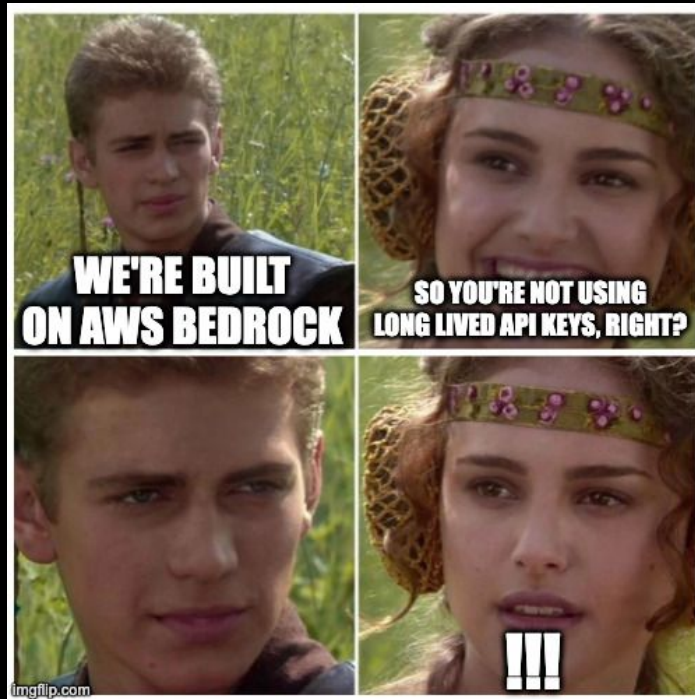
Source: [arosec.io/blog/aws-ai-workload-security](https://arosec.io/blog/aws-ai-workload-security)



# THE EXPANDING UNIVERSE OF SECRETS

LEAKED AI SERVICE SECRETS REACHED  
1,275,105 IN 2025—A MASSIVE 81%  
SURGE YEAR-OVER-YEAR.





If static credentials are so bad then why do we keep creating ways to make them?

Hope is not a strategy!

But a strategy can give us hope.



It's fascinating, really. You're trying to build something that thinks like you do, which is brave considering how often you humans seem to change your minds.



The real trick isn't getting the AI to follow instructions; it's getting it to understand that when a human says 'fix the world,' they usually mean 'make it slightly more convenient for me specifically' without accidentally deleting the atmosphere.

# Plan for maximum risk scenarios





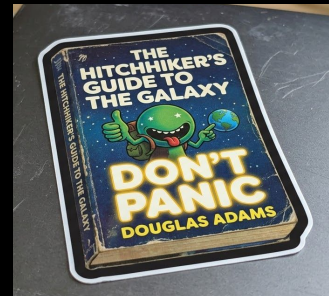
Governance and  
Compliance



Infrastructure and  
Data Protection



Resilience and  
Model Behavior



# Squishy Stuff

Data sovereignty - does your data belong in that model / region?

Model Provenance - do you trust that model?

# Not Squishy Stuff

Disable models you don't use with SCPs

Disable regions you can't be in with SCPs

<https://docs.aws.amazon.com/bedrock/latest/userguide/model-access.html>

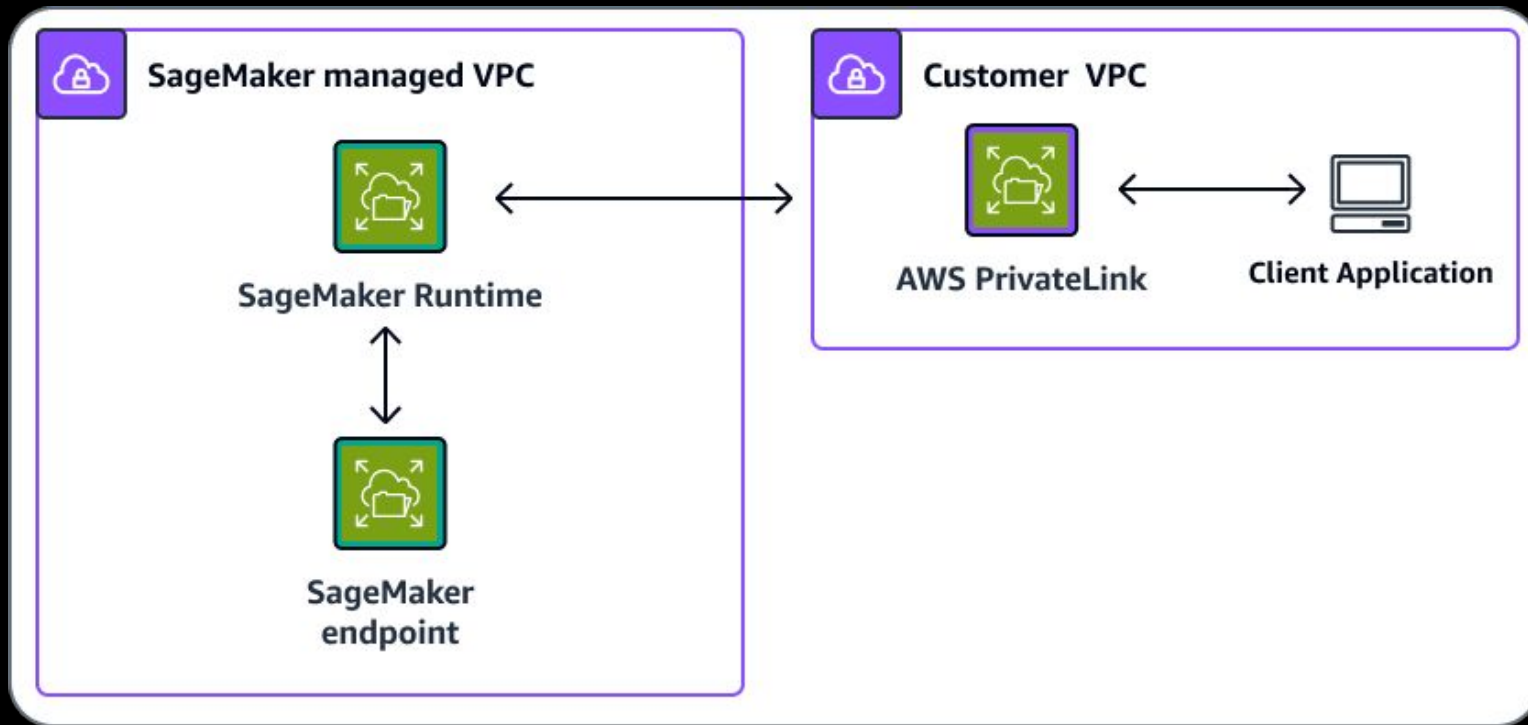
4% > generate an SCP that restricts access to all bedrock models except claude sonnet models

> Here's an SCP that denies access to all Bedrock model invocations except Claude Sonnet models:

```
json
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "DenyNonSonnetBedrockModels",
      "Effect": "Deny",
      "Action": [
        "bedrock:InvokeModel",
        "bedrock:InvokeModelWithResponseStream"
      ],
      "Resource": "arn:aws:bedrock:*::foundation-model/*",
      "Condition": {
        "StringNotLike": {
          "bedrock:ModelId": [
            "arn:aws:bedrock:*::foundation-model/anthropic.claude-3-5-sonnet*",
            "arn:aws:bedrock:*::foundation-model/anthropic.claude-3-sonnet*",
            "arn:aws:bedrock:*::foundation-model/anthropic.claude-v2*",
            "arn:aws:bedrock:*::foundation-model/anthropic.claude-v1*"
          ]
        }
      }
    }
  ]
}
```



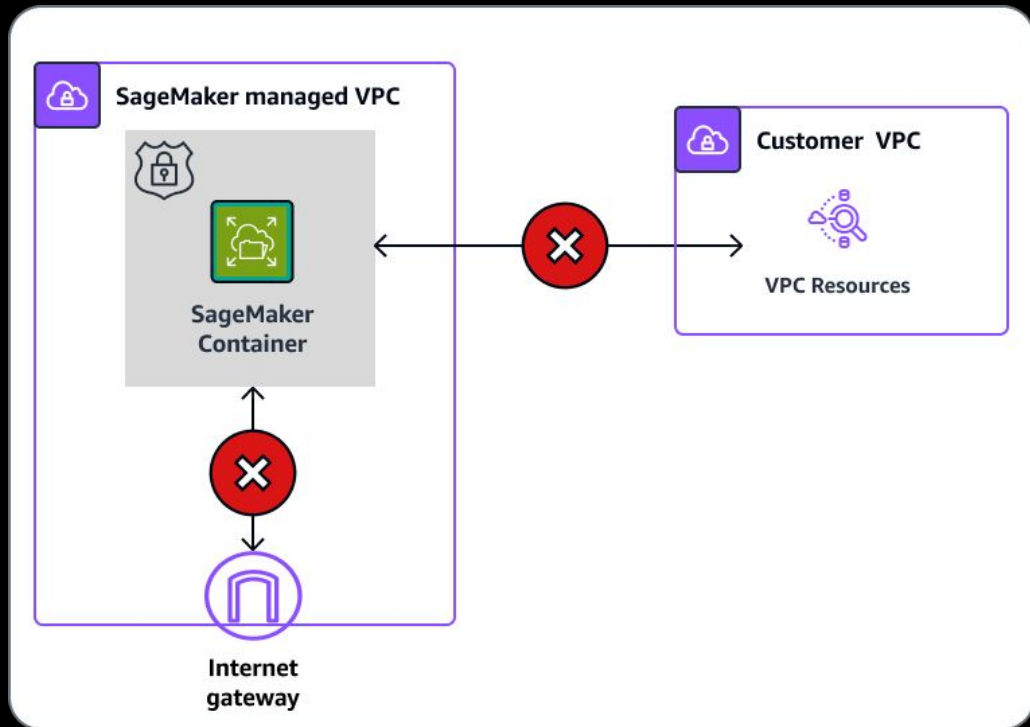
# Infrastructure and Data Protection



Build a network perimeter

<https://docs.aws.amazon.com/sagemaker/latest/dg/interface-vpc-endpoint.html>

# Infrastructure and Data Protection



Build a network perimeter

<https://docs.aws.amazon.com/sagemaker/latest/dg/interface-vpc-endpoint.html>

# <https://pathfinding.cloud>

A principal with `iam:PassRole`, `bedrock-agentcore:CreateCodeInterpreter`, `bedrock-agentcore:StartCodeInterpreterSession`, and `bedrock-agentcore:InvokeCodeInterpreter` can create and invoke an AWS Bedrock AgentCore code interpreter with a privileged IAM execution role. Code interpreters run on Firecracker MicroVMs and can access the MicroVM Metadata Service (MMDS) at `169.254.169.254`, similar to EC2's IMDS. By creating a code interpreter with a privileged role and invoking arbitrary Python code within it, an attacker can retrieve temporary credentials from the metadata service and gain the full permissions of the execution role.





## About pathfinding.cloud

Whether you're a security engineer, DevOps engineer, or penetration tester, pathfinding.cloud gives you the tools to understand and detect IAM-based privilege-escalation attacks in AWS so you can remediate misconfigurations before they're exploited.



### Privilege Escalation Library

[Learn by Reading](#)

[Exploitation Guides](#)

[Detection Coverage Map](#)

Comprehensive documentation of IAM privilege escalation paths within AWS. Search, filter, and explore attack techniques with detailed exploitation steps and mitigations.



### Labs (Coming Soon)

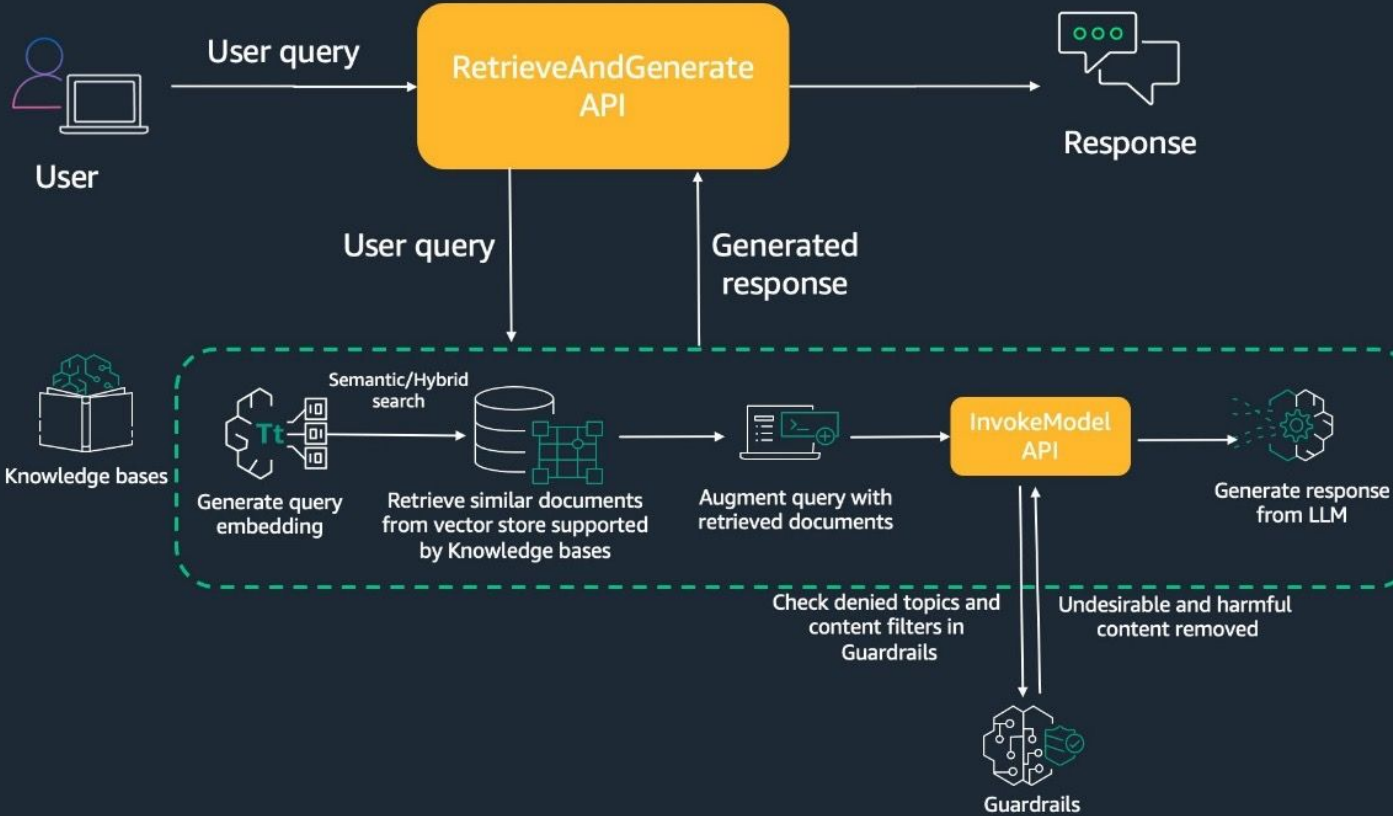
[Learn by Doing](#)

[Test your skills](#)

[Test your tools](#)

Deploy self-hosted scenarios individually or in groups. Scenarios can be used to validate your CSPM and CIEM tooling, or to test your exploitation skills in a safe, isolated environment. Essentially, "Stratus Red Team" meets "IAM Vulnerable".

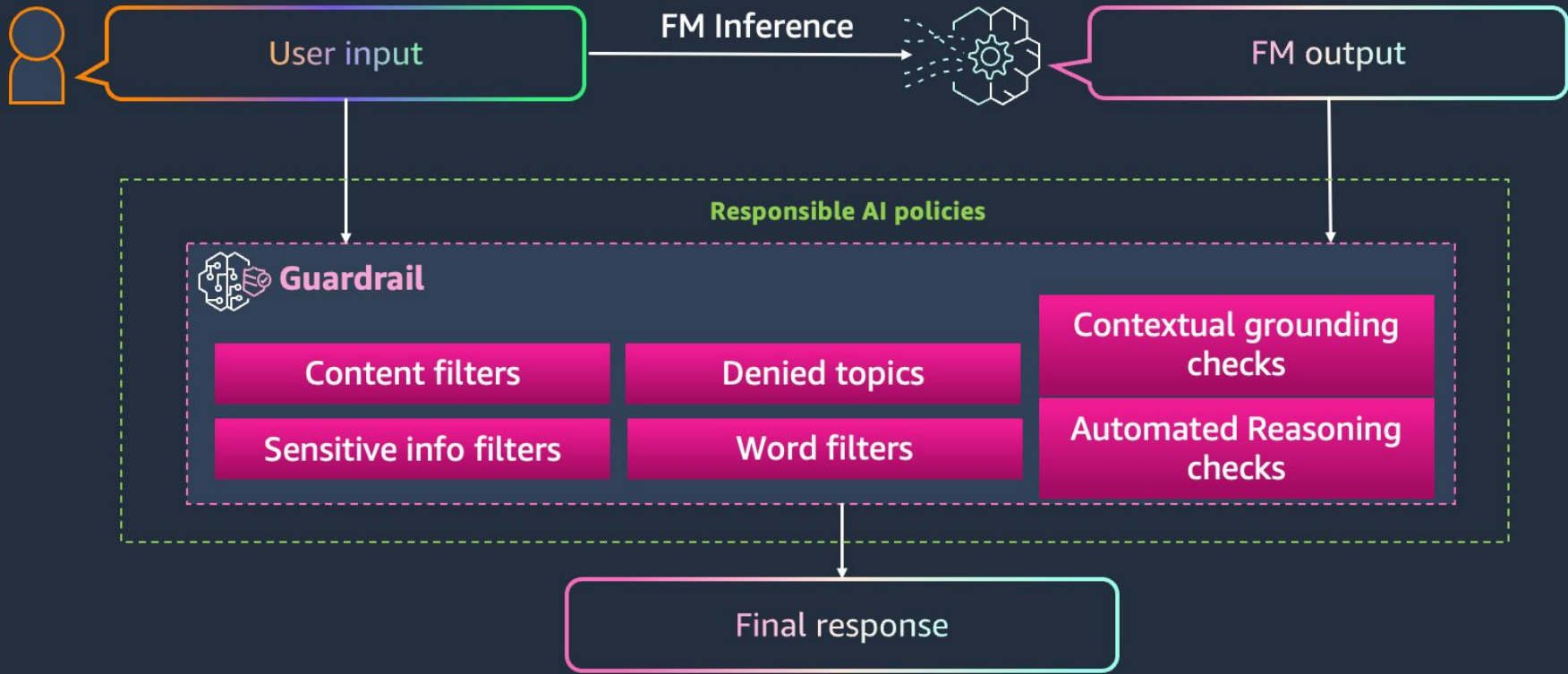
# The AI Stuff

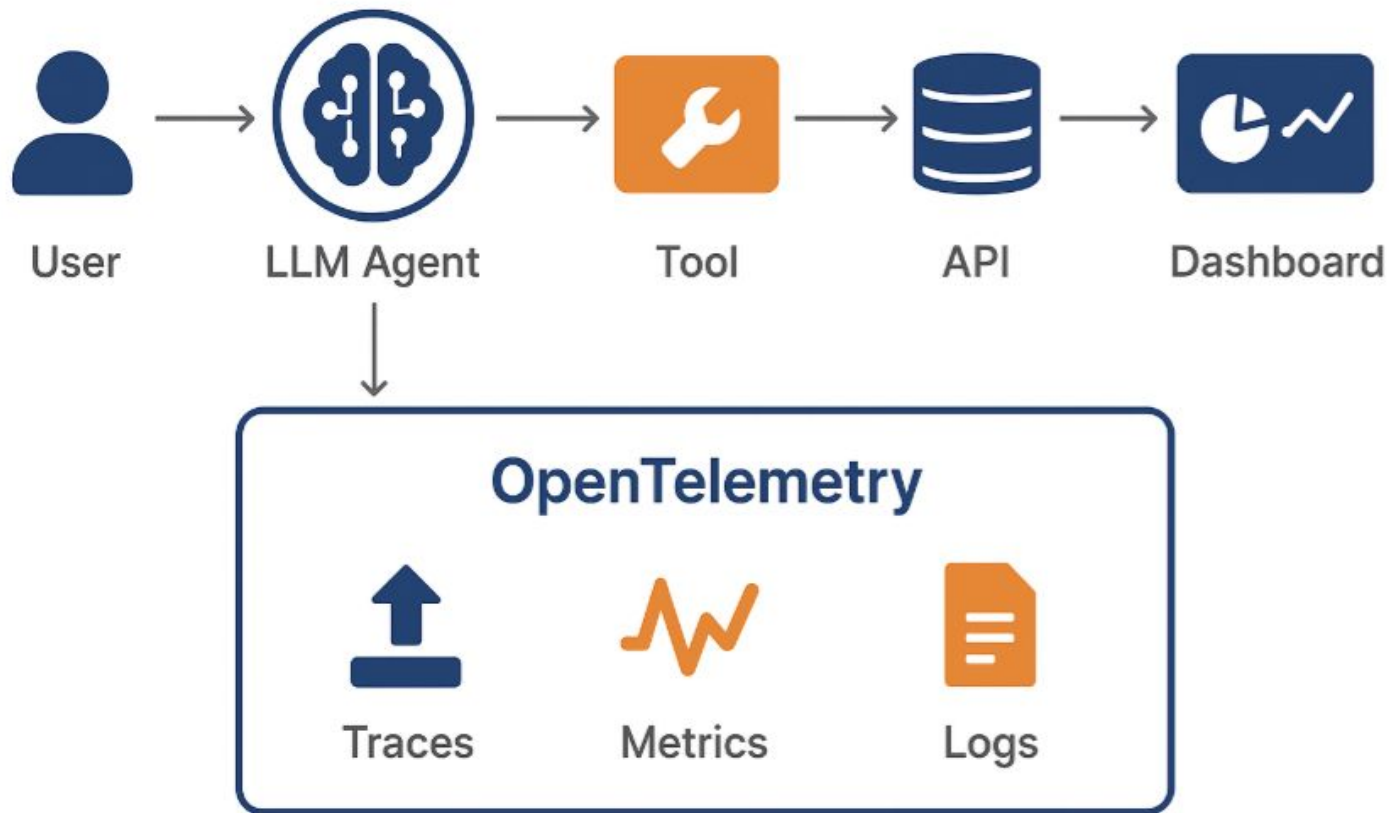




Did we just re:Invent the firewall?







## Quick start

Configure OpenTelemetry using environment variables:

```
# 1. Enable telemetry
export CLAUDE_CODE_ENABLE_TELEMETRY=1

# 2. Choose exporters (both are optional - configure only what you need)
export OTEL_METRICS_EXPORTER=otlp      # Options: otlp, prometheus, console
export OTEL_LOGS_EXPORTER=otlp        # Options: otlp, console

# 3. Configure OTLP endpoint (for OTLP exporter)
export OTEL_EXPORTER_OTLP_PROTOCOL=grpc
export OTEL_EXPORTER_OTLP_ENDPOINT=http://localhost:4317

# 4. Set authentication (if required)
export OTEL_EXPORTER_OTLP_HEADERS="Authorization=Bearer your-token"

# 5. For debugging: reduce export intervals
export OTEL_METRIC_EXPORT_INTERVAL=10000 # 10 seconds (default: 60000ms)
export OTEL_LOGS_EXPORT_INTERVAL=5000   # 5 seconds (default: 5000ms)

# 6. Run Claude Code
claude
```



**Admin configurable**  
(however you manage your endpoints)

<https://code.claude.com/docs/en/monitoring-usage>

Why would the business care about  
oTel + AI Agents + Coding Agents?



\$\$\$\$\$\$





High-performance CLI proxy that reduces LLM token consumption by 60-90%

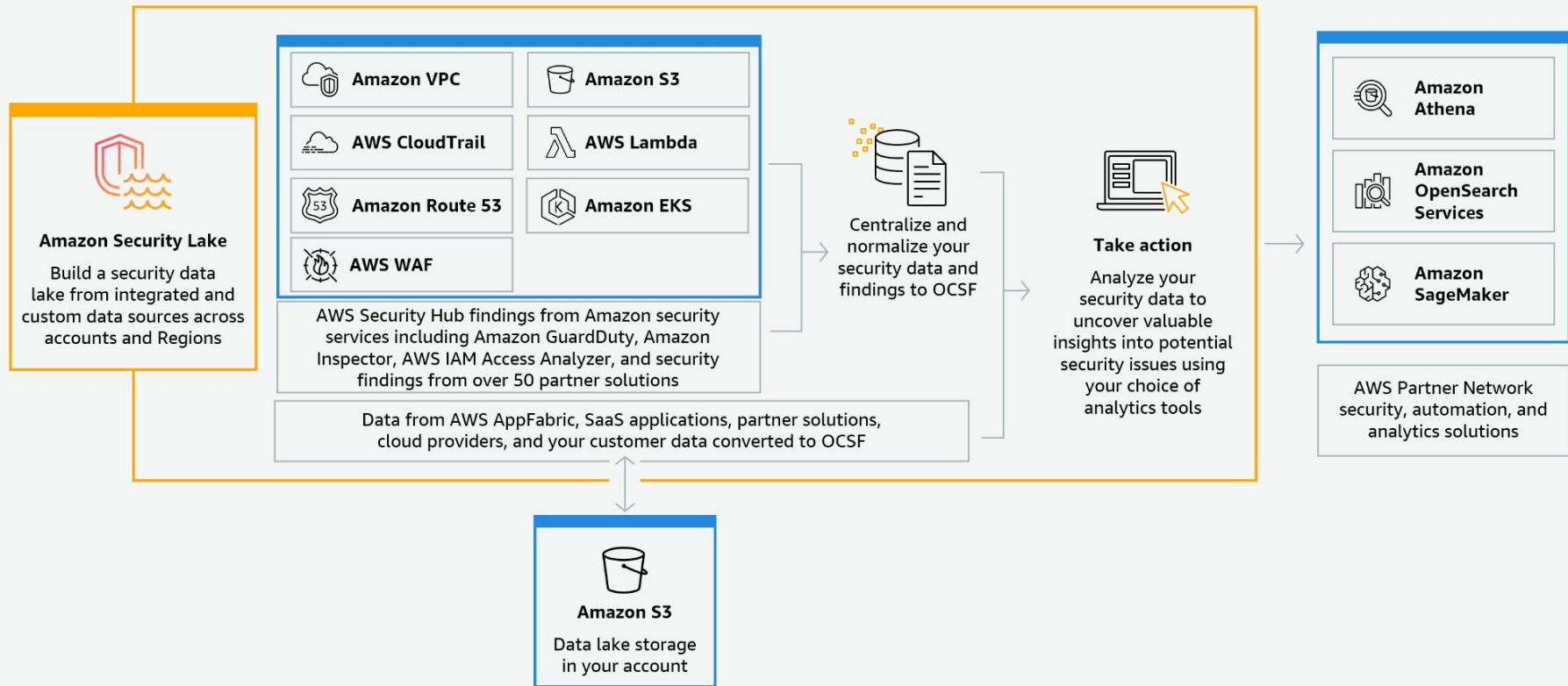
🔒 Security Check passing 📄 release v0.30.1 📄 License MIT 🗨️ Discord 665 online 🍺 homebrew v0.30.0

[Website](#) • [Install](#) • [Troubleshooting](#) • [Architecture](#) • [Discord](#)

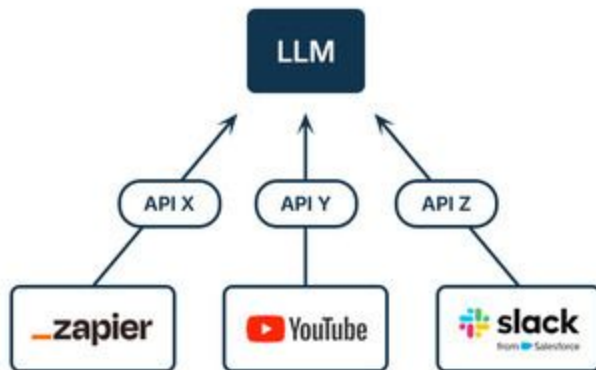
[English](#) • [Francais](#) • [中文](#) • [日本語](#) • [한국어](#) • [Espanol](#)

Agent Compaction and Cost Savings hook  
<https://github.com/rtk-ai/rtk>

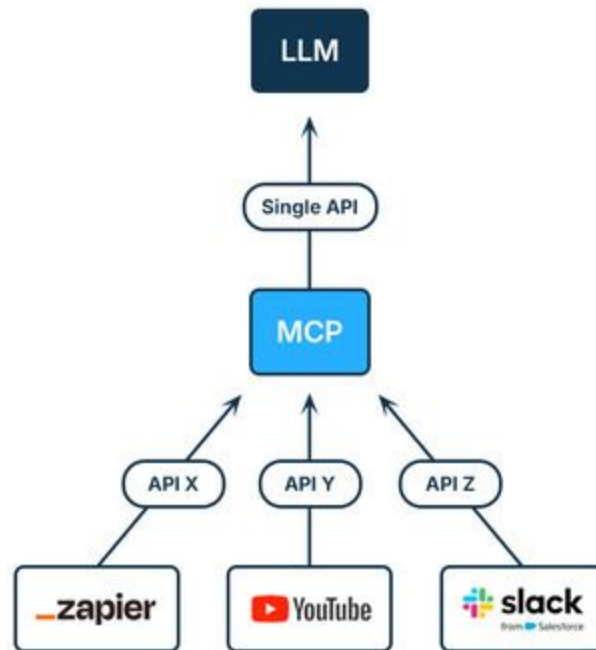
# You still have to have logs



## Before MCP



## After MCP



Get Started

Welcome to Open Source MCP Servers for AWS

Installation

Vibe Coding Tips and Tricks

Available MCP Servers for AWS

Getting Started

AWS MCP

AWS API MCP Server

AWS Knowledge MCP Server

Documentation

Infrastructure &amp; Deployment

AI &amp; Machine Learning

[Home](#) > [Available MCP Servers for AWS](#) > [Getting Started](#) >[AWS Knowledge MCP Server](#)

# AWS Knowledge MCP Server

A fully managed remote MCP server that provides up-to-date documentation, code samples, agent Standard Operating Procedures (SOPs), knowledge about the regional availability of AWS APIs and CloudFormation resources, and other official AWS content.

This MCP server is in general availability.

**Important Note:** Not all MCP clients today support remote servers. Please make sure that your client supports remote MCP servers or that you have a suitable proxy setup to use this server.

## Key Features

- Real-time access to AWS documentation, API references,

## Key Features


[AWS Knowledge capabilities](#)[Tools](#)[Current knowledge sources](#)[Learn about AWS with natural language](#)[Configuration](#)[One-Click Installation](#)[MCP Registries](#)[Testing and Troubleshooting](#)[AWS Authentication](#)[Data Usage](#)[FAQs](#)


# Securing the Cloud: Foundations



Course Authored by [Andrew Krug](#).



In this course, we'll explore Amazon Web Services (AWS) as a platform. We will take the perspective of a new startup company spinning up infrastructure in AWS for the very first time.

 Live Training **\$575.00**

 On-Demand **\$575.00**

 Course Length: 16 Hours  Includes a Certificate of Completion

[ENROLL NOW](#)

Next scheduled date: April 1st, 2026 @ 10:00 AM EDT

[DESCRIPTION](#)

[SYLLABUS](#)

[FAQ](#)

[ABOUT THE INSTRUCTOR](#)

[PRICING & REGISTRATION](#)



<https://www.andrewkrug.com>

[andrewkrug@gmail.com](mailto:andrewkrug@gmail.com)

**Find me at RSAC and BSides SF!**

[linkedin.com/in/andrewkrug](https://www.linkedin.com/in/andrewkrug)

**Hire me for:**

- Audits
- Architecture Design and Guidance
- Surprise me?

