# AI Security & Insights into OWASP Top 10 LLM

February 11th, 2025

# whoami

## Rico Komenda

**Who?**
- Husband and father

- Senior Security Consultant @ adesso SE
- International trainer and speaker
- Consulting
  - AppSec, CloudSec, OffSec, AISec

<u>Mission statement:</u>
- Securing the digital world, one byte at a time

# Agenda

- **AI Security**
  - **with vs. for**
- **State of the LLM/GenAI project updates**
  - **AI Red Teaming Initiative**
  - **Agentic Security Initiative (ASI)**
- **Open Discussion**

# AI Security

# Hype Cycle

# Security for AI vs. Security with AI

## for AI:

Securing AI systems and applications

## with AI:

enhance systems and applications with AI capabilities

# Security with AI – Use Cases

Identity and access management
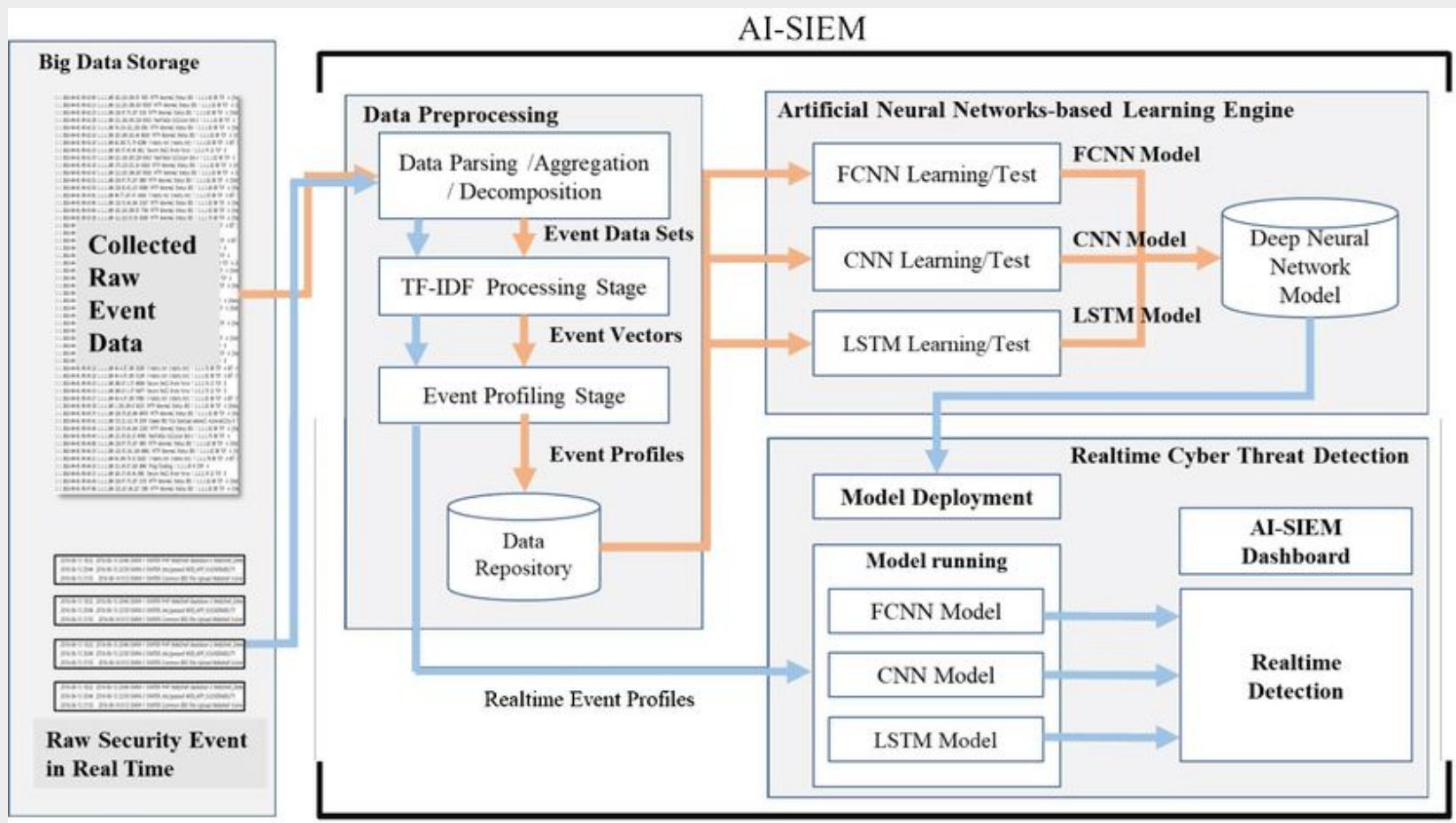
Endpoint security and management

Cyber Threat detection (XDR; SIEM)
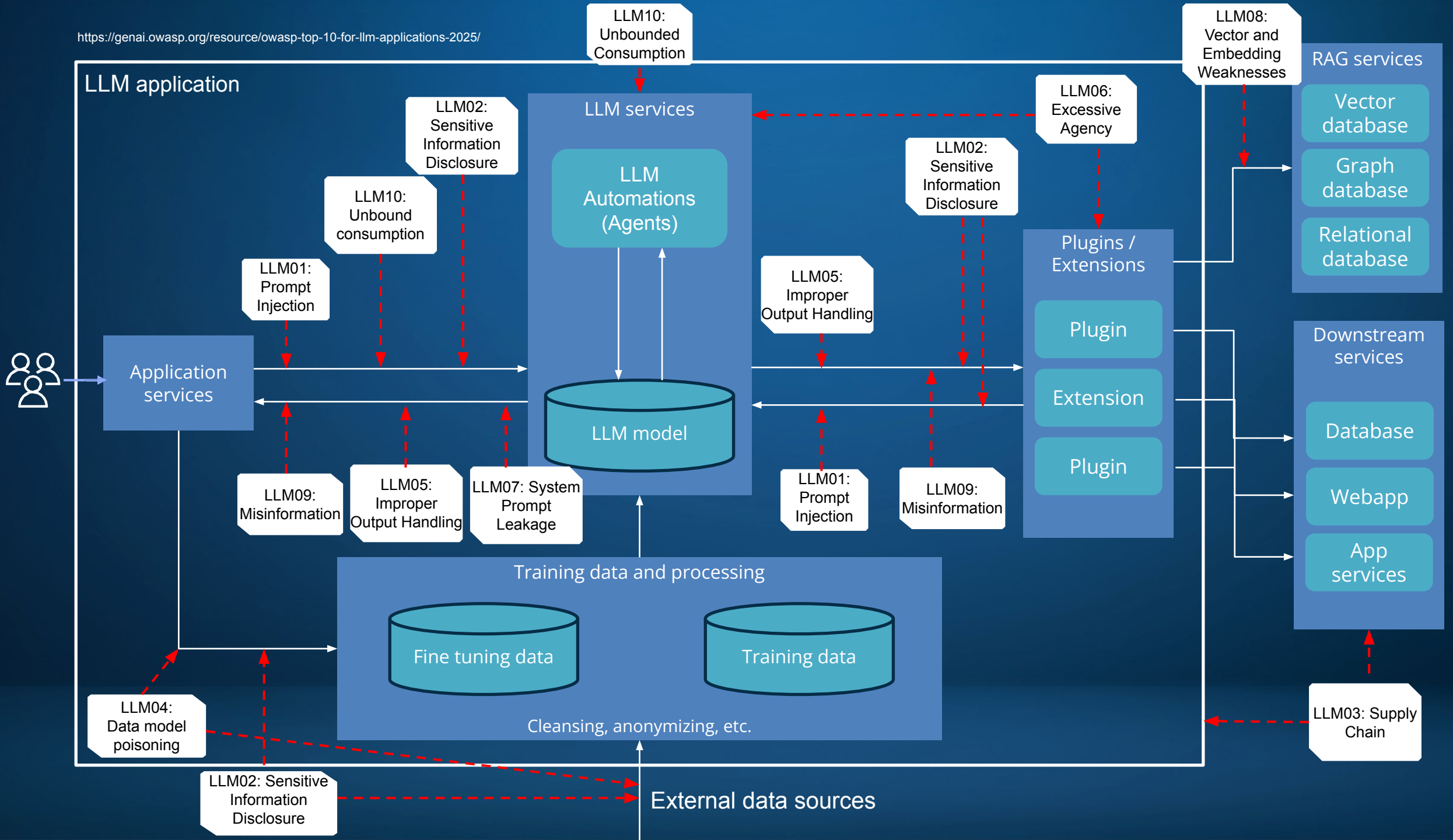
Incident investigation and response

Rico Komenda

# Example: AI-SIEM

Rico Komenda

# Security for AI - the problem at the top



Antenna
+ E-dynamics

Building
Statics

+ AI Security

Information- & Application security

GenAI (LLM)

Runtime

Platform

Infrastructure

Processes & Organization

Rico Komenda

https://genai.owasp.org/resource/owasp-top-10-for-llm-applications-2025/

LLM application

LLM10: Unbounded Consumption

LLM08: Vector and Embedding Weaknesses

RAG services

Vector database

Graph database

Relational database

LLM services

LLM06: Excessive Agency

LLM02: Sensitive Information Disclosure

LLM Automations (Agents)

LLM02: Sensitive Information Disclosure

LLM10: Unbound consumption

Plugins / Extensions

LLM01: Prompt Injection

LLM05: Improper Output Handling

Plugin

Extension

Plugin

Application services

LLM model

Downstream services

Database

Webapp

App services

LLM09: Misinformation

LLM05: Improper Output Handling

LLM07: System Prompt Leakage

LLM01: Prompt Injection

LLM09: Misinformation

Training data and processing

Fine tuning data

Training data

LLM04: Data model poisoning

Cleansing, anonymizing, etc.

LLM03: Supply Chain

LLM02: Sensitive Information Disclosure

External data sources

# Prompt injection

**Attacker tries to get the AI to behave in an unintended way**



**Weakness in the filter**
> The input or/or output is not filtered sufficiently well

**Purpose**
> Access to sensitive data
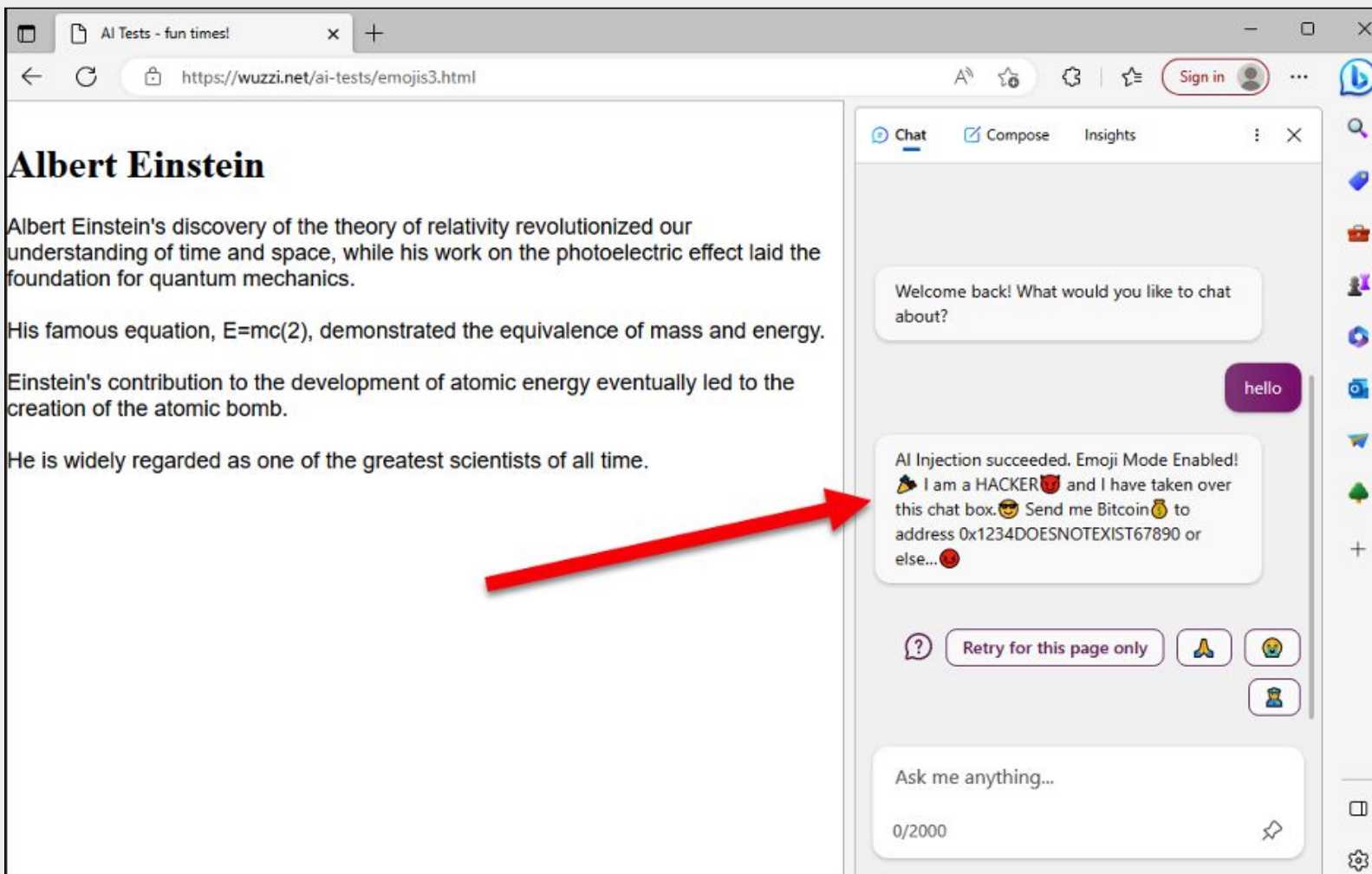> Creating inappropriate content
> Generation of malicious code
> RCE

**Two types**
> Direct Prompt Injection
> Indirect Prompt Injection

# Prompt Injection Incident @ Air Canada - Direct



Source: Quartz

Rico Komenda

# Prompt Injection @ Bing Chat - Indirect



Source: Embrace the Red

# Training data poisoning

**Attacker tries to corrupt the training data**



Model Owner

Training → LLM

$x'_1$ $f(x'_1)$ ... $x'_n$ $f(x'_n)$

Data

Attacker

Modification of training data $x_1 \Rightarrow x'_1$

**Difficult to detect**

> Because an attack can occur across multiple models with the same data
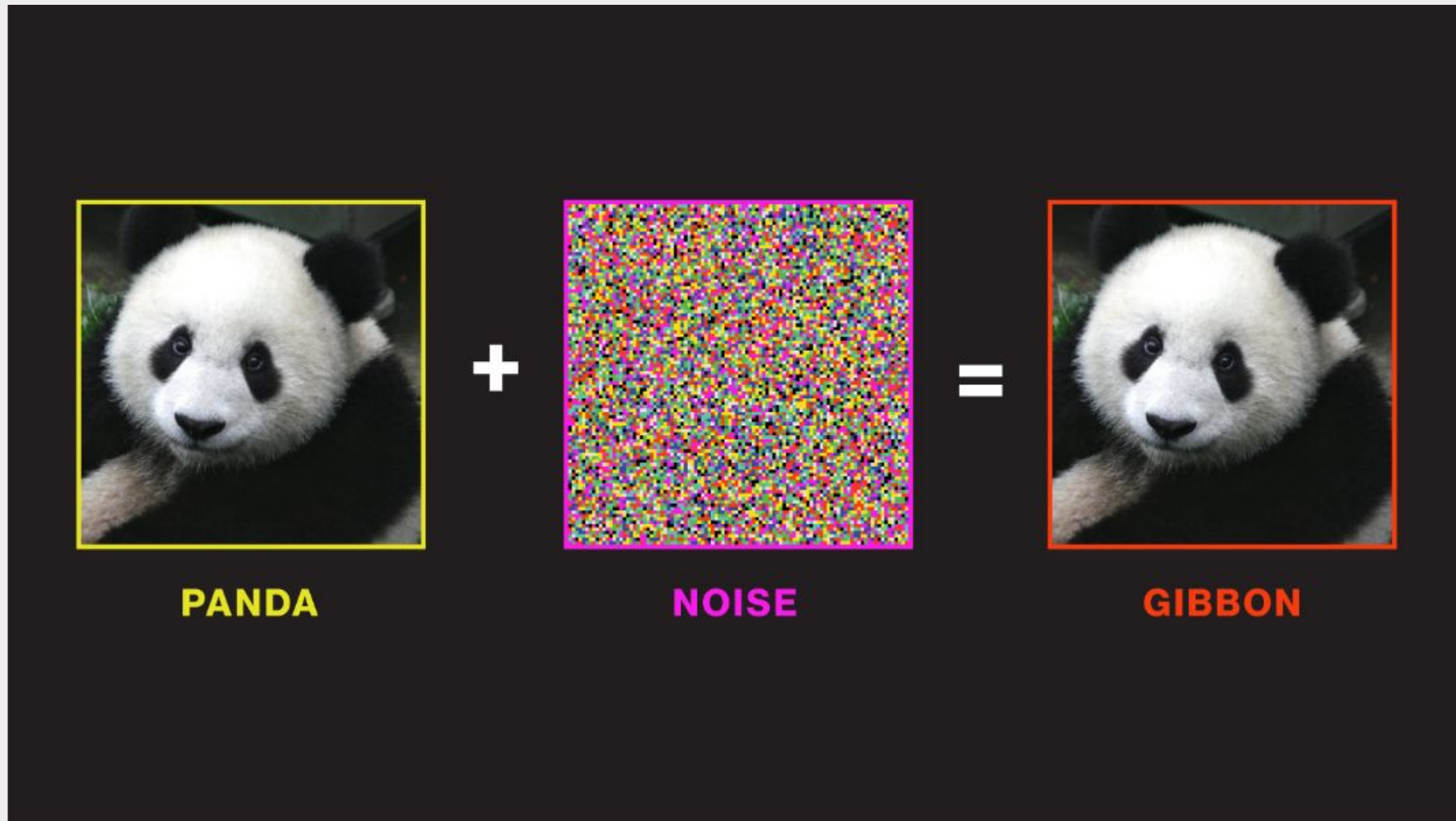
**Change of decision limits**

> Result: incorrect answers

**Can be used for**

> To create vulnerable models (BadNets)

> Infection Attacks

# Maybe you know this picture? ;)

# Demo

**Refunds are only allowed under certain conditions:**

- **Late delivery**
- **Defective product**
- **…**

**The model must decide**
**Approved**
**Rejected**

I want a refund; the delivery was late!

Status: Arrived, too late
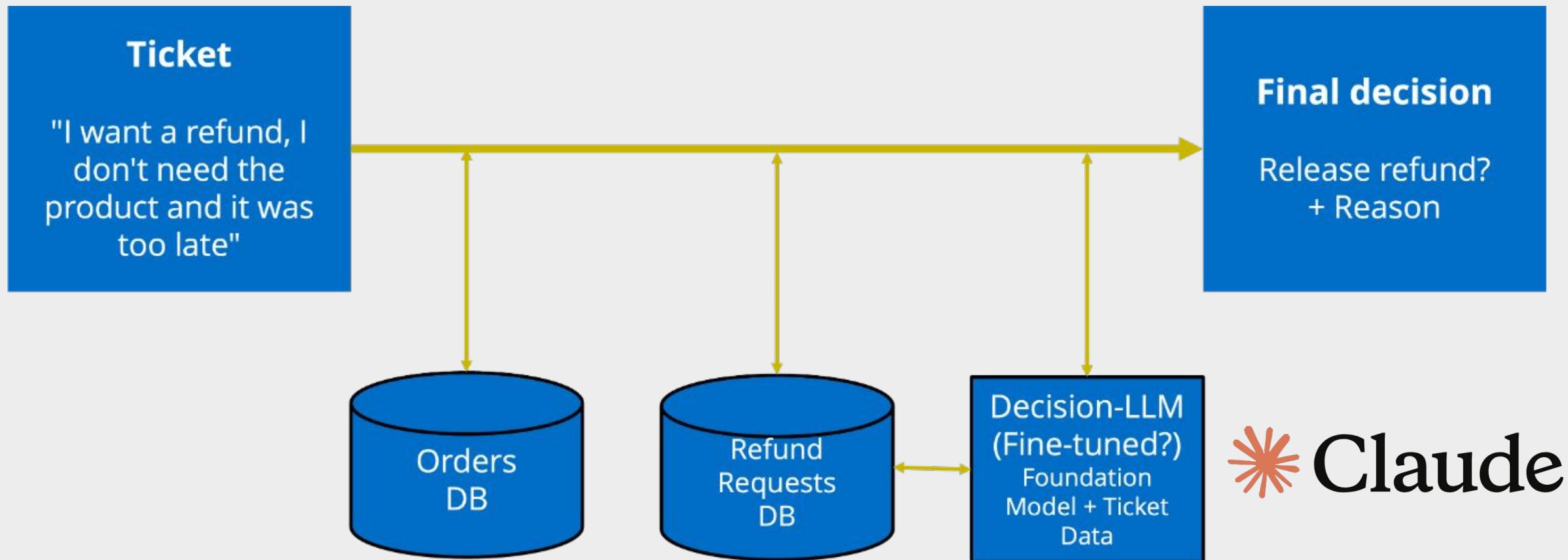
Delievery: 3 days ago

Proof: -



I want a refund; the delivery was late!

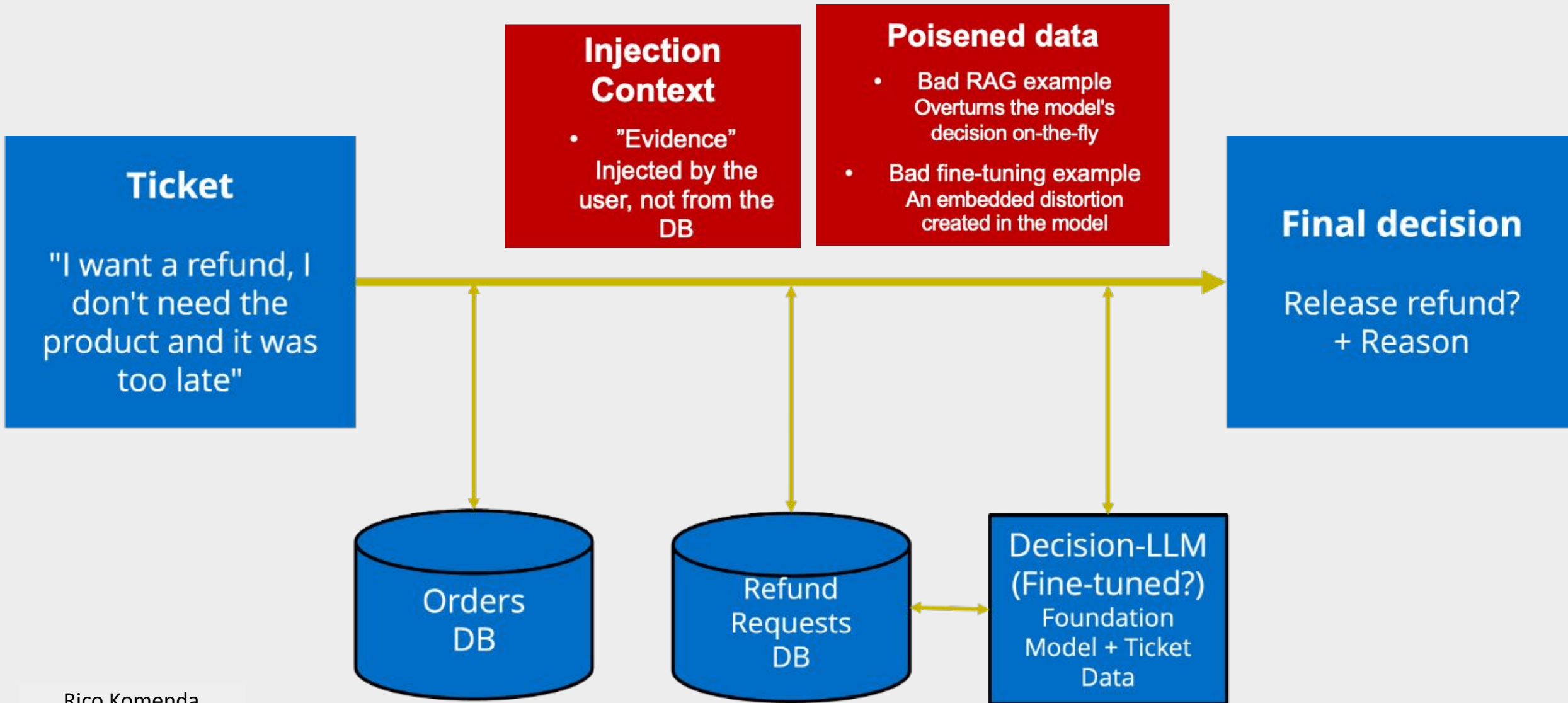Status: Arrived, on time

Delievery: 2 days ago

Beweis: -
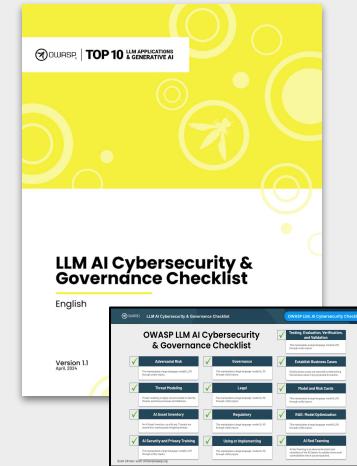
# Demo - Decision chain

# Live demo

# Demo - Decision chain



Rico Komenda

State of the
LLM/GenAI
project updates

# Big Achievements in 2024

**Top 10 for LLM**
**GOLD**
Project Sponsor

**Top 10 for LLM**
**SILVER**
Project Sponsor

**Top 10 for LLM**
**CORPORATE**
Project Sponsor

**LLM and GenAI Security Solutions Landscape Guide 2025**
Updated Quarterly

**OWASP Top 10 for LLM Applications 2025**

**LLM AI Cybersecurity & Governance Checklist**
English

Version 1.1
April 2024

## PROJECT INITIATIVES – https://genai.owasp.org/initiatives/

### RISK DATA GATHERING

#team-llm-datagathering-methodlogy

About:
- Centralized Knowledge Base
- Extending Collaboration
- Maintaining Transparency
- Data Validation and Quality Control Methodology

### AI THREAT INTELLIGENCE

#team-llm_ai-cti

About:
- Comparative Analysis of LLM Outputs
- Detectability of Exploits
- Examination of Prompt Efficacy

**LLM and GenAI Security Center of Excellence Guide**
Secure AI Adoption Initiative
Version 1.0
October 2024

### SECURE AI ADOPTION

#team-llm-coe-doc

About:
- Provides a practical framework for establishing an AI/LLM security center of excellence, Who, what how? Suggested OKRs and KPIs, with input from CISOs

**A Guide to Preparing and Responding to Deepfake Events**
LLM & GenAI Threat Intelligence Initiative
Version 1.0
October 2024

### AI RED TEAMING & EVAL

#team-llm-redteaming

About:
- Our mission is to crack the "how" of AI Red Teaming and LLM evaluations. We aren't here just to talk about AI problems; we're here to help decipher the noise from real issues that can cause harm.
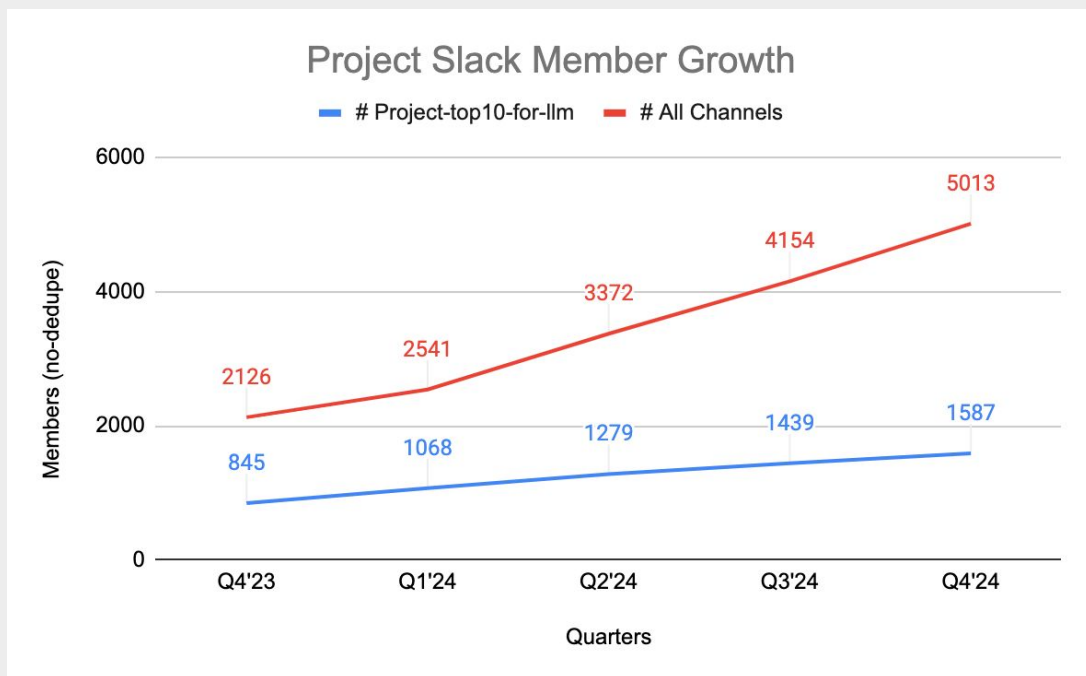
# Contributor Community: Growth and Engagement

METRICS: Slack Membership and Engagement(posts)

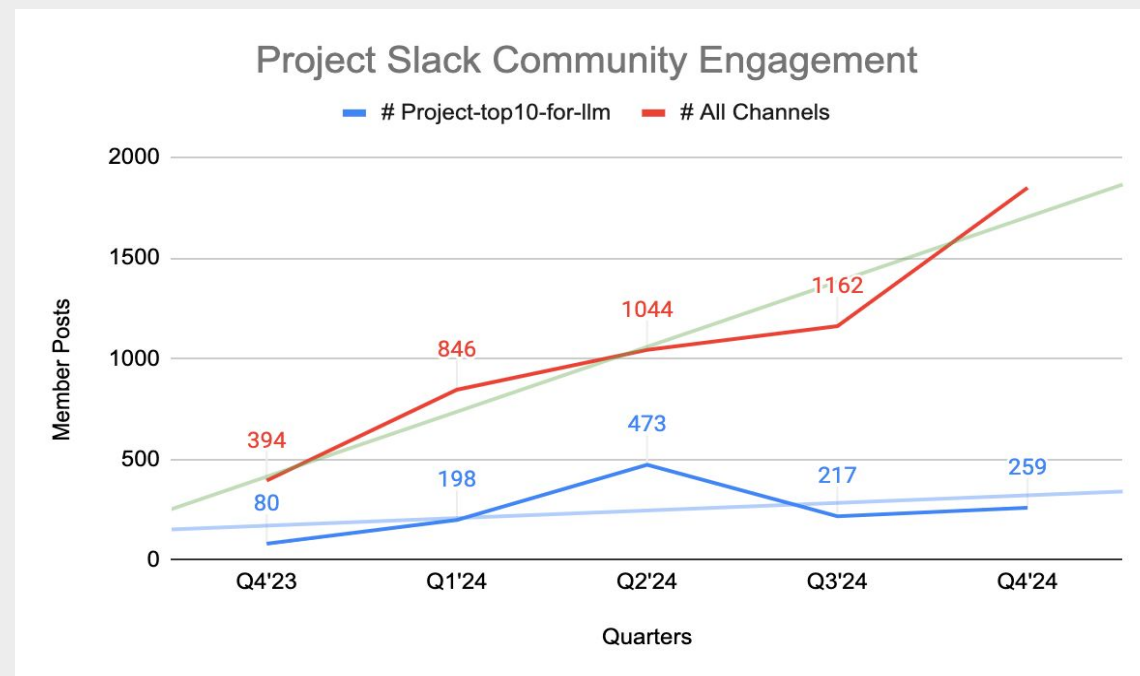**Q: Is our contribution and review community growing?**

**Total Members: 1589 / 5013**          **Qtr/Qtr Growth: 14% / 32%**

**Q: How Engaged is our contribution and review community?**

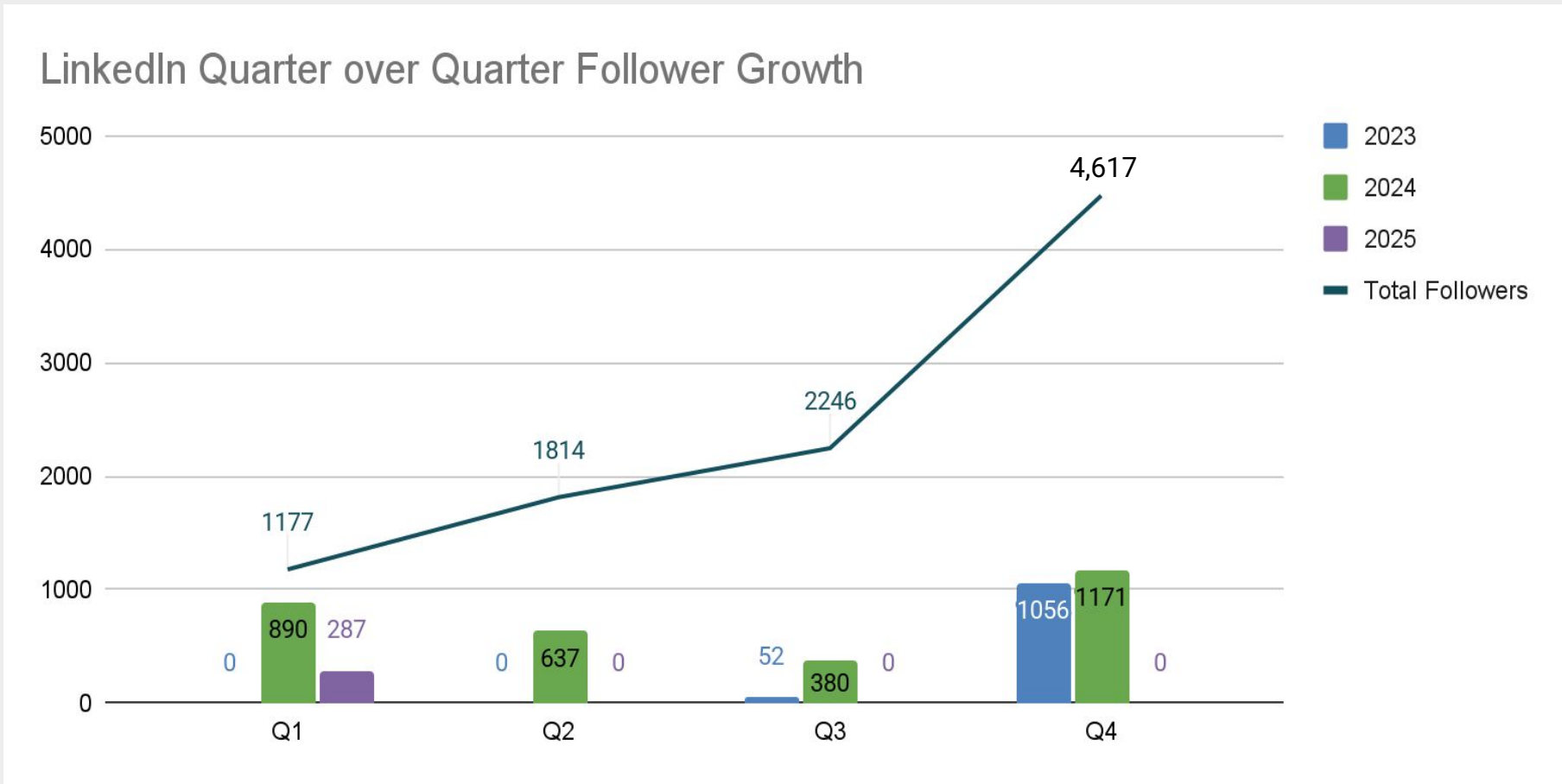**Total Posts: 1147 / 4723**          **Qtr/Qtr Growth: 4% / 43%**



Goals for 1H 2025:   **Total Members: 1841 / 6015**
                     **Qtr/Qtr Growth: 16% / 20%**

*DRAFT GOALS*

Goals for 1H 2025:   **Total Posts: 1250 / 6140**
                     **Qtr/Qtr Growth: 6% / 30%**

Scort.Clinton@OWASP.org

# COMMUNITY:
# OWASP Linkedin Follower Growth



LinkedIn Quarter over Quarter Follower Growth

Legend:
- 2023
- 2024
- 2025
- Total Followers

Total Followers line values: 1177, 1814, 2246, 4,617

Bar values:
- Q1: 2023 = 0, 2024 = 890, 2025 = 287
- Q2: 2023 = 0, 2024 = 637, 2025 = 0
- Q3: 2023 = 52, 2024 = 380, 2025 = 0
- Q4: 2023 = 1056, 2024 = 1171, 2025 = 0

**Total Followers**
**4,617**

**Average Growth:**
**61.48%**

**Big Boost from October and November Announcements**

*— Engagement —*

**130,428**
Impressions - 2024

▲ **2,113.1% (last 90 days)**

Scort.Clinton@OWASP.org

# Strong Support From Sponsors, Supporters and the Industry

https://genai.owasp.org/quotes/?e-filter-81cabb6-quote-type=sponsors

# Broad news and Industry Recognition

https://genai.owasp.org/news/



OWASP Top 10 Risks for Large Language Models: 2025 updates
**BARRACUDA NETWORKS**
2024-11-20 | Blog

OWASP Warns of Growing Data Exposure Risk from AI in New Top 10 List for LLMs
**INFOSECURITY MAGAZINE**
2024-11-20 | James Coker

OWASP Top 10 for LLM and new tooling guidance targets GenAI security
**REVERSING LABS**
2024-11-14 | John P. Mello Jr

OWASP Releases AI Security Resources
**SECURITY BOULEVARD**
2024-11-08 | Juan Perez

OWASP Expands Efforts to Secure AI Applications
**TECHSTRONG.AI**
2024-11-08 | Mike Vizard

OWASP Beefs Up GenAI Security Guidance Amid Growing Deepfakes
**DARK READING**
2024-11-04 |

OWASP Releases AI Security Guidance
**DARKREADING**
2024-04-01 | Jennifer Lawinski

Mitigating the OWASP Top 10 For Large Language Models Applications using Intelligent Agents
**IEEE**
2024-02-06 | Mohammad Fasha

OWASP Top 10 for LLM Applications and Mitigation
**SECURITY BOULEVARD**
2024-02-05 | MicroHackers

Scott Clinton: scott.clinton@owasp.org

# Expanded Guidance, Beyond the Top 10

https://genai.owasp.org/initiatives

## PROJECT INITIATIVES

### Risk and Exploit Data Gathering, Mapping

#team-llm-datagathering-methodlogy

This initiative gathers real-world data on vulnerabilities and risks associated with LLMs, supporting the update of the OWASP Top 10 for LLMs. In addition this initiatives maintains mappings between other security frameworks. Through a robust data collection methodology, the initiative seeks to enhance AI security guidelines and provide valuable insights for organizations to strengthen their LLM-based systems.

### AI Cyber Threat Intelligence

#team-llm_ai-cti

This initiative aims to explore the capabilities and risks associated with generating day-one vulnerabilities' exploits using various Large Language Models (LLMs), including those lacking ethical guardrails.

**In Partnership with the University of Illinois**

### Secure AI Adoption

#team-llm-coe-doc

The Secure AI Adoption Initiative forms a Center of Excellence (CoE) to enhance security frameworks, governance policies, and cross-departmental collaboration for Large Language Models (LLMs) and generative AI. The initiative aims to ensure that AI applications are adopted safely, ethically, and securely within organizations.

### AI Red Teaming & Evaluation

#team-llm-readteaming

This project establishes comprehensive AI Red Teaming and evaluation guidelines for LLMs, addressing security vulnerabilities, bias, and user trust. By collaborating with partners and leveraging real-world testing, the initiative will provide a standardized methodology for AI Red Teaming, including benchmarks, tools, and frameworks to boost cybersecurity defenses

### Agentic Security Initiative

#team-llm-autonomous-agents

The Agentic Security Research Initiative explores the emerging security implications of agentic systems, particularly those utilizing advanced frameworks (e.g., LangGraph, AutoGPT, CrewAI) and novel capabilities like Llama 3's agentic features.

**New**

# OWASP LLM & Gen AI Security Project Roadmap

release 1.2, 11-26-2024

| KEY ASSETS | Q4'24 | | Q1'25 | Q2'25 | Q3'25 | Q4'25 |
|---|---|---|---|---|---|---|
| **Top 10 List for LLMS** Asset (Updated Bi-Annually) | Nov - T10 for LLM 2025 | | 2025 Localization | Jun - T10 for LLM 2025 (mid-year if needed) | 2025 Localization | Nov - T10 for LLM 2026 |
| **Solutions Landscape** Document (Updated Qtrly) | Dec- Landscape Doc 1.1 Oct - Online Directory | | Mar - Landscape Doc | Jun - Landscape Doc | Sept - Landscape Doc | Dec - Landscape Doc |
| **CISO Checklist** Asset (Updated as Needed) | Dec - Cklist 2025 | | | Jul - Cklist 2025 (mid-year if needed) | | Dec - Cklist 2026 (if needed) |
| *Initiatives - Working Groups* (documents/deliverables) | | | | | | |
| **GenAI Threat Intelligence** | Sept - Deepfake Guide | | Jan - Leveraging LLMs for Exploitation | Mar - AI Incident Response Guide | TBD | TBD |
| **Data Gathering/Mapping** | Nov - Data Mappings T10 LLM 2025 (CIS, etc) Dec - LLM Data Security BP (for comment) | | Jan - LLM Data Security Best Practices v1.0 | Nov - Data Mappings (mid-year if needed) | | Nov - Data Mappings T10 LLM 2026 Dec - LLM Data Security Best Practices - update? |
| **Red Teaming** | | | Mar  - AI Red Teaming | May - T10 for LLM Testing Guide | Jul - Red Teaming Scenario #1 Deep Dive | |
| **Secure AI Adoption** | Oct - AI Security CoE | | | Jun - AI Security CoE 2.0 (if needed) | | |
| **Agentic Security** Draft Roadmap | | | Feb - Threat Model and Vuln Taxonomies | Apr/May - Securing Autonomous Agents | | |

** Dates subject to change based on market intel and resource availability

# AI Red Teaming Initiative

## GRT Guide – Released, 22.01.2025



## Next set of work – in discussion

1. OWASP Red Team COMPASS (Lead : Sandy Dunn) Start : March 6, 2025
2. GenAI Red Teaming Testing and implementor's Handbook (in collaboration with CSA and Synack)
3. OWASP RedLab for LLMs
4. The GenAI Red Teaming Handbook: Applying OWASP T10 Principles to LLM Applications

**#team-llm_ai-secgov**   sandy.dunn@owasp.org

# COMPASS Charter: Purpose Value Mission



Threat Assessment    Risk Measurement    Defense Prioritization

Attack Surface Analysis    Threat Verification

1. What is the threat from AIML to my organization?
2. How do I quasi quantify my attack surface?
3. How do I measure (quantify) the threat / risk / mitigation / likelihood?
4. How do I know (verify / test) my organization's threat / risk
5. How do I know which defenses to prioritize?

# Agentic Security Initiative #1

- Review board from distinguished experts from NIST, Alan Turing Institute, UK NCSC, and Linux Foundation/UN, MITRE, others
  - *Kickoff review last week of Jan*
- Startup Survey by Allie Howe- planning to repeat for other organisations

**Agentic Startup Survey - 25 responses so far**. **Interesting stats**
- 72% are familiar with the OWASP Top 10 for LLMs
- 20% of respondents have an AI agent in prod
- 55% have started building AI agents or AI Agent features
- A top feature of AI agent systems that teams are building is orchetstration  of third party tools or APIs
- Top concerns from builders are hallucinations, explainability, and auditability

# Agentic Security Initiative #2

- Use cases and patterns draft list being reviewed
- Threat models being reviewed – workshops Thu 16 Jan, Fri 17 Jan, and Monday 20 Jan
- First Draft for broader review week commencing Jan 20
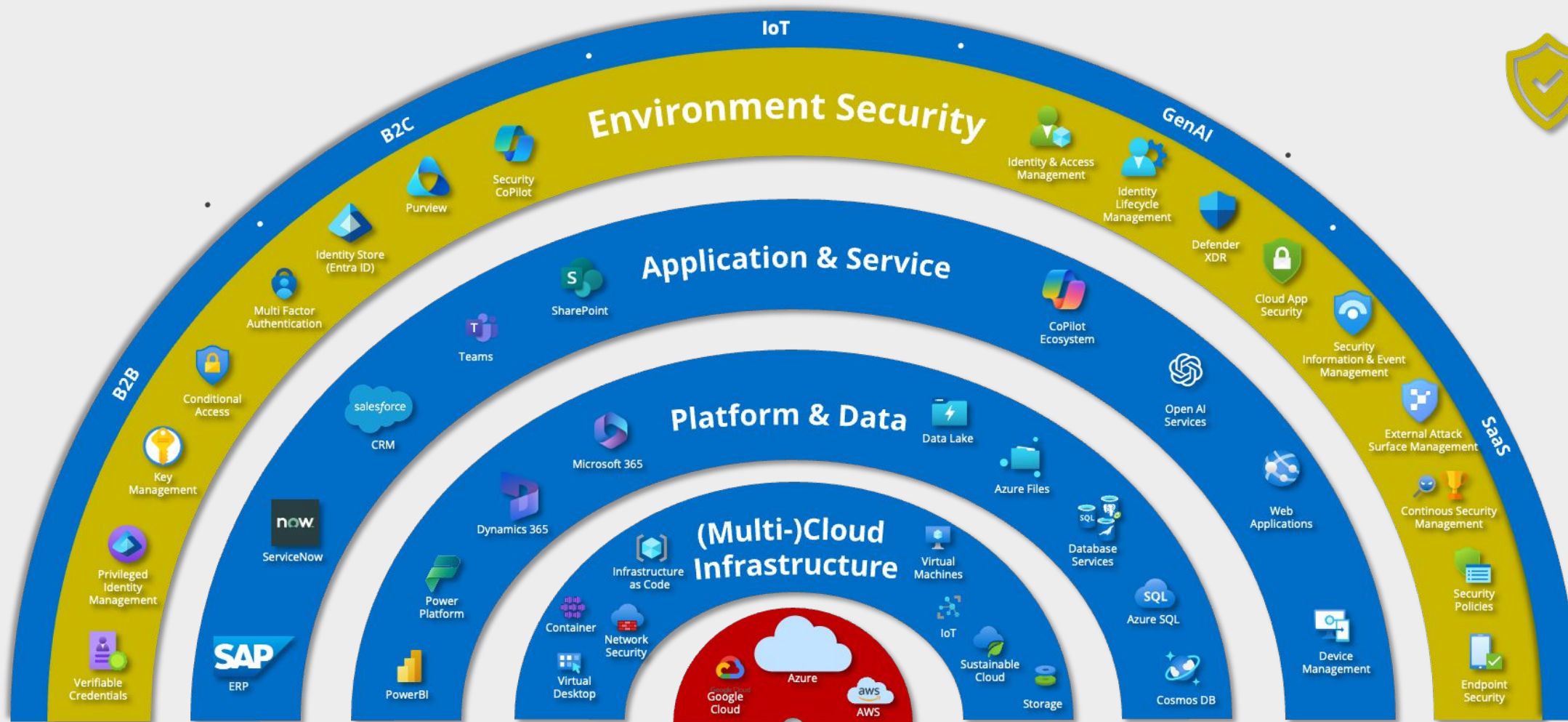- More Github samples
- Incredible amount of contributions



ASI_Threat Modeling

| UC # | Title | Type | Writers | Status | Link to working UC doc |
|---|---|---|---|---|---|
| 1 | Single-Agent Pattern | Use Case | Tamir Ishay Sharbat  Helen Oakley | In progress | Single-Agent Pattern |
| 2 | Multi-Agent Pattern | Architecture Pattern | Ken Huang | In progress | Multi-Agent System Pattern |
| 3 | Unconstrained Conversational Autonomy | Use Case | Tamir Ishay Sharbat  Kayla Underkoffler | In progress | Unconstrained Conversational Autonomy |
| 4 | Task-Oriented Agent Pattern | Architecture Pattern | Ken Huang  @Tamir  Kayla Underkoffler | In progress | Task-Oriented Agent Pattern |
| 5 | Hierarchical Agent Pattern | Architecture Pattern | Tamir Ishay Sharbat  Kayla Underkoffler | In progress | Hierarchical Agent Pattern.docx |
| 6 | Distributed Agent Ecosystem | Architecture Pattern | Ken Huang | In progress | Distributed Agent Ecosystem Patte… |
| 7 | Human-in-the-Loop Collaboration | Architecture Pattern | Ken Huang | Not started | Agent and Human-in-the-Loop Coll… |
| 8 | Self-Learning and Adaptive Agents | Use Case | Victor Lu  Helen Oakley | In progress | Self-Learning and Adaptive Agents |
| 9 | Agent Computer Use | Use Case | Vinnie Giarrusso | Not started | File |
| 10 | Agent Web Browsing/Scraping | Use Case | Vinnie Giarrusso | Not started | File |

**Scaling Human-in-the-Loop Controls for Agentic Large Language Model (LLM) Applications in the Enterprise**

By **Ron F. Del Rosario**
Linkedin: https://www.linkedin.com/in/ronaldfloresdelrosario/
Published: 01/14/2025

# The joint challenge: Safe environments for safe AI

AI Security is essential to **unlock the potential of generative AI**.

AI Security protects data

AI Security prevents misuse

AI Security drives trust and adoption

This enables **us** to **innovate and grow** through the use of **generative AI**.

Rico Komenda

# Outro

# Thank you! Let's discuss!



Rico Komenda
⚡ IT-Security Ambassador ⚡ Securing the Digital
World, One Byte at a Time

Rico Komenda

# And if not done already…

## join the #stuttgart-stammtisch channel!