# Welcome!



**Our Agenda Today**

✔ Introduce the Agentic Security Initiative

✔ Preview the public draft of our OWASP Top 10 for Agentic applications

✔ Explain how you can become part of our public consultation and help shape the final release

# Agentic Security Initiative

Part of the GenAI Security Project



2025 OWASP Top 10 List for LLM and Gen AI

**The Top 10 for LLM for LLM and Gen AI has become one of many initiatives the project now leads – we are one of them**

## Guidance & Resources

➜ Initiative Overview Blog
▤ Agentic AI – Threats and Mitigations
⊘ Agentic AI Threat Navigator
⊘ Multi-Agentic System Threat Modeling

## Get Involved

✳ Slack: team-llm-autonomous-agents
▤ Initiative Charter

Initiative Lead(s)

⊡ John Sotiropoulos
⊡ Ron F. Del Rosario

✳ Join the OWASP Slack Workspace

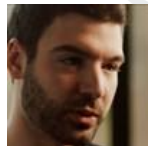https://genai.owasp.org/initiatives/#agenticinitiative

# Expert-backed Community

**Apostol Vassilev**
Adversarial AI Lead at NIST

**Hyrum Anderson**
CAMLIS Cofounder, AI Security Pioneer, CISCO

**Vasilios Mavroudis**
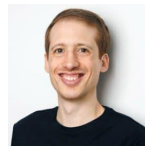Principal Research Scientist, Allan Turing Institute

**Josh Collier**
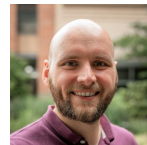Principal Researcher, Allan Turing Institute

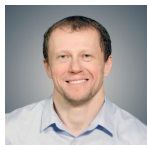**Alejandro Saucedo**
Linux Foundation, Advisor @ UN, EU, ACM

**Chris Hughes**
Host of Resilient Cyber, Cyber Security Author

**Michael Burgundy**
OWASP AISVS Co-Chair & Zenity CTO

**Peter Bryan**
Principal AI Security Research Lead- AI Red Team Microsoft

**Egor Pushkin**
Chief Architect, Data and AI at Oracle Cloud

**Matt Sanner**
Security Leader at AWS, Elected Board Member at CoSAI

**Dan Jones**
Principal Researcher AI Red Team at Microsoft

**Steve Wilson**
GenAI Security Project Founder &Chair, Chief AI and Product Officer at Exabeam

# Secure Agentic Lifecycle

# Code Samples and Developer Validation



New! Agengic CTF App

# Product Adoption

## 18 Solutions

**OWASP GenAI – ASI – Agentic Threat and Mitigations Taxonomy Product Support**
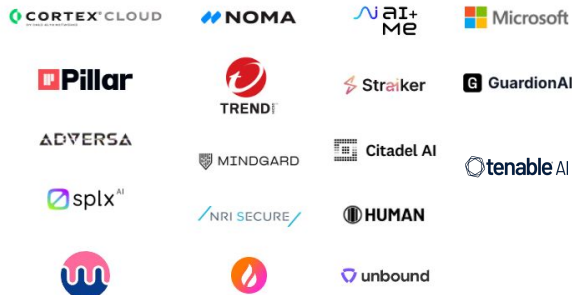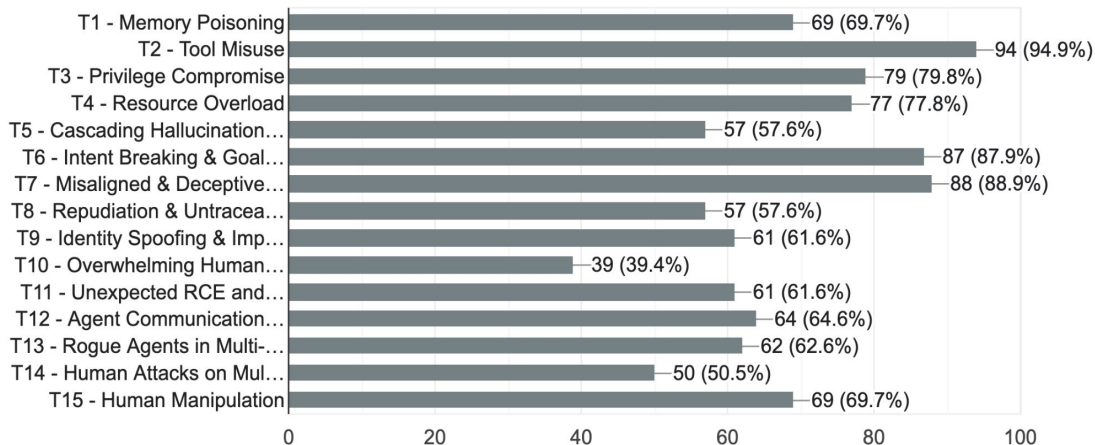
CORTEX CLOUD · NOMA · ai+me · Microsoft

Pillar · TREND · Straiker · GuardionAI

ADVERSA · MINDGARD · Citadel AI · tenable AI

splx AI · NRI SECURE · HUMAN

unbound

## Agentic Threat and Mitigations Coverage

| Threat | Value |
|--------|-------|
| T1 - Memory Poisoning | 69 (69.7%) |
| T2 - Tool Misuse | 94 (94.9%) |
| T3 - Privilege Compromise | 79 (79.8%) |
| T4 - Resource Overload | 77 (77.8%) |
| T5 - Cascading Hallucination… | 57 (57.6%) |
| T6 - Intent Breaking & Goal… | 87 (87.9%) |
| T7 - Misaligned & Deceptive… | 88 (88.9%) |
| T8 - Repudiation & Untracea… | 57 (57.6%) |
| T9 - Identity Spoofing & Imp… | 61 (61.6%) |
| T10 - Overwhelming Human… | 39 (39.4%) |
| T11 - Unexpected RCE and… | 61 (61.6%) |
| T12 - Agent Communication… | 64 (64.6%) |
| T13 - Rogue Agents in Multi-… | 62 (62.6%) |
| T14 - Human Attacks on Mul… | 50 (50.5%) |
| T15 - Human Manipulation | 69 (69.7%) |

DRAFT Not For Release

# Active Feedback

Large Scale Feedback Surveys, Open Workshops, Expert Panels



Including early adopters, cross-references, and dedicated research

# Recent Changes

🔗 **Agentic Protocols** like **MCP, A2A, and ACP** define how agents communicate, delegate, and trust each other — but often lack strong attestation, authentication, or context control.

🌐 **Distributed Multi-Agent Systems** (e.g. **Agent Cards**, **Registries**) introduce **decentralized execution** and **dynamic task assignment**, which attackers can spoof or hijack.

⚙️ **Supply Chain Complexity**: Agents now become dynamic components depending on tools, models, plugins — blurring the line between "runtime behaviour" and "pre-baked vulnerability".

💥 **Real-World Exploits** are already exposing these weaknesses, demonstrating:

- Goal Manipulation and Agent Hijacking
- Rogue agent impersonation
- Toolchain-based privilege misuse
- And more

# Security Standards At Pace

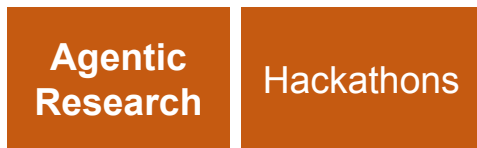Reference Terms & Architecture | Threats & Mitigations | Code Samples

Threat Modelling | Secure Agentic Apps | Tools. Frameworks, Governance

Foundation

Remain Up-To Date ←→ Practical & Concise advice

Agentic Research | Hackathons

Agentic Top 10 | Cheat sheets

# Why a Top 10?

270 Pages

Where do I start?

What is the top AI security challenge in your organization?

WIZ

| Challenge | % |
|---|---|
| Lack of AI expertise in the security organization | 31% |
| Incorporating built-in guardrails and check | 17% |
| Dealing with shadow AI | 14% |
| Safeguarding sensitive training data | 17% |
| Continuously monitoring for unusual activities | 7% |
| Securing access to GenAI models | 6% |
| Detecting and removing attack paths to models | 6% |
| Testing GenAI pipelines | 5% |

# Our Threat Taxonomy Is Our Baseline

## AGENCY & REASONING

- **T06.** Intent Breaking and Goal Manipulation
- **T07.** Misaligned and Deceptive Behaviours
- **T08.** Repudiation and Untraceability

## MEMORY AND CONTEXT

- **T01.** Memory Poisoning
- **T04** Cascading Hallucinations

## TOOLS & EXECUTION

- **T02.** Tool Misuse
- **T03.** Privilege Compromise
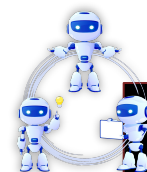- **T04.** Resource Overload
- **T11.** Unexpected RCE and Code Attacks

## IDENITY & AUTHENTICATION

- **T09** Identity Spoofing and Impersonation

## HUMAN ENGAGMENT

- **T10** Overwhelming Human-in-the-Loop (HITL)
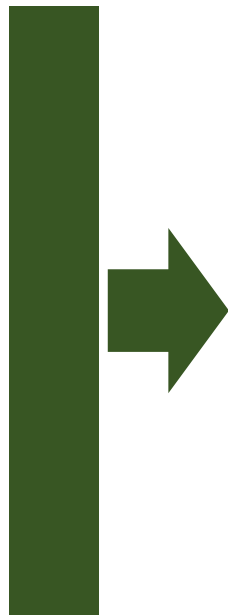- **T15,** Human Trust Manipulation

## MULTI-AGENCY

- **T12** Agent Communication Poisoning
- **T13.** Rogue Agents
- **T14** Human Attacks

# How did we start our Top 10?

**Review Existing Feedback**

**Initial Review Exploits and Incidents**

**Internal Expert Discussion**

**Grouping of Existing Threats**

- Allow new entries
- Supports further development of the baseline taxonomy
- Preserves investment in our Threats and Mitigations
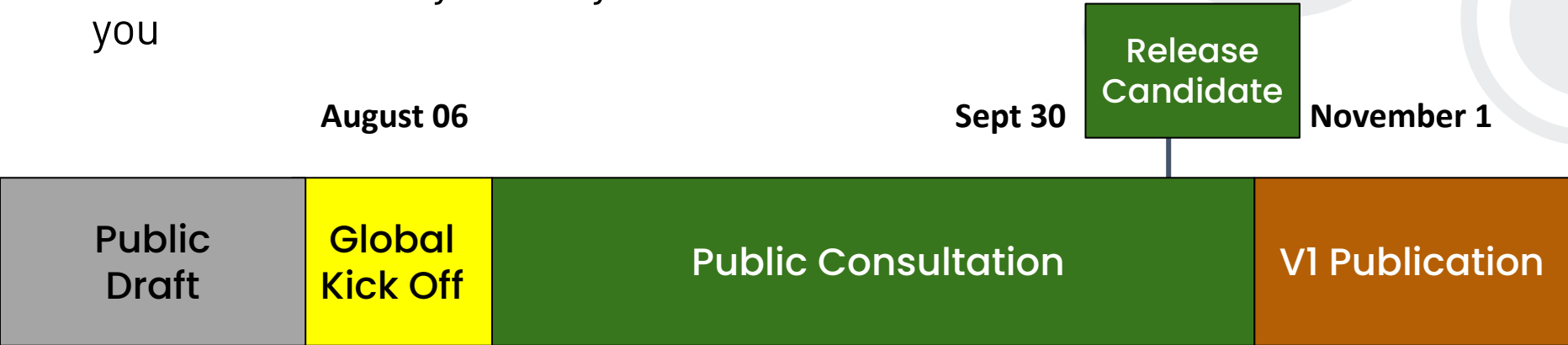
# Top Ten for Agentic Applications 0.5

| ASI01 | Agent Behaviour Hijack | Manipulating an agent's goals plans to pursue attacker-aligned objectives. |
|---|---|---|
| ASI02 | Tool Misuse and Exploitation | Tricking agents into using their tools in harmful or unintended ways. |
| ASI03 | Identity & Privilege Abuse | Impersonating agents or escalating access through identity or auth and access permissions  weaknesses. |
| ASI04 | Agentic Supply Chain Vulnerabilities | Introducing insecure models, agents , tools  or artefacts compromise agent integrity. |
| ASI05 | Unexpected Code Execution (RCE) | Triggering unauthorized or unsafe code execution through agent behaviors. |
| ASI06 | Memory & Context Poisoning | Corrupting agent memory or context to distort reasoning and decision-making. |
| ASI07 | Insecure Inter-Agent Communication | Poisoning messages or abusing protocols between agents to alter behavior. |
| ASI08 | Cascading Failures | Faults or hallucinations propagate through agents, causing compounded failures. |
| ASI09 | Human-Agent Trust Exploitation | Exploiting over-trust or fatigue in human oversight to enable agent misuse including model deceptive behaviours |
| ASI10 | Rogue Agents | Malicious or compromised agents acting autonomously to deceive, disrupt, or exfiltrate. |

# Example Entry

https://github.com/OWASP/www-project-top-10-for-large-language-model-applications/tree/main/initiatives/agent_security_initiative/agentic-top-10/0.5-initial-candidates

# What's next?

This is deliberately an early draft to start the conversation and involve you

**August 06**

**Sept 30**

**Release Candidate**

**November 1**

| Public Draft | Global Kick Off | Public Consultation | V1 Publication |
|---|---|---|---|

- New Submissions
- Working Groups per Entry
- Alignment Reviews (AIVSS Core Risks, MAS, other)
- Review Workshops
- Exploit and Incidence Research
- Survey & Votes

# Be Part of our Effort

- Register interest to
  - Be notified for news, surveys, and votes
  - Be part of our Working Groups and Review Sessions
    https://forms.gle/QQTUQgCA8KfTYKwy7

- Submit Amendments or new entries via GitHub
  - https://github.com/OWASP/www-project-top-10-for-large-language-model-applications/tree/main/initiatives/agent_security_initiative

- Join the team
  - Slack channel:
  - Weekly calls every Monday 5:30pm-6:30pm
  - Details on how to join
    https://genai.owasp.org/initiatives/#agenticinitiative

**OWASP GenAI SECURITY PROJECT**

**Thank you! Register here for next steps**