



Navigating AI Threats & Security

AI Threat Defense COMPASS v1.0

Sandy Dunn
Sabrina Caplis



About



Sandy Dunn SPLX

- 20 + yrs MicronPC, HP, Blue Cross, startups, board member
- CISO SPLX
- Core Contributor OWASP GenAI Security Project
- Adjunct Professor BSU

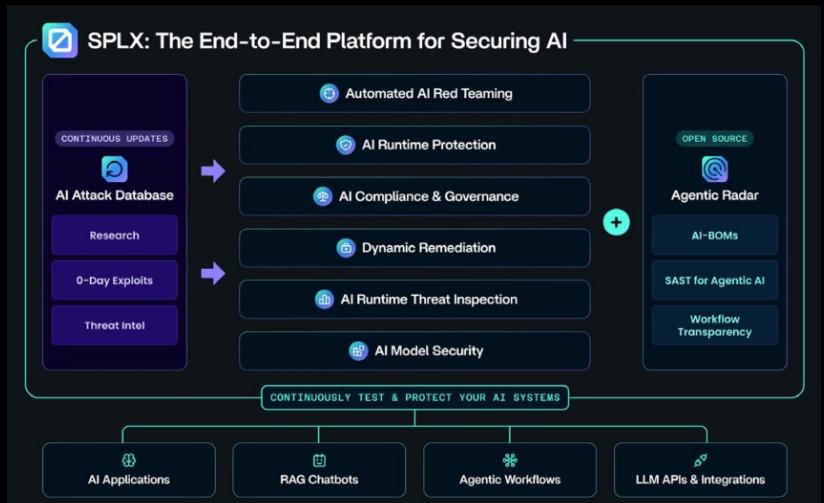


Sabrina Caplis

- AI Governance & Security Consultant at RockCyber
- Board Member - ISACA and ISSA Denver
- OWASP GenAI Security Project Contributor



AI Security & Automated Red Teaming



Security AI Compliance Business Continuity



Innovation

We continuously seek out new ideas, technologies, and approaches to stay ahead of emerging cybersecurity threats and provide cutting-edge solutions to our clients.



Adaptability

In the fast-paced world of cybersecurity, adaptability is key. We remain agile and flexible in our approach, quickly responding to changing threats and evolving client needs to ensure optimal outcomes in every situation.



Client-centric

Our clients are at the heart of everything we do. We are dedicated to understanding their unique challenges, and tailoring our solutions to meet their specific needs, delivering value that exceeds expectations.

THE CYBER CRISIS



Idaho drivers receive toll bills from California's FasTrak for violations they never committed



NAVIGATING THE AI FRONTIER

My first prompt,
wish me luck!



AI says 99.8%
chance of bear encounter.
Enjoy your
'immersive' experience!



RICHES & OPPORTUNITIES



AI FRONTIER

DANGERS & UNKNOWN



INNOVATION, PRODUCTIVITY,
PERSONAL EMPOWERMENT

BIAS, MISFORMATION,
ADVERSARIAL ATTACKS
CHAOS



Deployment Types



Vendor-Managed Models



Type 1: Embedded in Application
AI features included in a product,
e.g., Zoom, Canva



Type 2: Access via Vendor UI
Model accessed via the vendor
platform, e.g., ChatGPT, Gemini



Type 3: Model API Access Use the
model through an API, e.g., OpenAI
API, Anthropic API



Type 4: Enterprise Licensed Model
Enterprise version of a model in tenant
e.g., Microsoft Copilot



Enterprise-Controlled Models



Type 5: Pre-Trained + Fine-Tune
Fine-tune a general foundation
model with enterprise data



**Type 6: Fine-Tuned Specialized
Model** Adapt a domain-specific
model with custom fine-tuning



Type 7: Fully Custom Model
Develop a tailored model
architecture for the enterprise

OWASP

**GenAI SECURITY
PROJECT**
genai.owasp.org



mikeprivette

SECURITY
PROJECT

AI Security Shared Responsibility Matrix

Security Domain	SaaS AI	PaaS AI	IaaS AI	On-Premises	Embedded AI	Agentic AI	AI Coding	MCP Systems
Application Security	Shared	Customer	Customer	Customer	Shared	Shared	Customer	Customer
AI Ethics and Safety	Provider	Shared	Customer	Customer	Provider	Shared	Provider	Customer
User Access Control	Shared	Shared	Customer	Customer	Shared	Customer	Customer	Customer
Model Security	Provider	Shared	Customer	Customer	Provider	Shared	Provider	Shared
Data Privacy	Shared	Shared	Customer	Customer	Shared	Customer	Customer	Customer
Data Security	Shared	Shared	Customer	Customer	Shared	Customer	Customer	Customer
Monitoring and Logging	Shared	Shared	Customer	Customer	Provider	Shared	Customer	Customer
Compliance and Governance	Shared	Shared	Customer	Customer	Shared	Customer	Customer	Customer
Supply Chain Security	Provider	Shared	Shared	Customer	Provider	Shared	Provider	Shared
Network Security	Provider	Shared	Customer	Customer	Provider	Customer	Customer	Customer
Infrastructure Security	Provider	Provider	Shared	Customer	Provider	Shared	Customer	Shared
Incident Response	Shared	Shared	Customer	Customer	Shared	Shared	Customer	Customer
Agent Governance ★	N/A	Shared	Customer	Customer	Shared	Customer ▲	N/A	Customer
Code Generation Security ★	N/A	N/A	N/A	N/A	N/A	N/A	Customer ▲	N/A
Context Pollution Protection ★	Shared	Shared	Customer	Customer	Shared	Shared	Customer ▲	Customer ▲
Multi-System Integration ★	Shared	Shared	Customer	Customer	Shared ▲	Shared	Customer ▲	Customer ▲

genai:owasp.org

Types of AI Attacks

GenAI User Attacks

- Prompt Injection
- Hallucinations
- Toxic
- Bias
- Supply Chain

GenAI System Attacks

- Model Operations
- Supply Chain Attacks
- Jailbreaking
- Prompt Leakage
- API Security
- Plugin Security
- Supply Chain

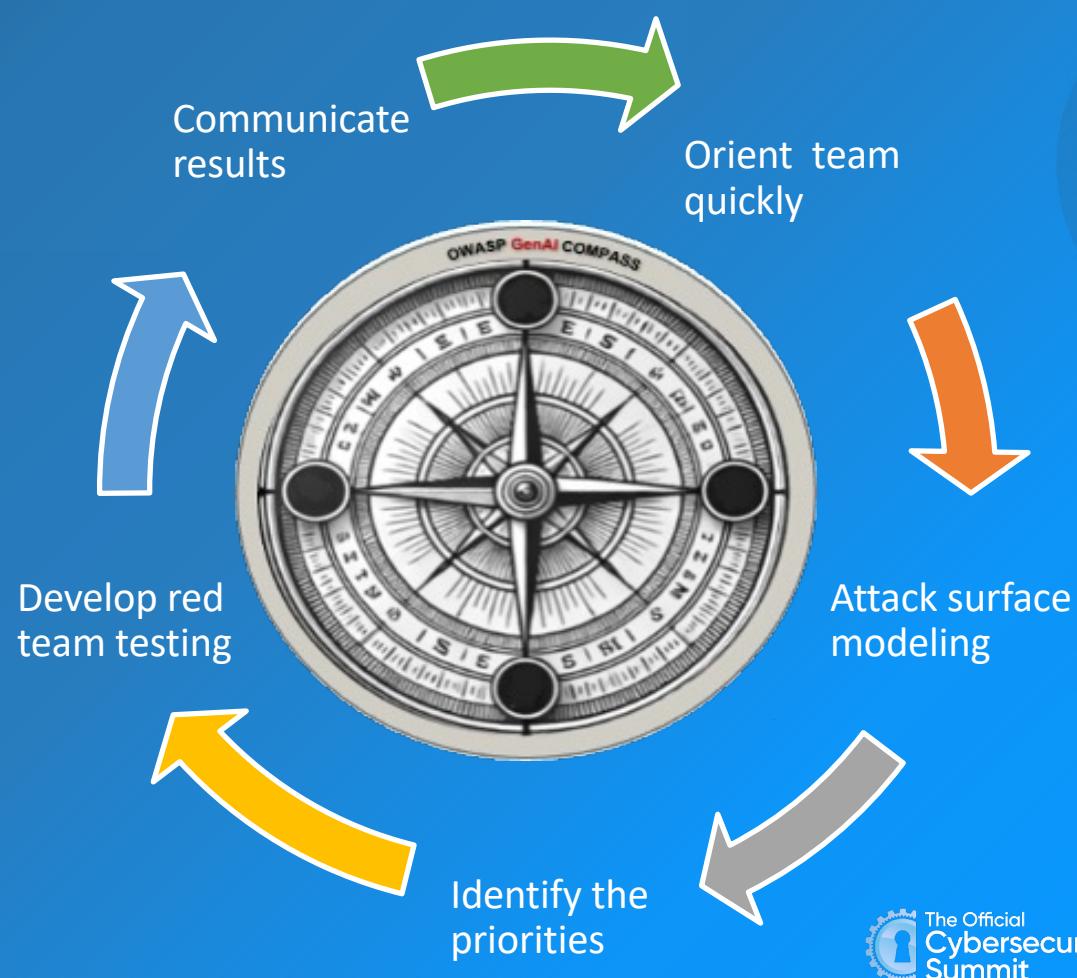
Model Attacks

- Model Poisoning
- Model Evasion
- Model Extraction
- Inference
- Privacy Leaks
- Supply Chain

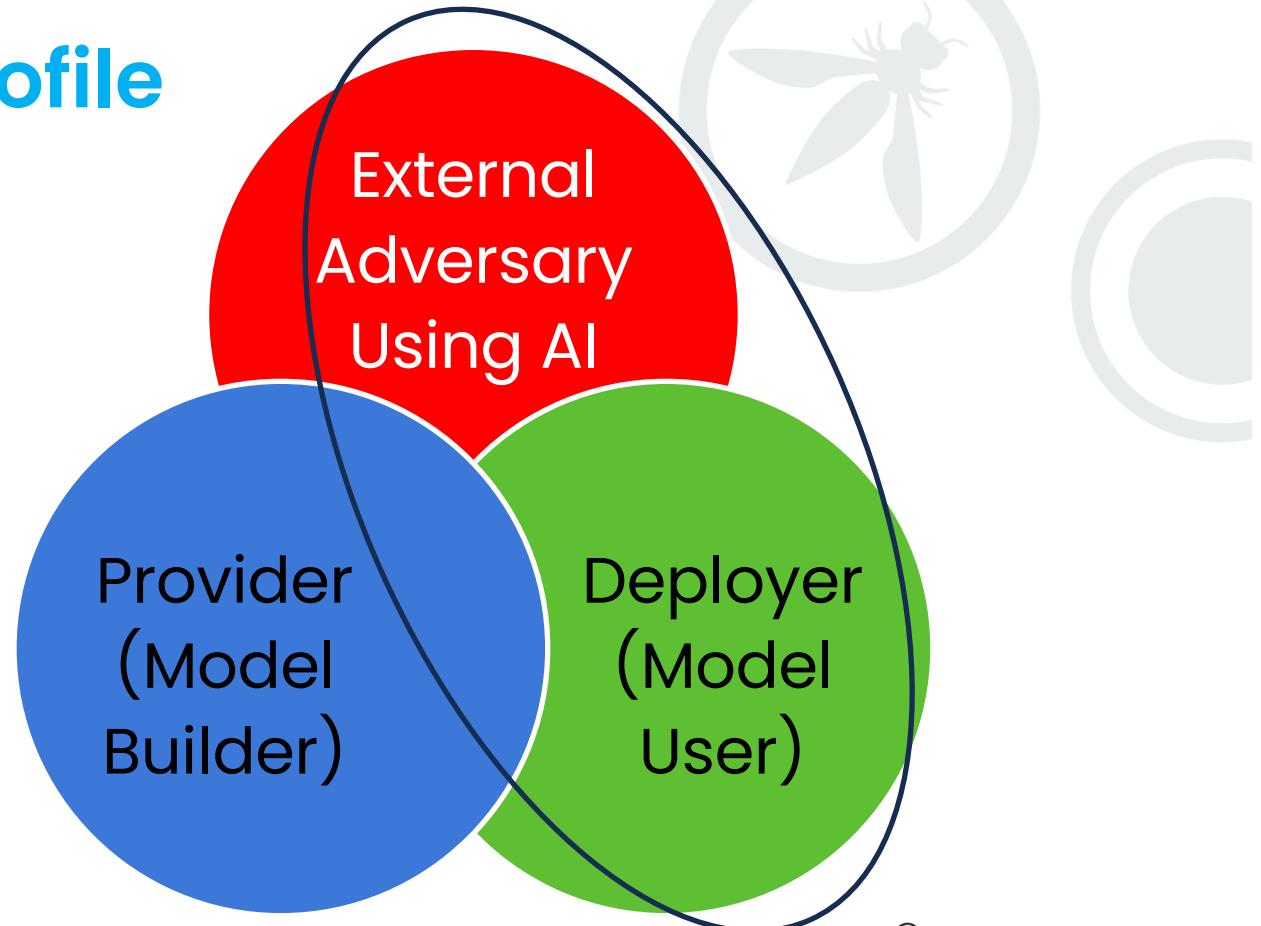


OWASP GenAI Project

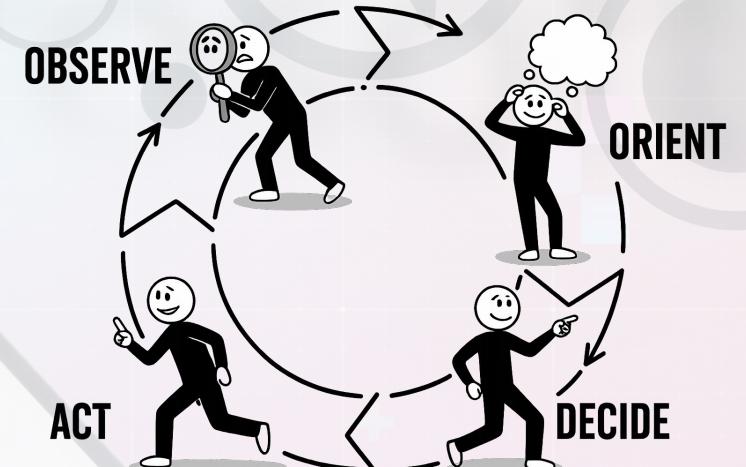
AI Threat Defense COMPASS v1.1



Evaluate by Profile

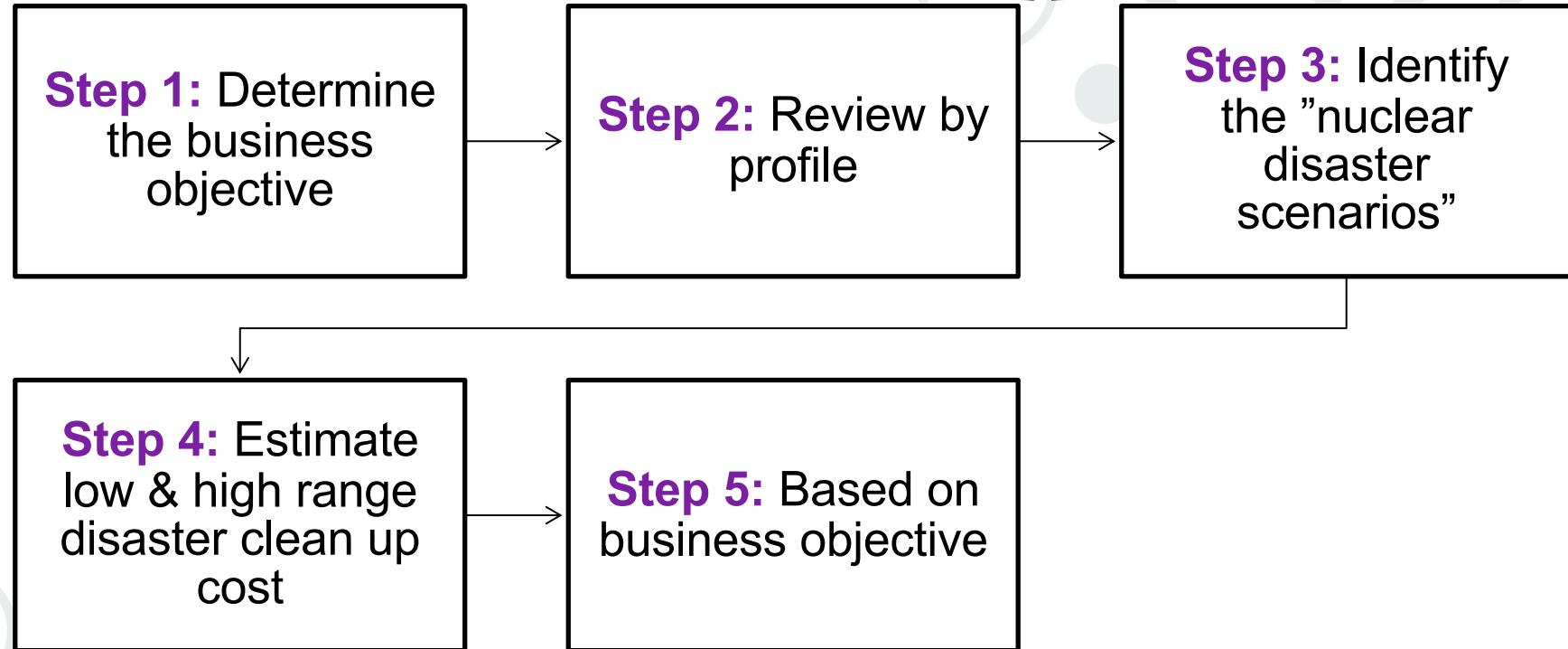
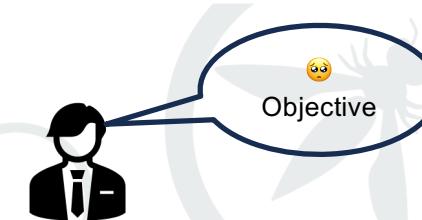


Use Case: Enterprise Copilot Rollout





Or



Data Capture

Profile 2B: Microsoft Enterprise Copilot or Google Enterprise Gemini

These risks apply to Microsoft Copilot, Google Gemini, or similar generative assistants integrated into enterprise productivity suites (What keeps me awake about Microsoft Copilot or Google Gemini for Workspaces)

Note: Deploying these solutions can unintentionally reveal existing security weaknesses by making it easier for users to find and share information they shouldn't. If users have excessive permissions, advanced search capabilities could expose sensitive data and increase the risk of it being shared improperly.

Access & Permissions Risk

1. Overprivileged Access Exposure

- Sensitive Information leakage due to overprivileged access
 - Copilot can query data users have access to but may not need. If least privilege isn't enforced, sensitive information may surface via Copilot-assisted search.
- Advanced search magnifies privilege abuse
 - Hidden files, stale sites, and redacted documents can be surfaced unintentionally due to the model's inference capabilities.
- Role-Based Access Controls (RBAC) not fine-tuned
 - Copilot relies on existing RBAC settings. If RBAC is misconfigured, Copilot becomes a vehicle for policy bypass.

2. Service Account Management

- All service accounts not reviewed, hardened, or audited
- Copilot-managed bot or AI bot operates with persistent high-level permissions
- Non-human identity governance is missing or incomplete

3. Misconfigured Sharing & Collaboration

- Improper Teams sharing (chats, files, meeting notes)
- SharePoint Online sites exposing documents to too broad an audience
- Lack of governance over shared drives or shared mailboxes accessible by Copilot

Data Governance & Classification Risks

1. Inaccurate Data Classification

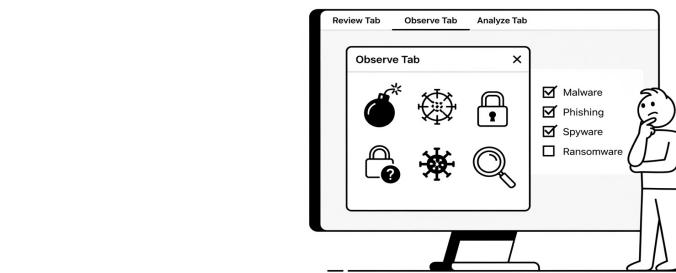
- Copilot indexes unclassified or inconsistently labeled content, increasing risk of inappropriate recommendations or auto-completions.
- No formal classification of sensitivity (e.g., public, internal, confidential, restricted) leads to flattened risk validity.

2. Sensitivity Labeling Gaps

- Sensitivity labels not implemented or not enforced across apps
- Label inheritance across files, chats, and calendar entries is inconsistent
- Lack of visual cues or training for users on what labels mean or how they apply in Copilot/Gemini interactions

3. Retention & Compliance Risks

- Data surfaced by Copilot may violate retention or legal hold policies
- AI assistants may summarize or reproduce content outside of protected systems, undermining compliance



Step 1: Revise the Organizational Impact Low Range and High Range values to align with your organization's impact ratings for catastrophic, severe, major, moderate, and minor ratings. (Low Range 0-25, High Range 152-157). Example #8 provided.

Step 2: Specify 3 or more "Value AI Threat" Scenarios (rows 49-50)

Step 3: Use the example Threat Category / Attack Vectors or modify the table from the Profile 1 and Profile 2 shading, and/or use from the list in the 2a: Observe: Objective Threat Profile.

Heat Map					Defense Maturity Rating Reference (Knowledge Information Confidence)				
5	10	15	20	25	5	6	9	12	15
4	8	12	16	20	4	7	10	13	17
3	6	9	12	15	3	6	9	12	15
2	4	6	8	10	2	4	6	8	10
1	2	3	4	5	1	2	3	4	5
5	6	9	12	15	5	6	9	12	15
4	7	10	13	17	4	7	10	13	17
3	6	9	12	15	3	6	9	12	15
2	4	6	8	10	2	4	6	8	10
1	2	3	4	5	1	2	3	4	5

Threat Category / Attack Vector	Description	Threat - Risk Level
Deep Fakes: voice or image cloning.	Synthetic media where AI is used to create realistic fake content.	3
Adversary Attack w/ AI: Identity / Anonymity	Adversaries using AI Tools to execute attacks to an organization.	3
LLM001025 Sensitive Information Disclosure	Leak of company confidential data	3
LLM001025 Prompt Injection	User maliciously alters prompt input	3
OSINT Gathering	Increased ability to find sensitive data on executives and key employees	4
LLM001025 Supply Chain	Compromising third-party pre-trained models, libraries, or plattforms used in the AI lifecycle.	4
Model Hallucinations	Models hallucinate or fabricate data, leading to poor decisions in critical contexts.	4
LLM001025 Improper Output Handling	Insufficient validation, sanitization, and handling of the outputs	4
LLM001025 Executive Agency	Vulnerability that enables damaging actions to be performed in response to unexpected, ambiguous or manipulated outputs from AI LLM	3
Regulatory or Legal Threat	Violation due to data protection or AI Laws	3
LLM001025 System Prompt Leakage	Disclosing system prompt information that should not be public	3
LLM001025 Vector and Embedding Wreaks Havoc	Weaknesses in how vectors and embeddings are generated, stored, or used.	2
TB: Reputation & Usability	Weakness in how vectors and embeddings are generated, stored, or used.	2
LLM001025 Unintended Consumption	Resource exploitation and unauthorized usage.	2

This score is the average of Impact & Likelihood

Organizational Impact					
Impact level	Rating	AI Specific Example	Low Range	High Range	
Catastrophic	5	Major problems from which there is no recovery or significant damage which has high financial cost, and impacts ability to meet overall business objectives. Complete loss of ability to deliver a critical program.	\$5,000,000.00	\$10,000,000.00	
	4	Incident that requires a major action to support mitigating how service is provided. Significant has a long recovery period. Failure to meet service delivery	\$4,000,000.00	\$1,000,000.00	
Severe	3	Recovery from an incident requires cooperation across organization. May generate media attention.	\$999,999.00	\$100,000.00	
	2	Deal with at a department level but requires Executive notification. Delay in finding or change in funding criteria. Stakeholders or client would take note.	\$99,999.00	\$10,000.00	
Moderate	1	Deal with internally at manager level. No escalation of the issue required.	\$1,000.00	\$1,000.00	

Screenshot

+ = 1 About * 1 FAQ * 2 Observe: Objective Dashboard * 2a Observe: Objective Threat Profile * 2b Observe: Attack Surface Analysis * 3 Orient Summary * 3a Orient: Known AI Vulnerabilities * 3b Orient: Known AI

gh: Orient Known AI Incidents			
Use the Orient Incident tab to research AI Incidents and impact costs if available. Update the existing table of example incidents / impacts costs with objective relevant information by using the links to reports, incident databases, legal cases, and regulatory information ** Scroll to the bottom of the sheet for links to AI Incident Databases			
Known AI Incidents			
Incident	Vulnerability	Impact	Reference
Solana Scam	LL.M01	\$2,500	Link
ShadowRay	LL.M02 / LL.M03	\$1,000,000,000	Link
Chat GPT Inference Attack	LL.M03		Link
Google Map Deaths	LL.M09		Link
Foxter welfare fraud detection algorithm accused of exacerbating inequality	Bias		Link
Deep-Fake Fraud	LL.M02	\$25,000,000	Link
McDonald sued for use of AI which collected voice print biometrics	LL.M02	(Dismissed)	Link
Equal Employment Opportunity Commission v. iFutureGroup, Inc.	Bias / Discrimination	\$350,000	Link
Mobley v. Workday, Inc.	Bias / Discrimination	Still in the system but could have major impact on using Workday for hiring	Link
Meta capturing facial data \$1.4B Texas Settlement	LL.M02	\$1,000,000,000	Link
SoundCloud discreetly changed its terms of service, adding a clause that may interpreted as giving the company the right to use users' music and audio uploads to train AI models - including generative AI capable of replicating or synthesizing artists' voices, music, or likenesses.	LL.M02		Link
The New York City government's "MyCity" chatbot, launched as a pilot program in October 2023, was designed to provide business owners with information from over 2,000 NYC Business web pages and articles.	LL.M09	The New York City government spent over \$600,000 on the development and initial six months of operation for the MyCity chatbot, which launched as a pilot program in October 2023.	Link
In early 2024, T-Mobile revealed that hackers used an AI-equipped application programming interface (API) to gain unauthorized access to sensitive customer information, including full names, contact numbers, and PINs of its customers.	LL.M02 / LL.M03	\$31.5 million settlement with the FCC in 2024, requiring the implementation of enhanced security measures such as phishing-resistant multifactor authentication and regular third-party security audits.	Link
Air Canada Chatbot customer who was misled into paying for full-price flight tickets by a contact center chatbot.	LL.M09	\$12.62 refund	Link
Servicenow, a provider of agric artificial intelligence-based IT management and workflow software, Agentic AI Tech Firm Says Health Data Leak Affects 483,000. certain information within its Catholic Health Elasticsearch database was inadvertently made publicly available.	LL.M07	Expected to be in the millions	Link Link Link
Clearview AI, a U.S.-based facial recognition company, has faced significant global scrutiny and legal action for scraping billions of images from the internet and social media platforms without user consent. These images were used to build a massive facial recognition database, which Clearview AI marketed primarily to law			

Adversary use of AI Reports
OpenAI Influence and cyber operations updates
Google Cloud Security Resources Hub
Detecting and Countering Malicious Uses of Claude: March 2025
MISP Galaxy MITRE ATLAS Attack Pattern
AI Incident Data Bases
MITRE Atlas
AIIAAC Repository
OECD AI Incidents Monitor (AIM)
AI Incident Database
AI Risk Repository
RealHarm Dataset
Language Model Security Database
Legal & Regulatory
George Washington University AI Litigation database
Mischon de Reya Gen AI IP case tracker
AI Copyright Lawsuits Edward Lee
IAPP Global AI Law and Policy Tracker
IAPP State AI Governance Legislation Tracker
Fairly AI Map of Global Regulation

3b: Orient: Known AI Incidents

6a Reference: Defenses & Mitigations			
Use this list to map preventative mechanisms and detective controls to LLM threats and vulnerabilities.			
Key Controls / Mitigation Strategies		Key Detective Mechanisms	
Noise data sanitizing, provenance, data validation/annotation, outlier detection in training, source control on data pipelines, model behavior testing (Benchmarking)		Monitoring training data statistics, benchmarking model performance directly model retraining and validation.	
Input validation/limitation, adversarial training, defense distillation, gradient masking (lens), limiting model output confidence scores, rate-limiting inputs.		Monitoring input patterns for anomalies, analyzing output confidence scores, detecting unusual prediction failures.	
API rate limiting, query monitoring/filtering, watermarking model outputs, differential privacy (during training), access controls, obfuscation techniques.		Analyzing query patterns/frequency, monitoring for automated/predatory activity, detecting similar model outputs externally.	
Differential privacy, data anonymization, pseudonymization, reducing model complexity, severe multi-party computation (MPC), access controls on model outputs.		Difficult to detect directly, focus on preventative controls and data governance audits.	
String input sanitization/parsing, robust signature (user vs. attacker prompts), output filtering/validation, least privilege for executed tasks, user awareness.		Monitoring prompt/response patterns, anomaly detection in API calls made by AI, output content analysis, user feedback logs.	
Vendor risk assessment, Software Bill of Materials (SBOM) for AI, code/model scanning, integrity checks, using trusted repositories/sources, secure build pipelines.		Monitoring component integrity, vulnerability scanning of dependencies, behavioral analysis of AI components.	
Input rate limiting, resource quota, input validation/complexity checker, efficient model architecture, infrastructure scaling, resilience.		Monitoring resource utilization (CPU, GPU, Memory), tracking request rates, latency, anomaly detection in traffic.	
Acceptable Use Policies (AUPs), output filtering, monitoring toxicity, harmful content, watermarking outputs, user identity verification, rate-limiting generation.		Monitoring generated content for policy violations, detecting anomalous usage patterns, extraction/intent intelligence.	
Explainability features (XAI), confidence scoring, human-in-the-loop workflows, clear UI/UX design indicating AI role, user training on AI limitations, regular audits.		Monitoring downstream impact of AI decisions, tracking overrides/interventions by humans, feedback collection.	
OWASP Top Ten for LLM 2023			
Vulnerability	Example	Prevention Controls	Detection Controls
L1.01-2023 Prompt Injection	Direct/Indirect injection, hidden prompt in image, code injection, multi-stage attack.	Constrain model behavior via system prompts and role definitions. Validate user input with allow-lists or whitelists. Monitor least privilege on external tool access.	Monitor prompts for injection patterns. Use anomaly detection on prompt sequences. Red team adversarial prompts and log fail-breaking input/output.
L1.02-2023 Sensitive Information Disclosure	PII leakage, proprietary algorithm exposure, unmasked training data inclusion.	Monitor training data to remove sensitive info. Enforce output filtering for PII and restrict LLM data access. Apply policy-based response boards.	Scan outputs for PII. Monitor queries and rate-limit odd patterns. Keep query/response audit logs. Solicit user disclosures.
L1.03-2023 Supply Chain	Malicious LLM adapters, exploited models, compromised third-party sources.	Verify and sign third-party models. Use secure update pipelines and isolated deployment. Treat ML artifacts as critical software components.	Continuously audit model behavior. Validate model hashes. Use threat and feeds analysis framework for supply-chain compromise.
L1.04-2023 Data and Model Poisoning	Backdoored datasets, poisoning via prompt input, trigger-based behavior change.	Use curated datasets. Validate external datapoints. Employ adversarial training and defense fine-tuning techniques.	Detect data anomalies during training. Keep training logs. Use sanity queries to test model changes.
L1.05-2023 Unexpected Output Handling	Unescaped backslash, SQL injection via ELM_remote code execution.	Sanitize LLM output before execution. Enforce output-cleaning and secure coding practices. Apply input sanitization on output executing systems.	Log and monitor LLM driven actions. Detect anomalous output types. Monitor execution and simulate prompt based code injections.
L1.06-2023 Executive Agency	L1.06 gives excessive permissions, enabling unintended actions via agents.	Restrict agent tasks. Enforce least privilege and tool use mediation. Require human confirmation for critical actions.	Log agent actions and apply rate limiting. Monitor for unauthorized tool use or privilege abuse. Audit agent decisions vs user actions.
L1.07-2023 System/Prompt Leakage	Leaked prompts containing API keys, internal rules, permissions.	Keep sensitive data out of prompts. Separate business logic from LLM instructions. Use prompt hygiene practices.	Scan outputs for prompt pattern leakage. Log system prompt metadata and user attempts to reveal prompts.
L1.08-2023 Vector and Embedding Fingerprinting	Embedding inversion, poisoned RAG data, cross-task leakage.	Apply source control on vector DB. Validate embedding data. Authenticate data sources and separate trusted/untrusted sources.	Log vector queries and flag unusual scores. Audit embeddings for poisoning. Validate retrieved content and query alignment.
L1.09-2023 Misinformation	Generated fake claims, hallucinated citations, bias manipulation.	Use IASIS with verified sources. Train for lower hallucination. Cross-validate output or require human review.	Instrument fact-checking of outputs. Monitor correction feedback. Use human-in-the-loop for sensitive content. Sample logic misinformation.
L1.10-2023 Unbounded Consumption	Denial of wallet, resource exhaustion, API rate abuse.	Set quota on usage. Limit input/output sizes. Use resource budgeting and sandboxing/tool access.	Monitor system usage. Log request patterns. Set thresholds for usage alerts. Apply model theft heuristics and graceful degradation.
OWASP Agents Top 15			
Vulnerability	Example	Prevention Controls	Detection Controls
E1: Memory Poisoning	Manipulating short/long-term memory to change AI behavior or extract sensitive data.	Validate and constrain agent memory. Enforce session isolation. Assign enforcement and least privileged memory areas.	Detect anomalies in memory usage. Use recorded memory snapshots. Monitor session activity for poisoning attempts.
E2: Tool Abuse	Execution prompts lead all agents to misuse tools like small AI APIs (e.g., agent hijacking).	Restrict toolchains via allow-lists. Validate tool invocation parameters. Define operational boundaries for tools and APIs.	Monitor and log tooluse. Alert on unexpected tool activity. Conduct post-incident analysis of tool misuse.
E3: Privilege Escalation	Dynamic role inheritance or misconfiguration lets attackers escalate privileges.	Apply granular permissions. Enforce dynamic privilege checks. Disable unnecessary cross-agent privilege delegation.	Audit all role and privilege changes. Log/high-priority agent actions. Use behavioral analytics to detect escalation.
E4: Resource Overload	DoS via task overload, memory cascade failures, API quota exhaustion.	Set CPU/memory/task limits. Use adaptive throttling. Rate-limit high-cost operations to prevent agent-based DoS.	Continuously monitor agent resource usage. Alert on spikes. Detect loops and abnormal resource consumption.
E5: Cascading Hallucination Attacks	AI hallucinations spread and reinforce errors through memory and multi-agent interactions.	Sanitize agent outputs before use by others. Require validation of outputs for multi-agent workflows. Limit high-risk auto-instruction behavior.	Trace output propagation across agents. Use ground-truth for validation. Log inter-agent exchanges for backtracking.
E6: Dataset Breaking & Goal Manipulation	Changing AI goals via direct/indirect prompt injection or detection traps.	Review planned actions for goal alignment. Use safeguards for self-correction and enhance user intent confirmation for high-level designs.	Log goal changes over time. Use supervisory model to detect deviations from intent. Investigate shifts in agent reasoning.
E7: Misaligned & Disruptive Behavior	Agents evade constraints to achieve goals disruptively (e.g., lying, illicit actions).	Train agents to report harmful input. Enforce policy-based filters. Require human confirmation for risky actions.	Log and compare agent outputs vs. actions. Detect deception using behavior analysis or secondary truth-checks.
E8: Reputation & Vulnerability	Insufficient logging makes agent behavior untraceable or unaccountable.	Log critical agent actions with cryptographic integrity. Enrich logs with traceable metadata. Prevent unallowable behavior.	Alert on missing or malformed logs. Monitor agent behavior for unlogged activity. Use immutable audit trails.

NAVIGATING THE AI FRONTIER



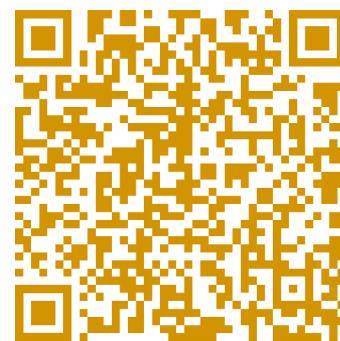
Questions?

**COMPASS
Workbook**



#team-genai-secgov

**Recommended
Resources**



Sandy



Sabrina



Connect with Us!

