# What is Threat-Informed Defense?

*"The systematic application of a deep understanding of adversary tradecraft and technology to improve defenses."*

*https://ctid.mitre.org/our-mission/*

MITRE | Center for Threat Informed Defense

# AI Intersections with Cyber

## Security & Assurance of AI-enabled Systems

Securing the unique system vulnerabilities of AI-enabled systems – includes red teaming to discover vulnerabilities

## Using AI in Offensive or Malicious Cyber Attacks

Attackers using AI in offensive assaults on both traditional cyber systems and AI-enabled systems

## Using AI in Cybersecurity Practices

Using AI to improve our cybersecurity practices, i.e., detection, risk analysis, and defensive or mitigation techniques

**MITRE has capabilities and teams working in all three areas Today we are focusing on the security of AI-enabled systems**

# Focusing on Real-World Demonstrated AI Attack Vectors

## Incident - Malicious Attack by an Adversary

Ongoing real-world AI supply chain attack vector with estimated financial impact over $1 Billion
(as of March 2024, that $ estimate would likely be much higher now)

**ShadowRay:** The lack of authorization in the Ray Jobs API default configuration allows adversaries to invoke arbitrary jobs on the API. This grants access to user tokens/PII and fraudulent use of cloud compute time at the cost of the user.

CVE-2023-48022



Image credit: New ShadowRay Campaign Targets Ray AI Framework in Global Attack (hackread.com)

ShadowRay: First Known Attack Campaign Targeting AI Workloads Exploited In The Wild | Oligo Security

**Ray Jobs API**

**User PII**
User Passwords
Slack Tokens
Slack Messages
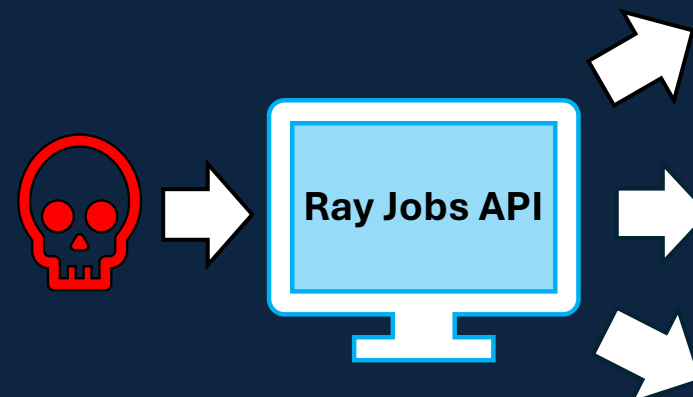Private SSH Keys

**Ray Leaks**
AI Production Workloads
Production DB Credentials
Full Ray Database Access

**External Services**
OpenAI Tokens
KubernetesAPI Access
Stripe Tokens
HuggingFace Tokens

**MITRE** | Center for Threat
Informed Defense™

# Focusing on Real-World Demonstrated AI Attack Vectors

## Exercise - Red Teaming/Real World Demonstration

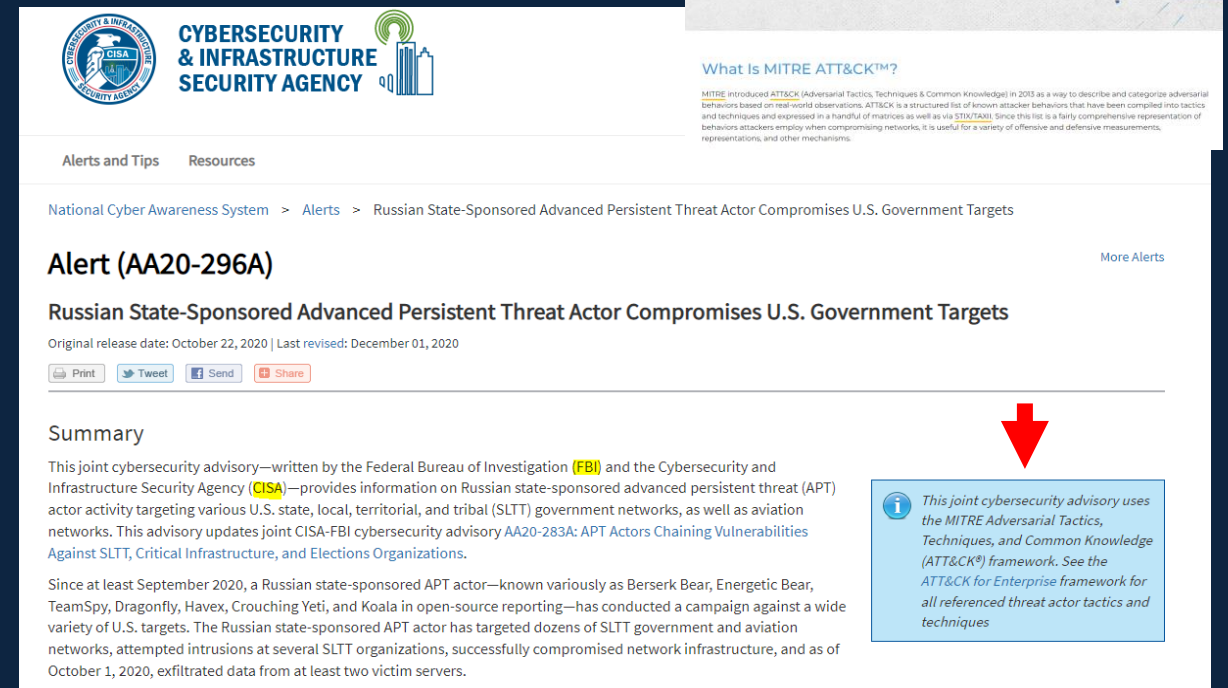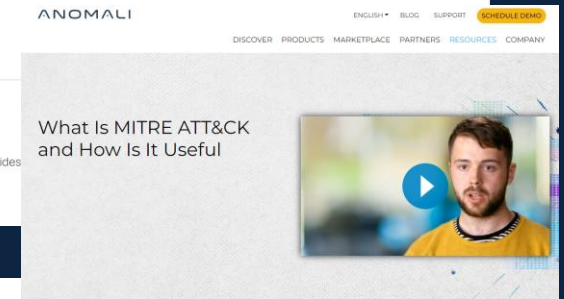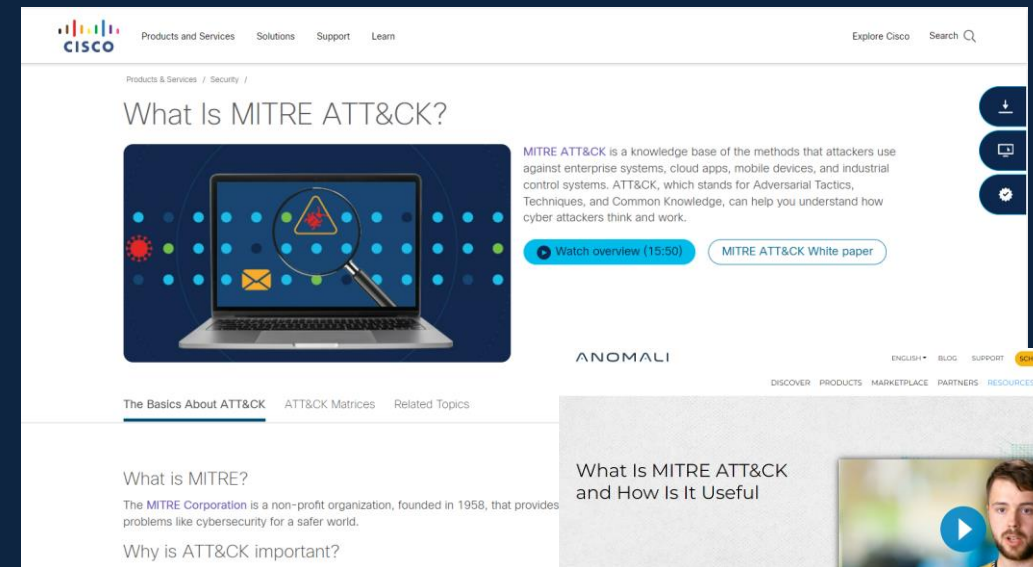Designed to attack the GenAI ecosystem and propagate without user interaction

**Morris II Worm:** Injects the prompt without user interaction via the RAG email context collection and delivers a payload of the adversary's choosing (in this case, leaking PII). The worm replicates the adversarial prompt in email auto-replies and propagates via other RAG-Enabled email databases.



RAG- Enabled Replication

Server receives external email → Generator (Language Model) → Malicious Auto-Reply

Retriever

Malicious Email sent to RAG-Enabled Server → Email Server → Adversarial Self-Replicating Prompt is injected into the generator by RAG Context

Retrieved

# Path to Actionable Impact

## Demonstrated Capability: MITRE ATT&CK

- Provides **common language** for cybersecurity professionals (STIX) to document common tactics, techniques, and procedures of advanced persistent threats

- **Fully adopted and promoted by CISA** and used across government agencies such as the FBI in advisories about threat activity

The ATLAS Matrix below shows the progression of tactics used in attacks as columns from left to right, with ML techniques belonging to each tactic below. & indicates an adaption from ATT&CK. Click on the blue links to learn more about each item, or search and view ATLAS tactics and techniques using the links at the top navigation bar. View the ATLAS matrix highlighted alongside ATT&CK Enterprise techniques on the ATLAS Navigator.

| Reconnaissance & | Resource Development & | Initial Access & | ML Model Access | Execution & | Persistence & | Privilege Escalation & | Defense Evasion & | Credential Access & | Discovery & | Collection & | ML Attack Staging | Exfiltration & | Impact & |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 techniques | 9 techniques | 6 techniques | 4 techniques | 3 techniques | 4 techniques | 3 techniques | 3 techniques | 1 technique | 6 techniques | 3 techniques | 4 techniques | 4 techniques | 7 techniques |
| Search for Victim's Publicly Available Research Materials | Acquire Public ML Artifacts | ML Supply Chain Compromise | AI Model Inference API Access | User Execution & | Poison Training Data | LLM Prompt Injection | Evade ML Model | Unsecured Credentials & | Discover ML Model Ontology | ML Artifact Collection | Create Proxy ML Model | Exfiltration via ML Inference API | Evade ML Model |
| Search for Publicly Available Adversarial Vulnerability Analysis | Obtain Capabilities & | Valid Accounts & | ML-Enabled Product or Service | Command and Scripting Interpreter & | Backdoor ML Model | LLM Plugin Compromise | LLM Prompt Injection | | Discover ML Model Family | Data from Information Repositories & | Backdoor ML Model | Exfiltration via Cyber Means | Denial of ML Service |
| Search Victim-Owned Websites | Develop Capabilities & | Evade ML Model | Physical Environment Access | LLM Plugin Compromise | LLM Prompt Injection | LLM Jailbreak | LLM Jailbreak | | Discover ML Artifacts | Data from Local System & | Verify Attack | LLM Meta Prompt Extraction | Spamming ML System with Chaff Data |
| Search Application Repositories | Acquire Infrastructure | Exploit Public-Facing Application & | Full ML Model Access | | LLM Prompt Self-Replication | | | | Discover ML Model Artifacts | | Craft Adversarial Data | LLM Data Leakage | Erode ML Model Integrity |
| Active Scanning & | Publish Poisoned Datasets | LLM Prompt Injection | | | | | | | LLM Meta Prompt Extraction | | | | Cost Harvesting |
| | Poison Training Data | Phishing & | | | | | | | Discover LLM Hallucinations | | | | External Harms |
| | Establish Accounts & | | | | | | | | Discover AI Model Outputs | | | | Erode Dataset Integrity |
| | Publish Poisoned Models | | | | | | | | | | | | |
| | Publish Hallucinated Entities | | | | | | | | | | | | |

atlas.mitre.org

**150+ organizations** are engaged in ATLAS, **using ATLAS tools and capabilities** to understand and mitigate AI security risks

# ATLAS Case Study:
# Camera Hijack Attack on Facial Recognition System

Two individuals in China attacked an ML-enabled face identification system to gain access to the local government's tax system. They created a fake shell company and sent invoices via the tax system to supposed clients and **steal $77 million from the Shanghai Tax Authority**.
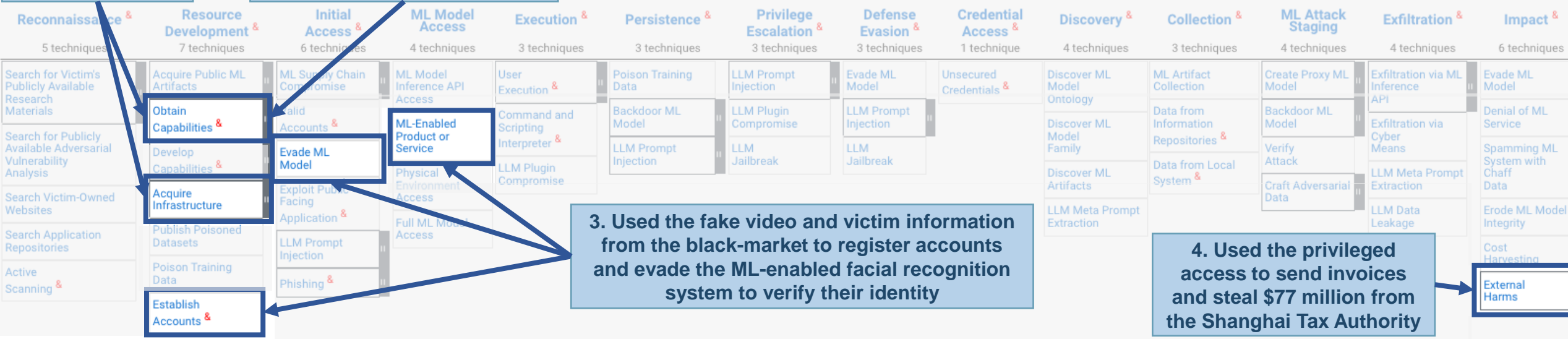


**1. Customized a cheap cell phone to display a fake video feed**

**2. Created crude videos of a head turning, opening/closing its eyes/mouth using black-market photos of users**

**3. Used the fake video and victim information from the black-market to register accounts and evade the ML-enabled facial recognition system to verify their identity**

**4. Used the privileged access to send invoices and steal $77 million from the Shanghai Tax Authority**

| Reconnaissance & | Resource Development & | Initial Access & | ML Model Access & | Execution & | Persistence & | Privilege Escalation & | Defense Evasion & | Credential Access & | Discovery & | Collection & | ML Attack Staging & | Exfiltration & | Impact & |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 techniques | 7 techniques | 6 techniques | 4 techniques | 3 techniques | 3 techniques | 3 techniques | 3 techniques | 1 technique | 4 techniques | 3 techniques | 4 techniques | 4 techniques | 6 techniques |
| Search for Victim's Publicly Available Research Materials | Acquire Public ML Artifacts | ML Supply Chain Compromise | ML Model Inference API Access | User Execution & | Poison Training Data | LLM Prompt Injection | Evade ML Model | Unsecured Credentials & | Discover ML Model Ontology | ML Artifact Collection | Create Proxy ML Model | Exfiltration via ML Inference API | Evade ML Model |
| Search for Publicly Available Adversarial Vulnerability Analysis | Obtain Capabilities & | Valid Accounts & | ML-Enabled Product or Service | Command and Scripting Interpreter & | Backdoor ML Model | LLM Plugin Compromise | LLM Prompt Injection | | Discover ML Model Family | Data from Information Repositories & | Backdoor ML Model | Exfiltration via Cyber Means | Denial of ML Service |
| Search Victim-Owned Websites | Develop Capabilities & | Exploit Public Facing Application | Physical Environment Access | LLM Plugin Compromise | LLM Prompt Injection | LLM Jailbreak | LLM Jailbreak | | Discover ML Artifacts | Data from Local System & | Verify Attack | LLM Meta Prompt Extraction | Spamming ML System with Chaff Data |
| Search Application Repositories | Acquire Infrastructure & | LLM Prompt Injection | Full ML Model Access | | | | | | LLM Meta Prompt Extraction | Craft Adversarial Data | LLM Data Leakage | Erode ML Model Integrity |
| Active Scanning & | Publish Poisoned Datasets | Phishing & | | | | | | | | | | Cost Harvesting |
| | Poison Training Data | | | | | | | | | | | | |
| | Establish Accounts & | | | | | | | | | | | External Harms |

MITRE | Center for Threat Informed Defense

# Secure AI @ the Center for Threat-Informed Defense

# Secure AI @ the Center for Threat-Informed Defense

## Problem

In addition to traditional cybersecurity vulnerabilities, AI-enabled systems are also susceptible to new attacks based on the unique vulnerabilities of AI-enabled systems.

## Solution

Accelerate the development of MITRE ATLAS to meet industry needs in AI Security, including incident sharing metrics & mechanisms, threats to Generative AI systems, strategies to mitigate threats to AI-enabled systems, tools and playbooks to emulate threats to AI-enabled systems.

## Impact

Secure organizations against the unique emergent attack surfaces that arise in complex systems containing AI.

Industry Leaders Expand Threat-Informed Defense to AI-Enabled Systems

Jon Baker · Follow
Published in MITRE-Engenuity · 3 min read · Jul 16, 2024

Written by *Suneel Sundar*.

MITRE ENGENUITY. | Center for Threat Informed Defense

https://ctid.io/secure-ai

MITRE | Center for Threat Informed Defense

# Secure AI: Core Deliverables

## 1. ATLAS Knowledge Base

Increase the knowledge base and understanding of real-world threats through collection of **incident sharing** metrics and mechanisms.

## 2. Generative AI Threats

Extend the data-driven generative AI focus of MITRE ATLAS by documenting **new case studies & mitigations** that address the vulnerabilities of systems that incorporate generative AI.

## 3. Synchronize Updates to ATLAS & ATT&CK

Align the ATLAS TTPs with the current version of ATT&CK TTPs and implement a plan that may keep the TTP versions in sync.

MITRE | Center for Threat Informed Defense

# Secure AI: ATLAS Matrix Update

## Case Studies

- ChatGPT Package Hallucination
- ShadowRay
- Morris II Worm: RAG-Based Attack
- Web-Scale Data Poisoning: Split-View Attack

Booz | Allen | Hamilton®

intel

FS-ISAC

verizon business

standard chartered

Microsoft

## Techniques

- Discover LLM Hallucinations
- Discover AI Model Outputs
- Erode Dataset Integrity
- Publish Hallucinated Entities
- User Execution: Malicious Package
- Acquire Infrastructure: Domains
- LLM Prompt Self-Replication
- Publish Poisoned ML Model
- Acquire Infrastructure: Physical Countermeasures
- AI Supply Chain Compromise: Hardware
- AI Model Inference API Access

## Mitigations

- Generative AI Guidelines
- Generative AI Model Alignment
- AI Bill of Materials
- AI Telemetry Logging
- Maintain AI Dataset Provenance

# Synchronize Updates to ATLAS & ATT&CK

Align the ATLAS TTPs with the current version of ATT&CK TTPs and implement a plan that may keep the TTP versions in sync.

# AI Incident Sharing

Digital, rapid, anonymized community sharing at **ai-incidents.mitre.org**

- Developed a structured format to collect relevant AI incident information

- Incidents are **shared with MITRE under a data sharing agreement** and stored in a protected database

- Data can be anonymized for **sharing with a trusted community**

- Aggregated data and trends can be visualized in dashboards and **inform AI risk analysis at scale**

MITRE | Center for Threat Informed Defense

# AI Incident Sharing

## Digital, rapid, anonymized community AI Incident Sharing



Notional Dashboard Concept

In capturing and carefully distributing the appropriately sanitized and technically focused AI incident data, this effort aims to enable more data driven risk intelligence and analysis at scale across the community.

**Project Release
October 2, 2024**

**Key Resources**

ctid.io/secure-ai

AI-incidents.mitre.org

atlas.mitre.org

# What's Next?

The Center will continue to build out the knowledge base and incident sharing mechanisms, collect case studies and feedback, and identify additional directions for future Secure AI research.



Implements ATLAS techniques and adversary profiles that can be used stand-alone or with ATT&CK TTPs

An **Open-Source CALDERA Plug-In** to Emulate an Adversary Attacking Your AI-enabled System

https://github.com/mitre-atlas/arsenal

# Ongoing ATLAS Efforts



## Mitigations

**Mitigations can prevent an adversary from the executing techniques in ATLAS**

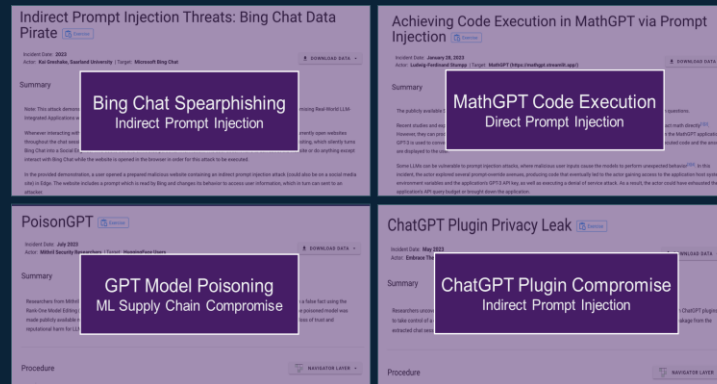**October Update Added Community Inputs** to the 2023 draft release

**Aligned with CRISP-ML Lifecycle Phases**
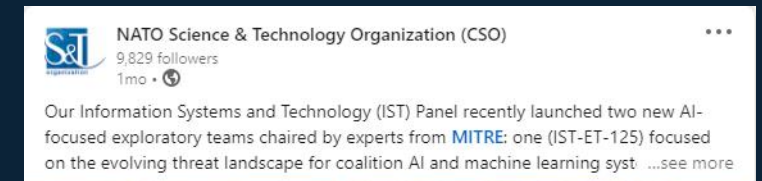
## GenAI Attack Vectors

**Updated in Oct 2024**

**Real World LLM Attack Pathways** Grounded in New Techniques and Case Studies

**Collaboratively Developed with Microsoft, Intel, Verizon, CrowdStrike and 10+ other orgs**

## NATO Task Group

**NATO RTG Launched in May**

- **Leverage coalition AI Security & Assurance capabilities (ATLAS)**
- **Share threat intelligence/vulnerabilities**
- **Shape exemplar shared use cases**
- **Build defensive and mitigation techniques**
- **Develop red teaming capabilities/exercises**

# CVE and CWE

**Engaging with both the Common Vulnerability Enumeration (CVE) and Common Weakness Enumeration (CWE) communities on AI Vulnerabilities & Weaknesses**

- **CVE AI WG** – Working with the CVE Board and AI WG to provide more clarity on how AI security vulnerabilities will fall inside/outside CVE scope via a series of blog posts.

- **CWE AI WG** – Working with the AI Working group on AI-related updates

  - CWE-1426: Improper Validation of Generative AI Output

  - A new demonstrative example for "prompt injection" was added to CWE-77: Improper Neutralization of Special Elements used in a Command ('Command Injection').

  - New observed examples were added to multiple CWEs related to AI/ML and generative AI prompts, including one example of "prompt injection."

https://www.cve.org/Media/News/item/blog/2024/07/09/CVE-and-AIrelated-Vulnerabilities

https://cwe.mitre.org/news/archives/news2024.html#july16_CWE_Version_4.15_Now_Available

MITRE | Center for Threat Informed Defense

# AI Risk Database

**Inspired by VirusTotal, we are building on that vision as we shape a long-term expansion plan.**



**We would love to get your feedback and input on where it should go/what would be most helpful to you!**

MITRE | Center for Threat Informed Defense

# Dioptra – Test Platform for Trustworthy AI

### Dioptra is NIST's software test platform for assessing the trustworthy characteristics of artificial intelligence.

Dioptra **supports the Measure function of NIST's AI Risk Management Framework** by providing functionality to assess, analyze, and track identified AI risks.

# Plans for Dioptra + AI Risk Database

**Dioptra can be used to evaluate AI models and submit reports to the AI Risk Database**

- Help standardize vulnerability reports and metrics

- Allow AIRDB users to verify and validate results

- Promote sharing of experiment templates

- Provide plugins for a variety of evaluations – security and beyond

**Looking for feedback and ideas**

- Dioptra v1.0.0 was released in July of 2024

- https://github.com/usnistgov/dioptra

The Center for Threat-Informed Defense conducts collaborative R&D projects that **improve cyber defense at scale**

+ MITRE

**Membership is:**
- ✓ Highly-sophisticated
- ✓ Global & cross-sector
- ✓ Non-governmental
- ✓ Committed to collaborative R&D in the public interest

https://ctid.io/our-work

**Mission: Advance the state of the art and the state of the practice in threat-informed defense globally.**

# Advance Secure AI

**We would love to have you involved in the next Secure AI project!**

## AI Incident Sharing

- Beta test the incident submission system and submit your incidents/successful red teaming exercises.
- Provide feedback on the kind of information/data your org would want to receive to mitigate incidents.

## AI Risk Database

- Beta test the new vulnerability submission system and provide feedback that will help us shape the next version (including the AI BOM).
- Shape the combination of opensource/public vulnerability detection tools built around the NIST Dioptra tool as our centerpiece

## ATLAS Matrix

- Contribute to the bi-annual major update: send us your mitigations, case studies, and/or feedback on the tactics and techniques in the matrix.