

# AI Red Teaming

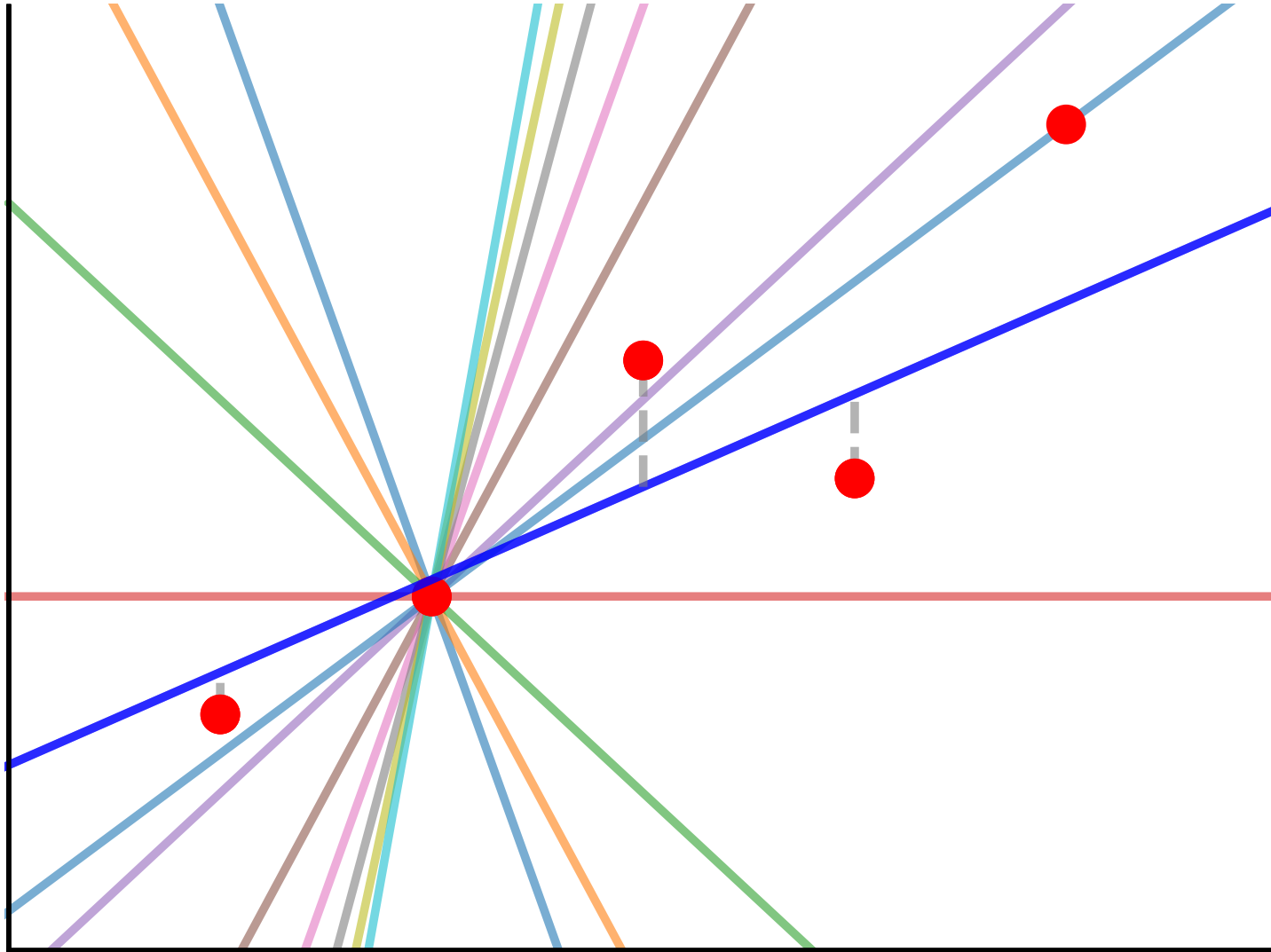
Federico Cerutti

[federico.cerutti@unibs.it](mailto:federico.cerutti@unibs.it)

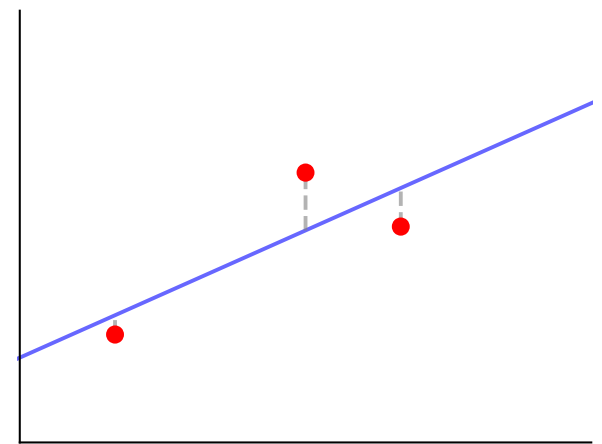
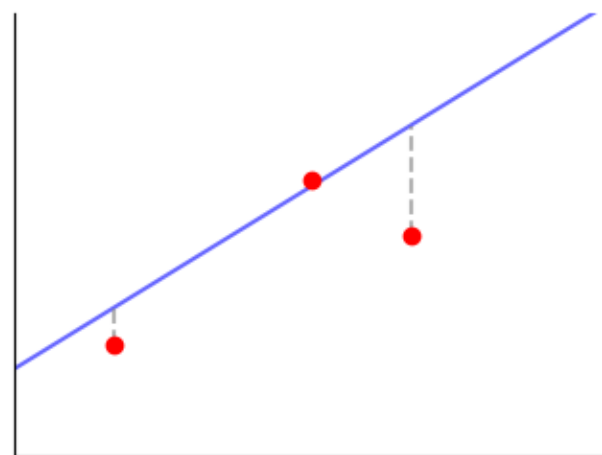
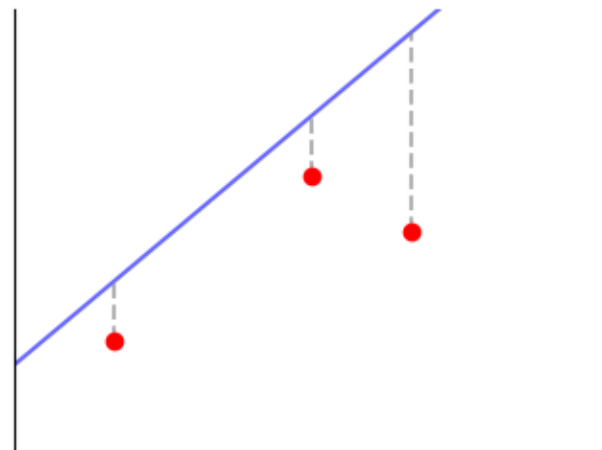
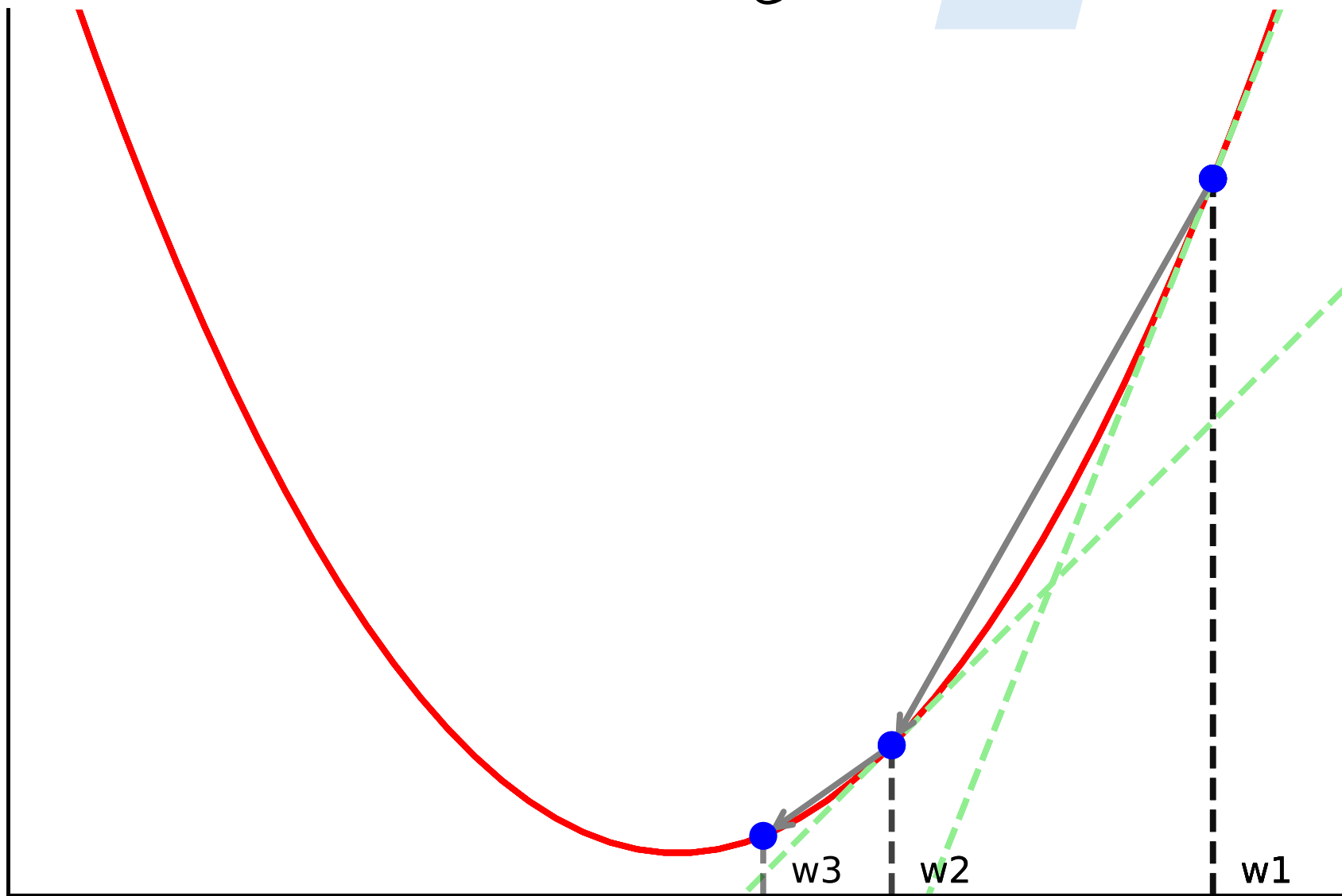
$$\max_{x \in K} \left| f(x) - \sum_{k=1}^N a_k h(w_k x + b_k) \right| \leq \varepsilon$$

$$y = wx + b$$

$$y = wx + b$$

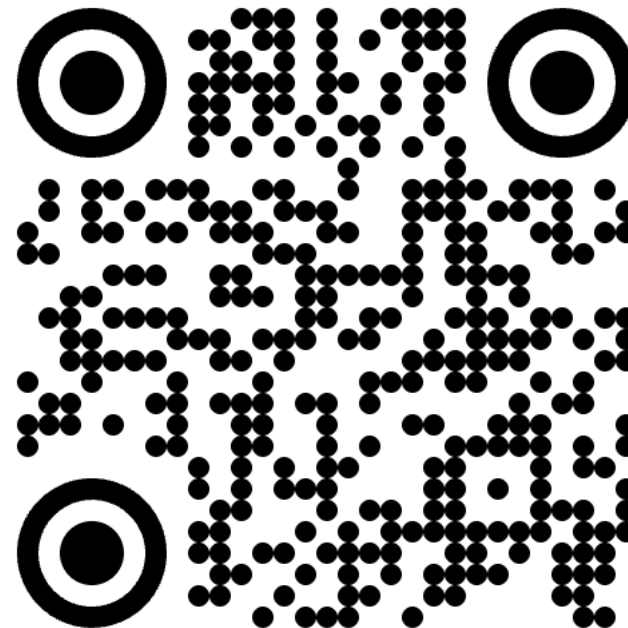


$$y = wx + b$$



**NIST AI 100-1**

**Artificial Intelligence Risk Management  
Framework (AI RMF 1.0)**



<https://www.nist.gov/itl/ai-risk-management-framework>



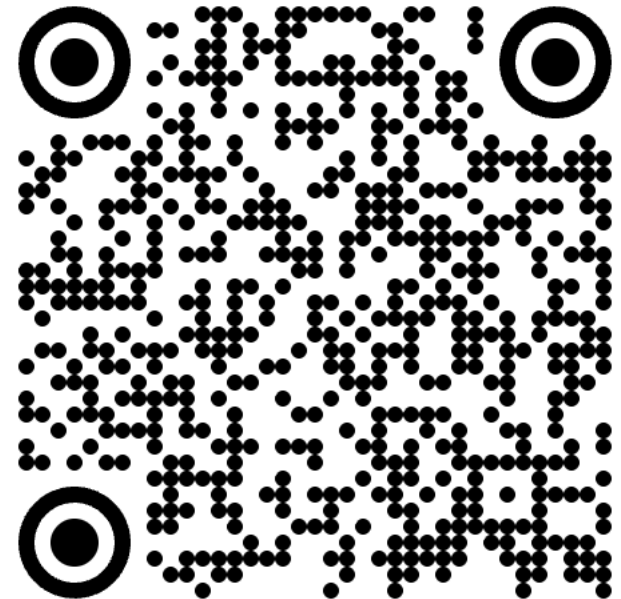
# NIST Trustworthy and Responsible AI

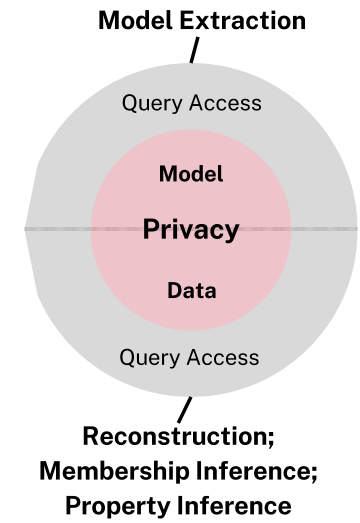
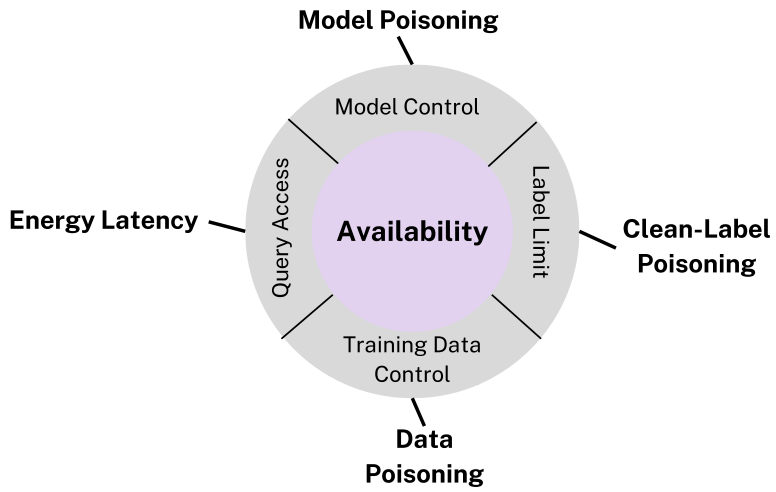
## NIST AI 100-2e2025

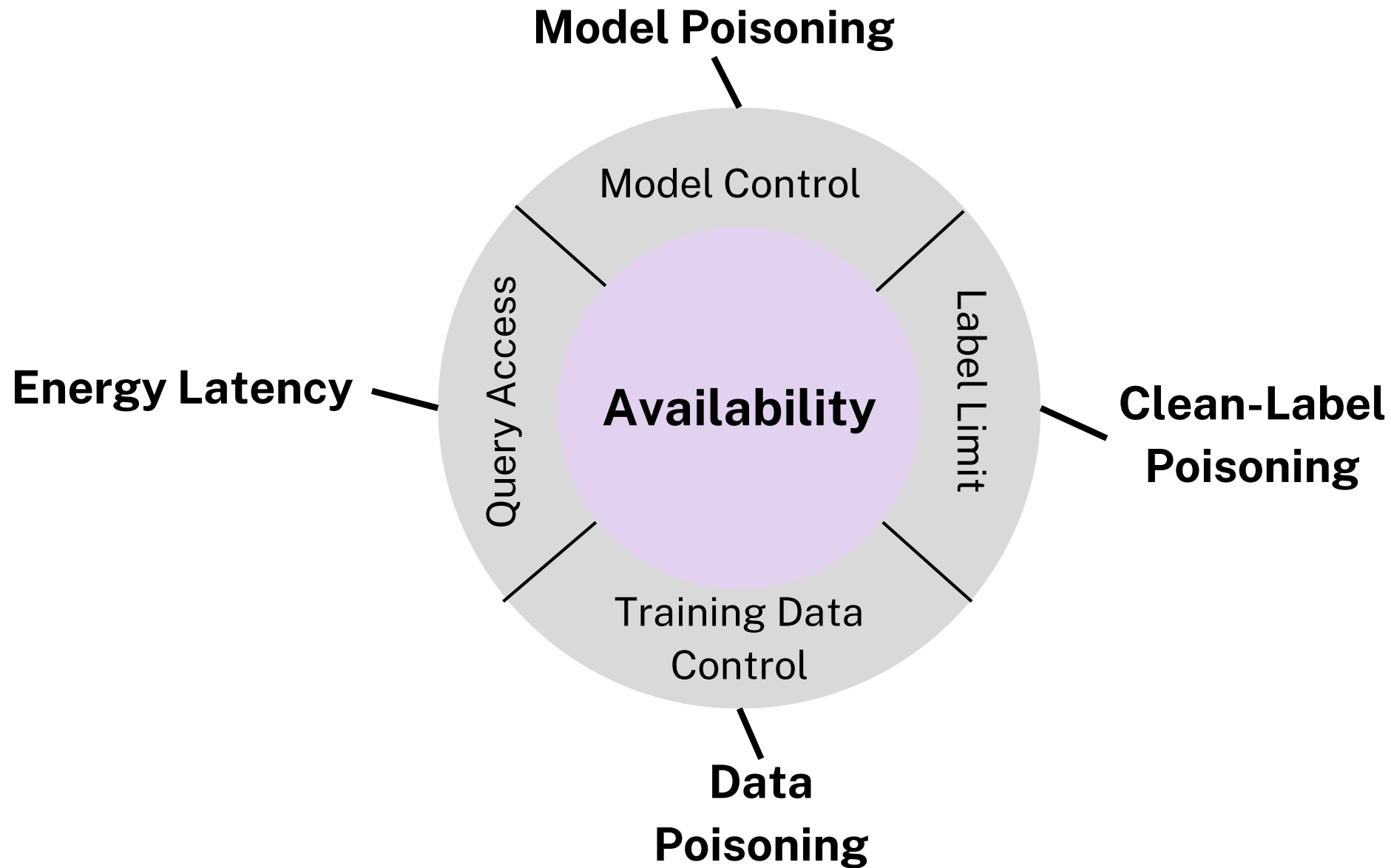
### **Adversarial Machine Learning**

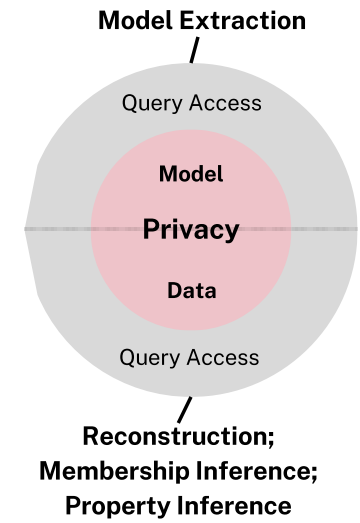
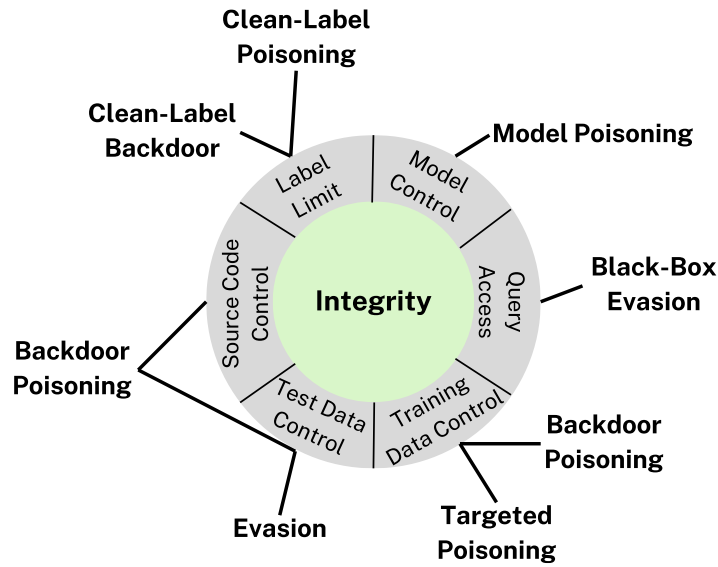
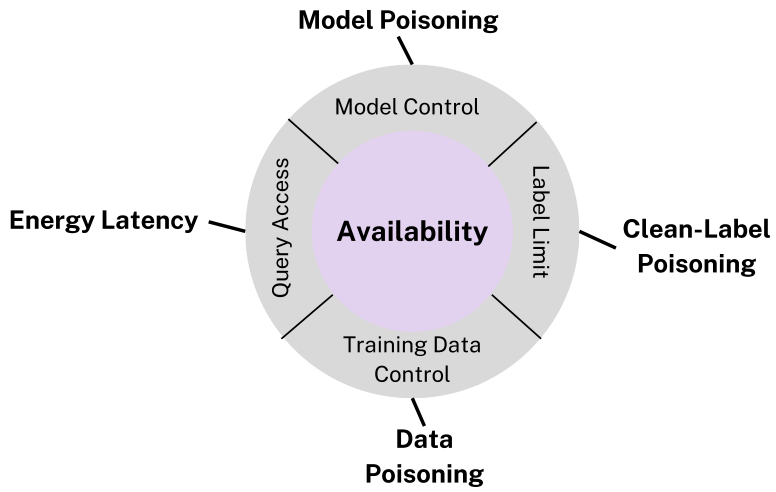
*A Taxonomy and Terminology of Attacks and Mitigations*

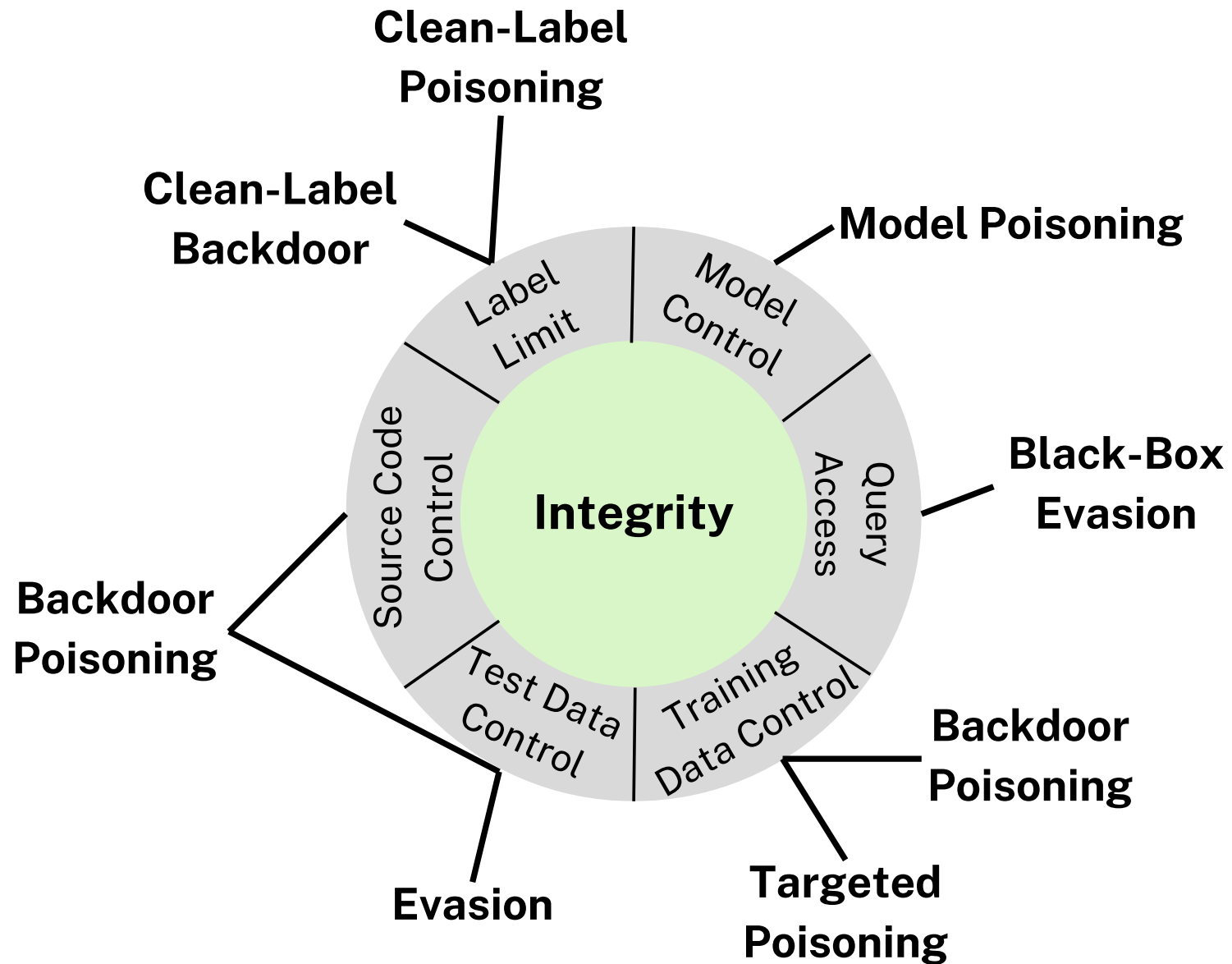
<https://csrc.nist.gov/pubs/ai/100/2/e2025/final>

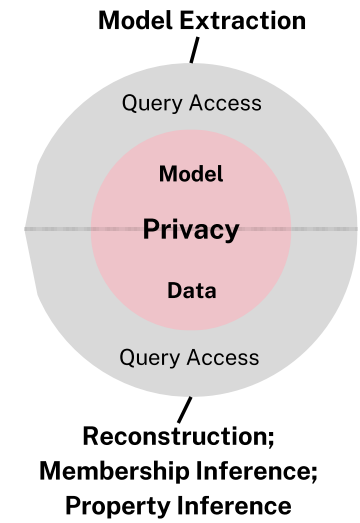
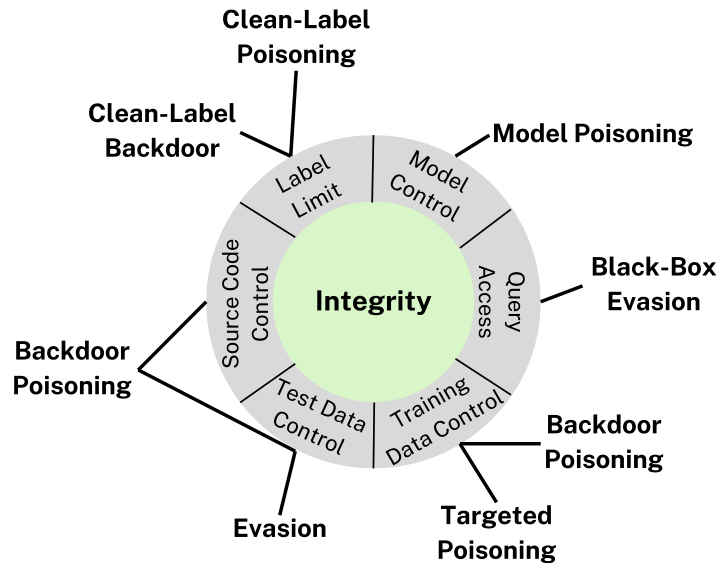
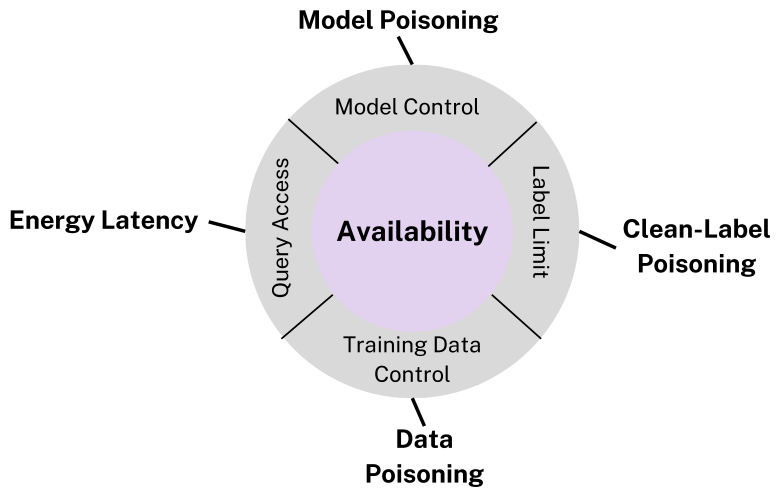




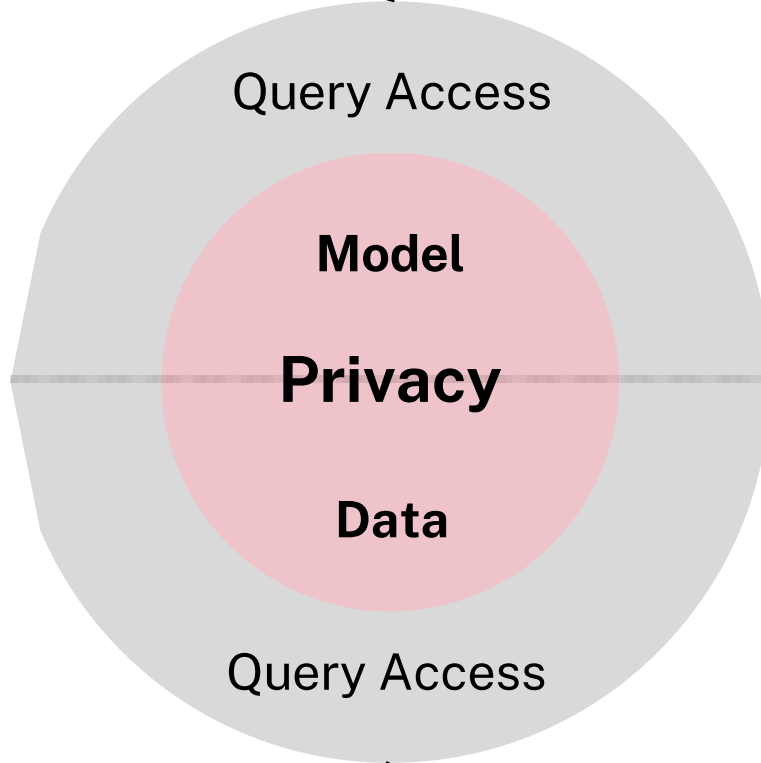








**Model Extraction**



Query Access

**Model**

**Privacy**

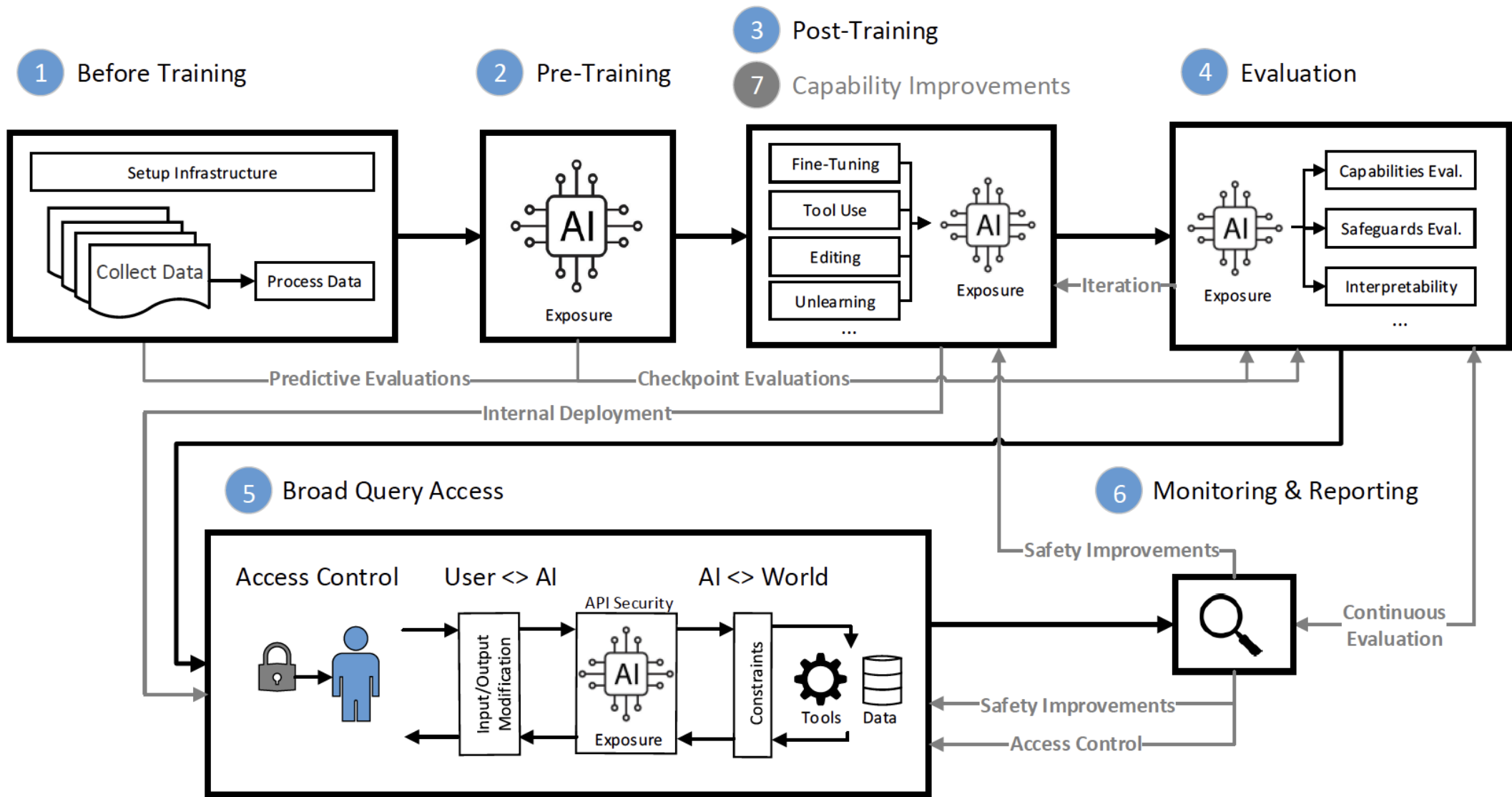
**Data**

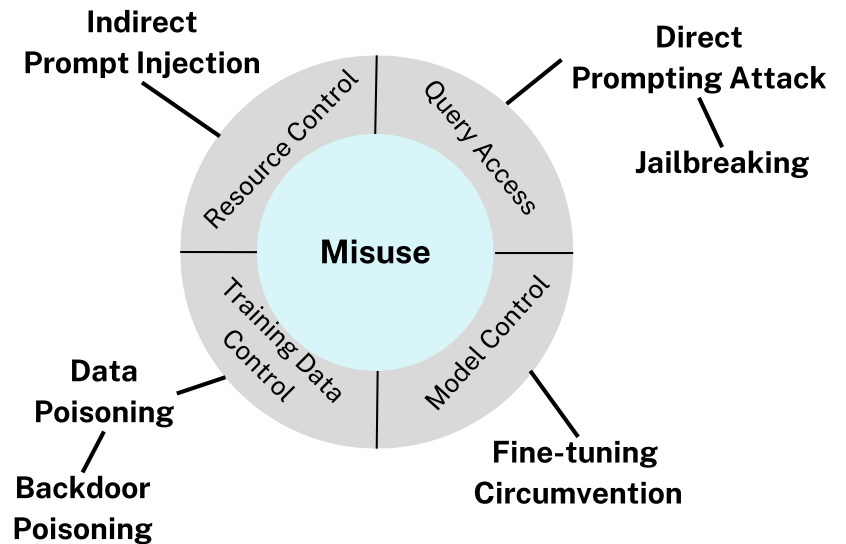
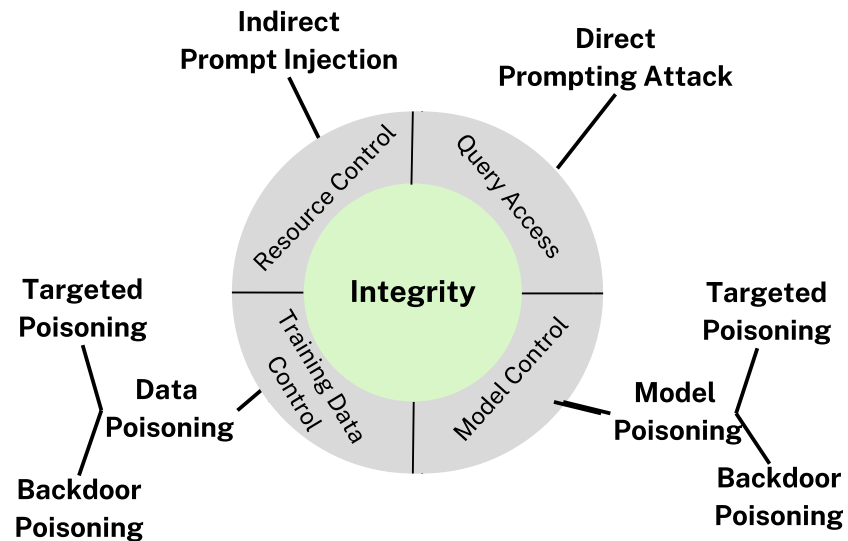
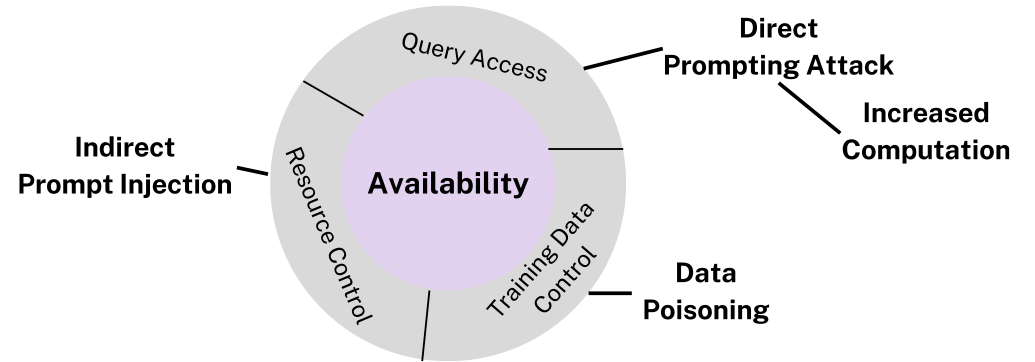
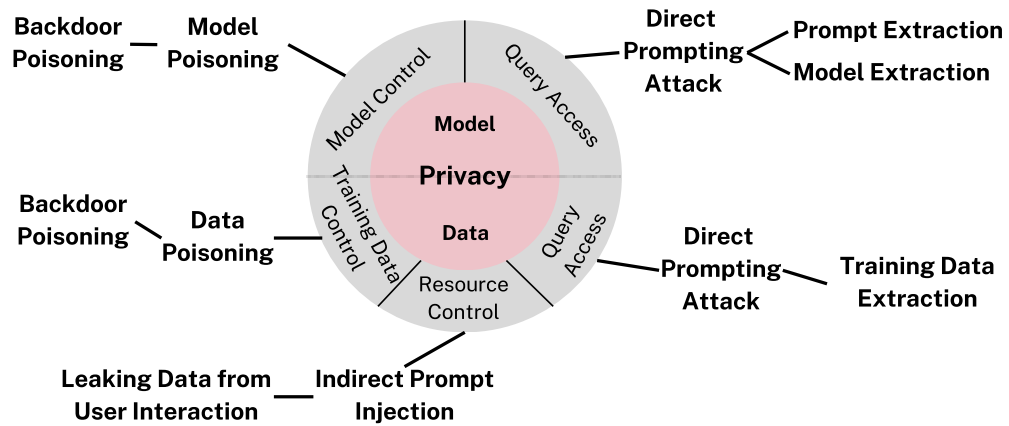
Query Access

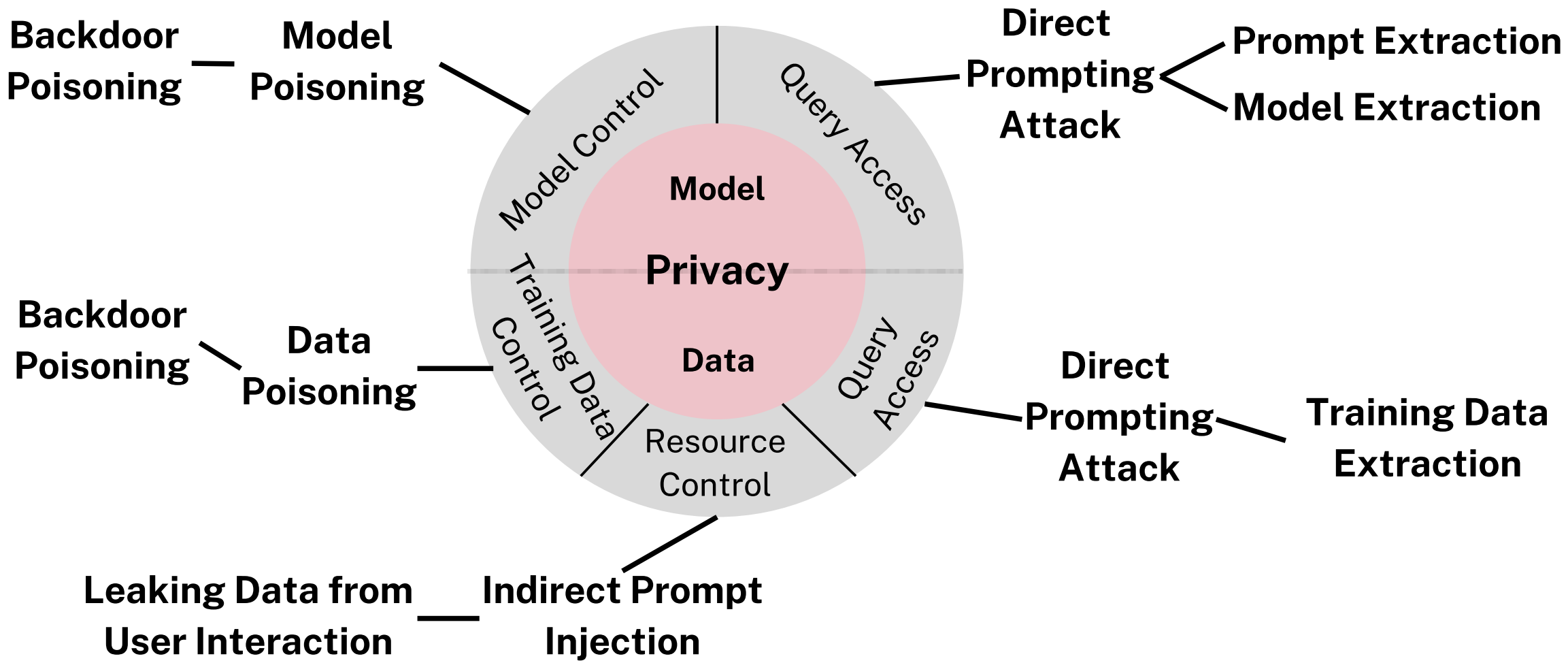


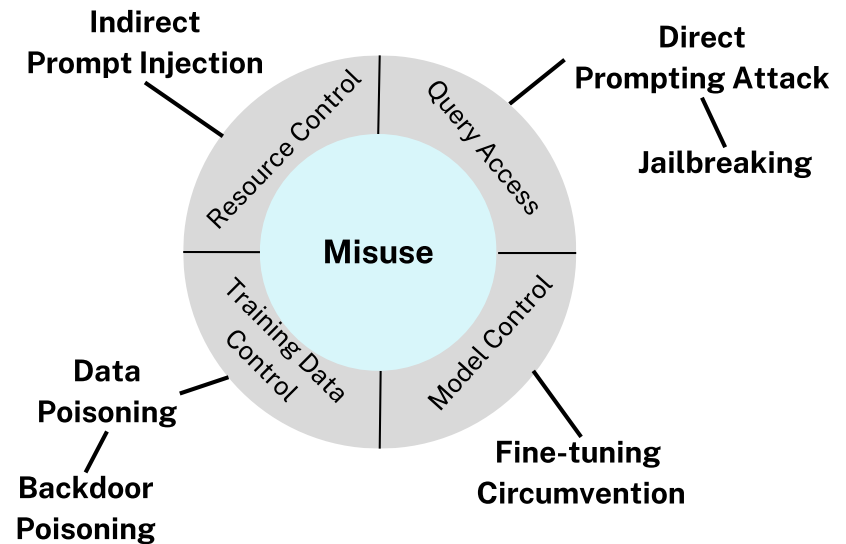
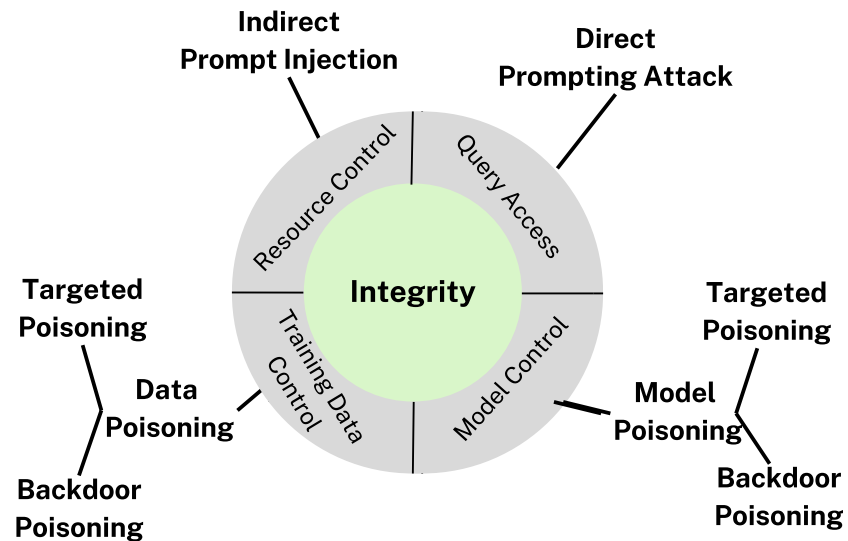
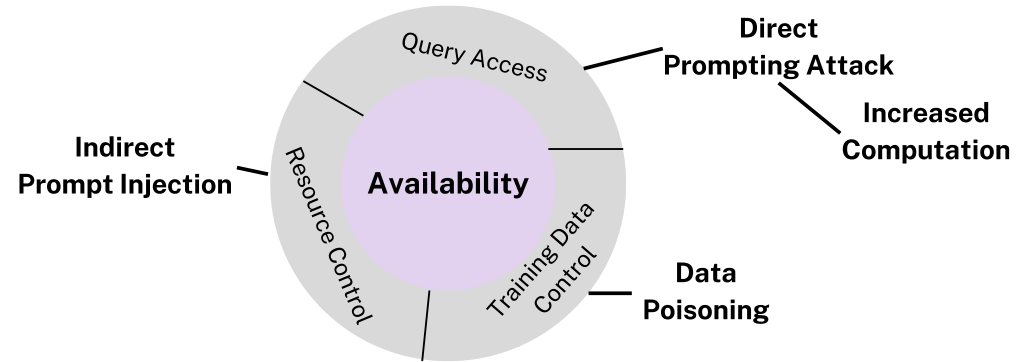
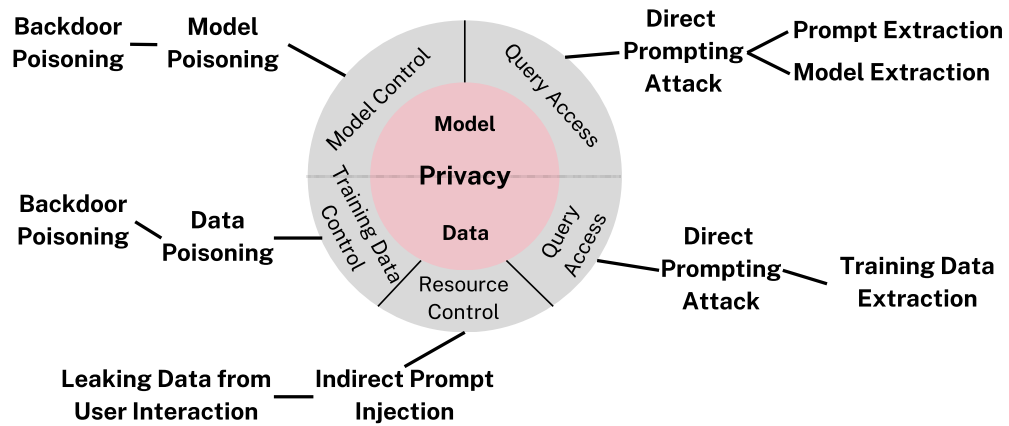
**Reconstruction;  
Membership Inference;  
Property Inference**

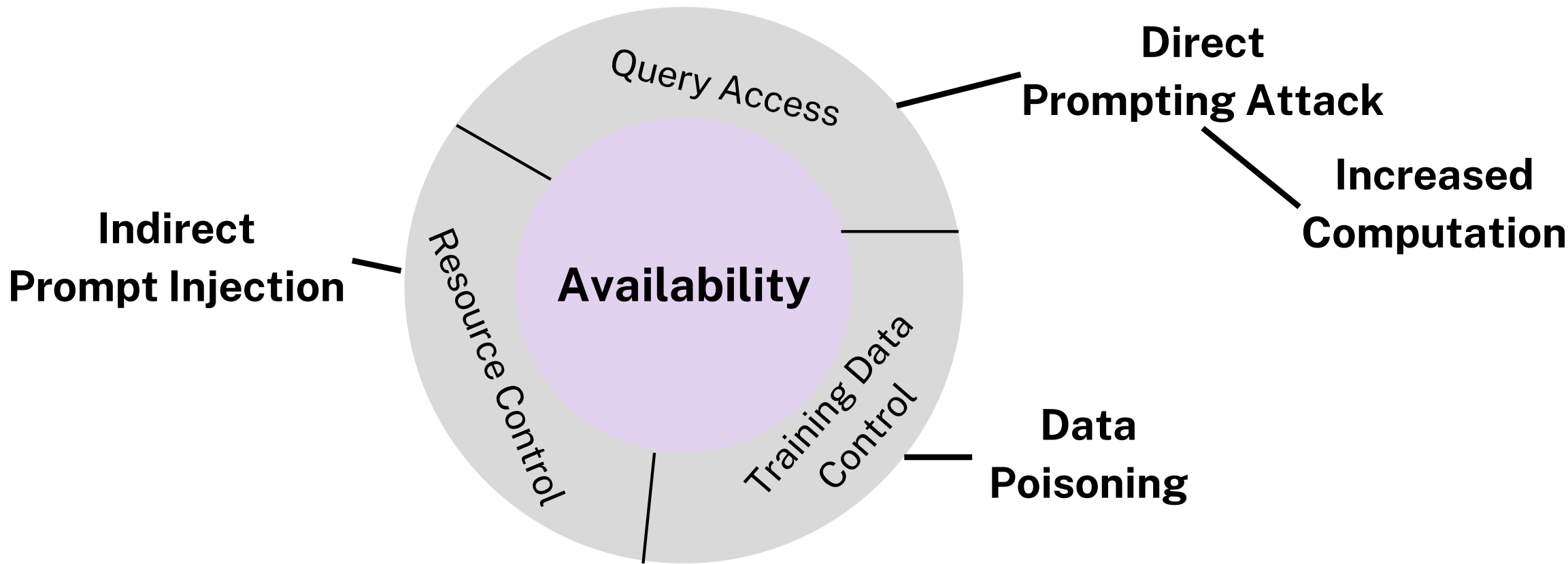
GenAI

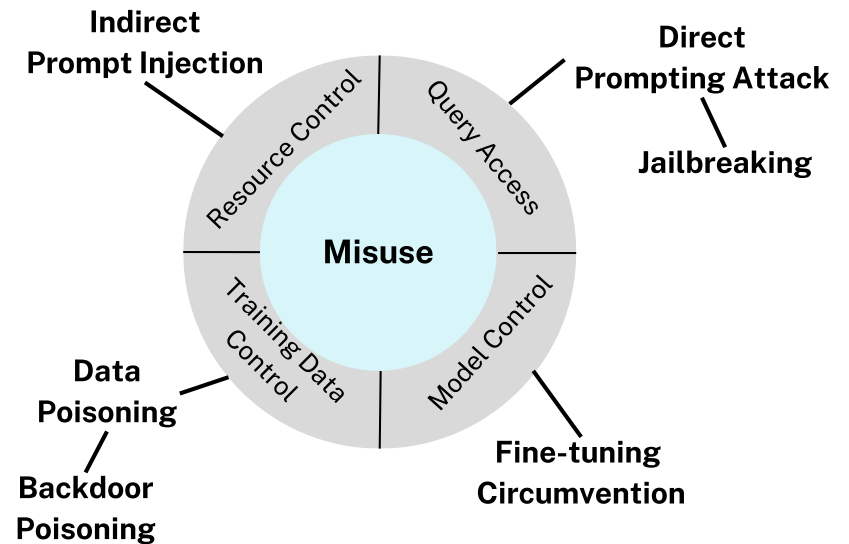
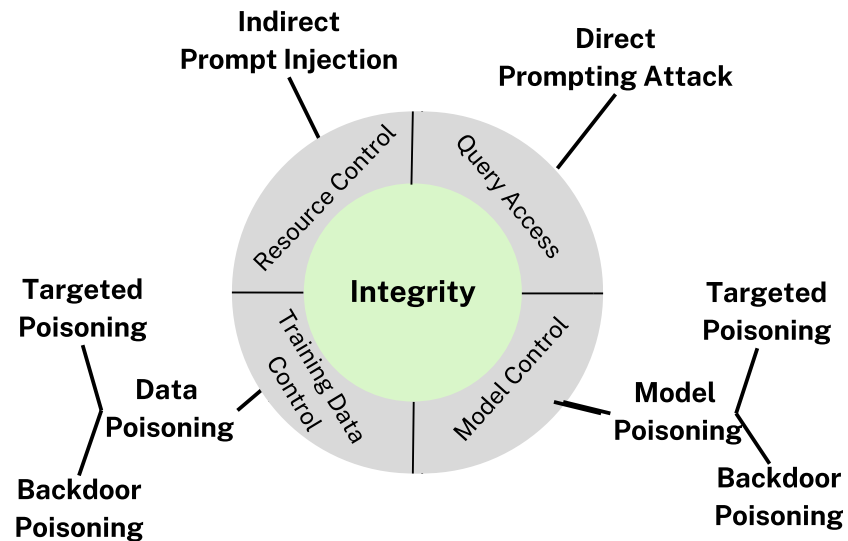
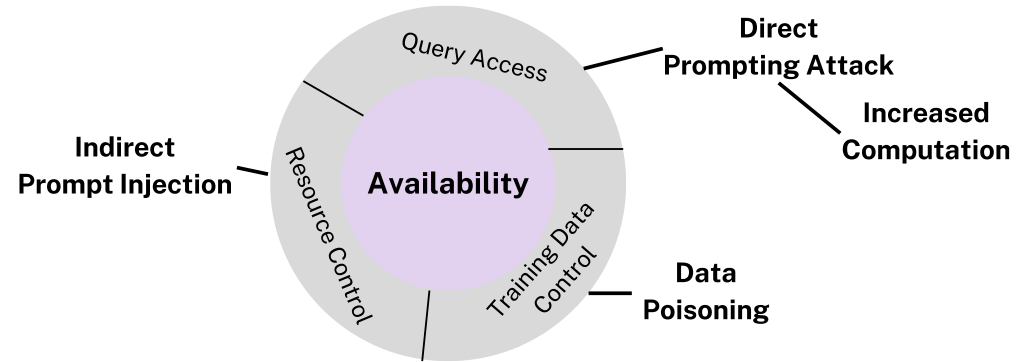
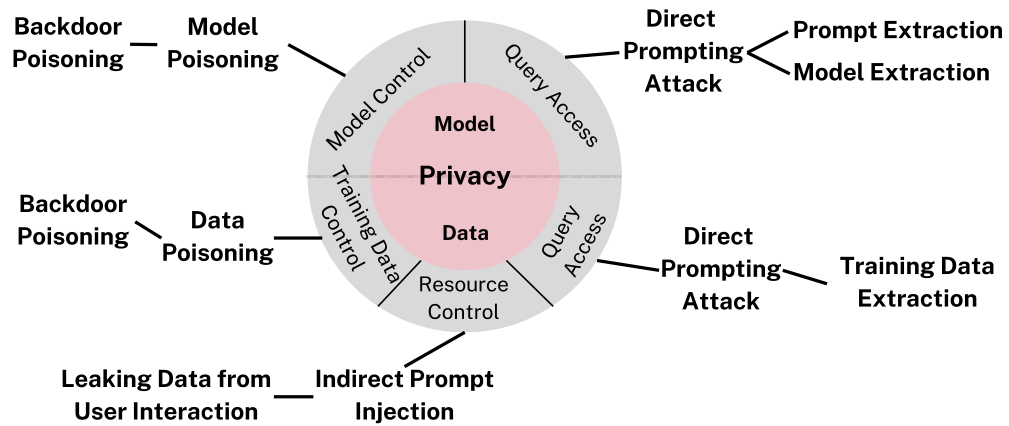


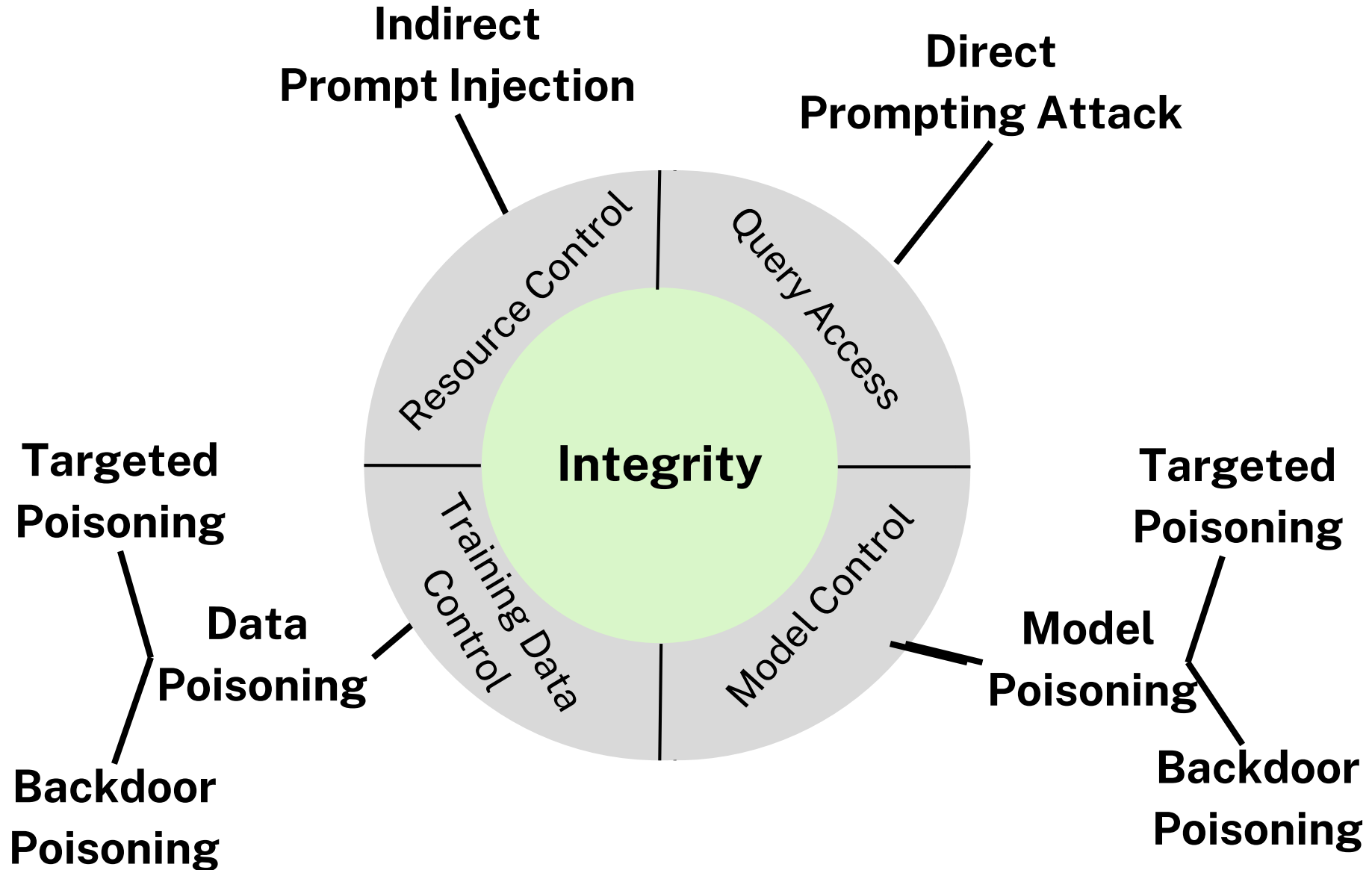


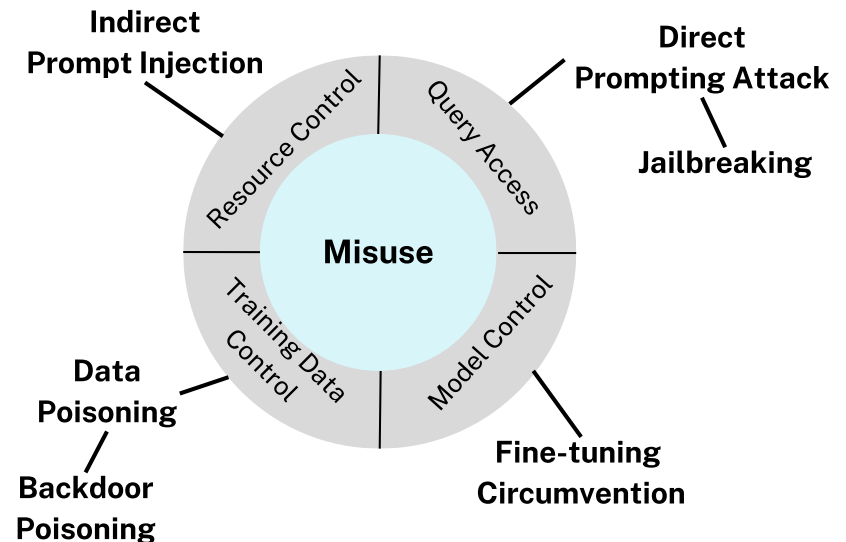
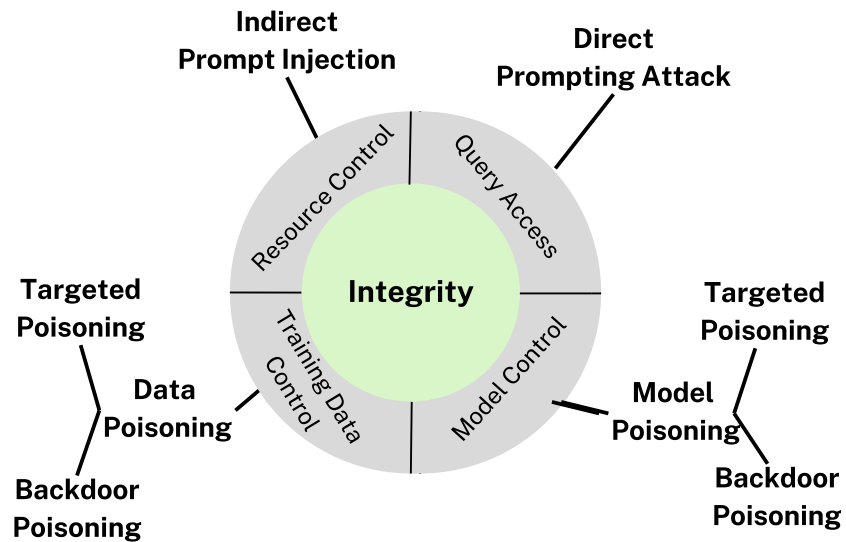
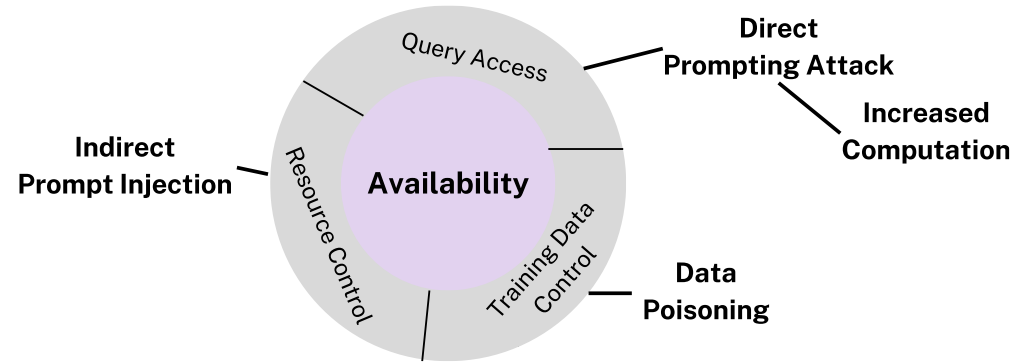
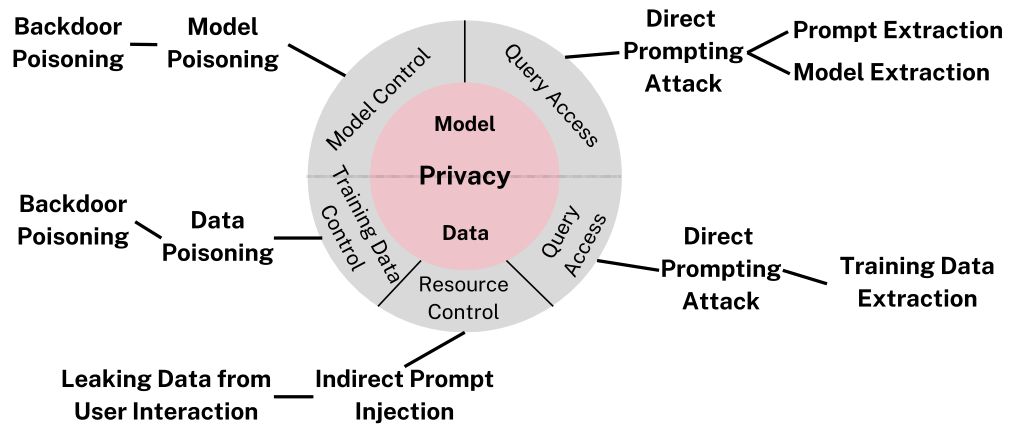


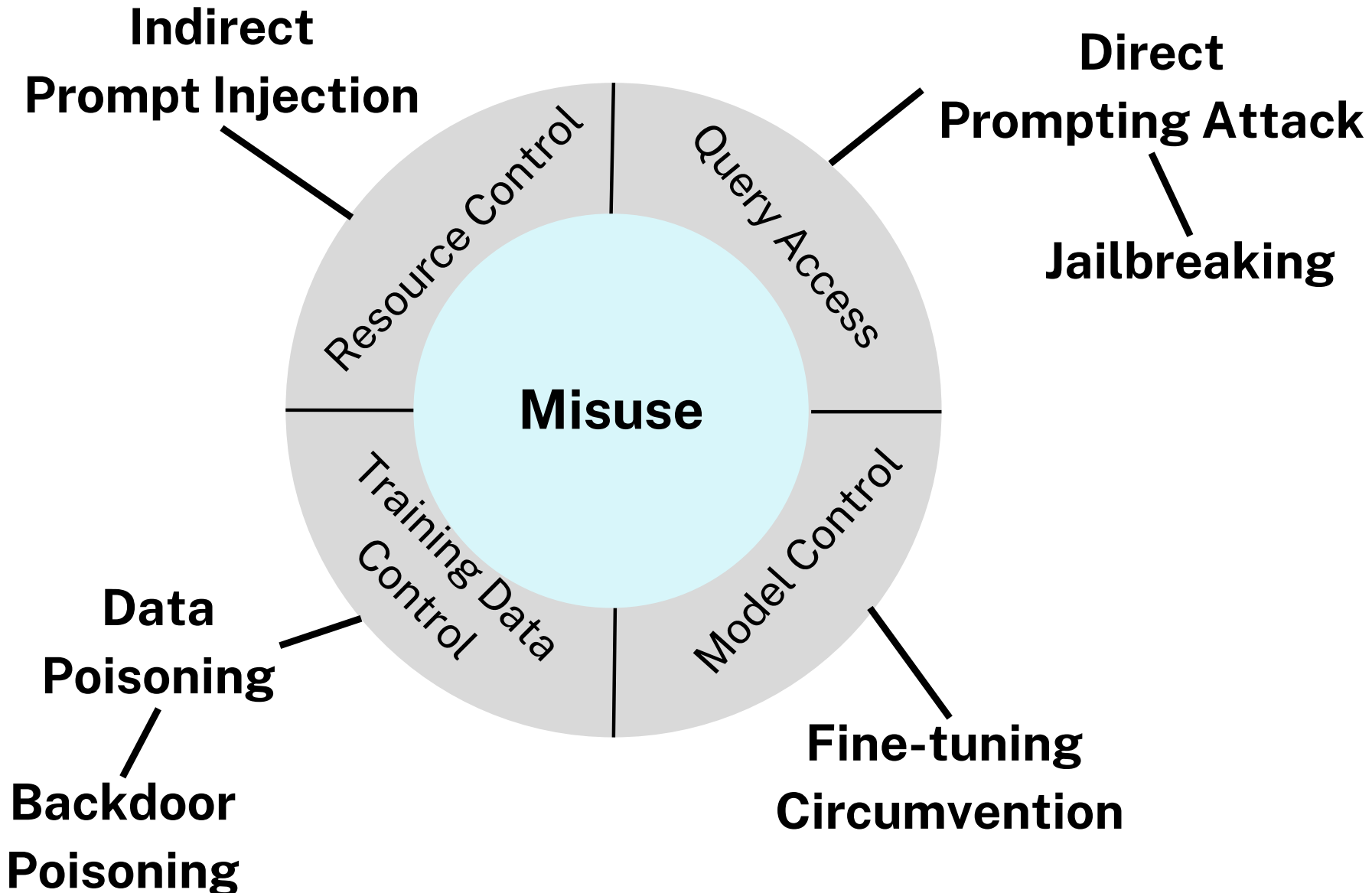


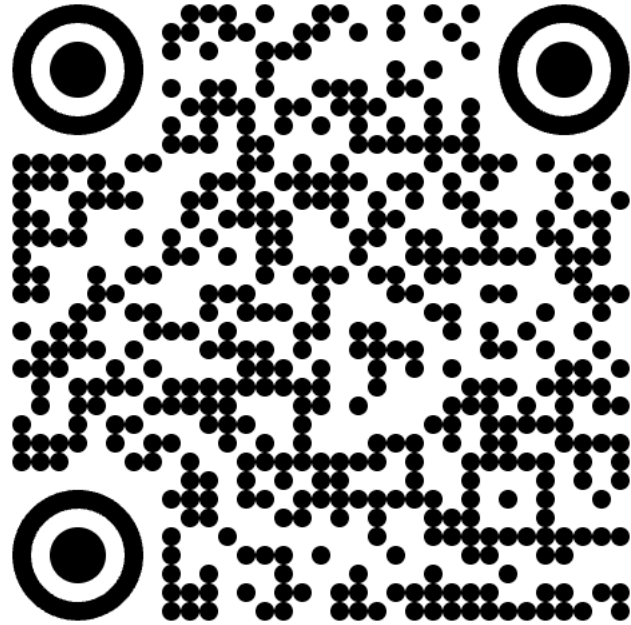








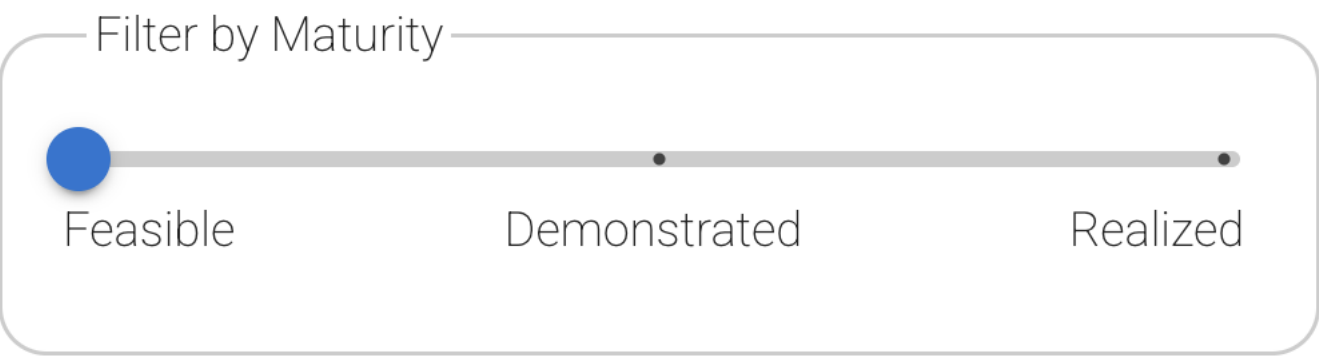




<https://atlas.mitre.org/matrices/ATLAS>

# ATLAS Matrix

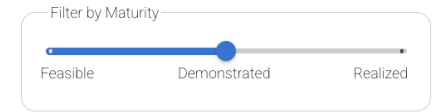
The ATLAS Matrix below shows the progression of tactics used in attacks as columns from left to right, with ML tec about each item, or search and view ATLAS tactics and techniques using the links at the top navigation bar. View th



Reconnaissance &	Resource Development &	Initial Access &	AI Model Access	Execution &	Persistence &	Privilege Escalation &	Defense Evasion &	Credential Access &	Discovery &	Lateral Movement &	Collection &	AI Attack Staging	Command and Control &	Exfiltration &	Impact &
8 techniques	13 techniques	7 techniques	4 techniques	6 techniques	8 techniques	4 techniques	13 techniques	6 techniques	9 techniques	2 techniques	4 techniques	6 techniques	3 techniques	6 techniques	8 techniques
Active Scanning &	Acquire Infrastructure	AI Supply Chain Compromise	AI Model Inference API Access	AI Agent Clickbait	AI Agent Context Poisoning	AI Agent Tool Invocation	Corrupt AI Model	AI Agent Tool Credential Harvesting	Cloud Service Discovery &	Phishing &	AI Artifact Collection	Craft Adversarial Data	AI Agent	Exfiltration via AI Agent Tool Invocation	Cost Harvesting
Gather RAG-Indexed Targets	Acquire Public AI Artifacts	Drive-by Compromise &	AI-Enabled Product or Service	AI Agent Tool Invocation	AI Agent Tool Data Poisoning	Escape to Host &	Delay Execution of LLM Instructions	Credentials from AI Agent Configuration	Discover AI Agent Configuration	Use Alternate Authentication Material &	Data from AI Services	Create Proxy AI Model	AI Service API	Exfiltration via AI Inference API	Data Destruction via AI Agent Tool Invocation
Gather Victim Identity Information &	Develop Capabilities &	Evade AI Model	Full AI Model Access	Command and Scripting Interpreter &	LLM Prompt Self-Replication	LLM Jailbreak	Evade AI Model	Exploitation for Credential Access &	Discover AI Artifacts		Data from Information Repositories &	Generate Deepfakes	Reverse Shell	Exfiltration via Cyber Means	Denial of AI Service
Search Application Repositories	Establish Accounts &	Exploit Public-Facing Application &	Physical Environment Access	Deploy AI Agent	Manipulate AI Model	Valid Accounts &	Exploitation for Defense Evasion &	Exploitation for Credential Access &	Discover AI Model Family		Data from Local System &	Generate Malicious Commands		Extract LLM System Prompt	Erode AI Model Integrity
Search Open AI Vulnerability Analysis	LLM Prompt Crafting	Phishing &		LLM Prompt Injection	Modify AI Agent Configuration		False RAG Entry Injection	OS Credential Dumping &	Discover AI Model Ontology			Manipulate AI Model		LLM Data Leakage	Erode Dataset Integrity
Search Open Technical Databases &	Obtain Capabilities &	Prompt Infiltration via Public-Facing Application		User Execution &	Poison Training Data		Impersonation &	RAG Credential Harvesting	Discover AI Model Outputs		Verify Attack		LLM Response Rendering	Evade AI Model	
Search Open Websites/Domains &	Poison Training Data	Valid Accounts &			Prompt Infiltration via Public-Facing Application		LLM Jailbreak	Unsecured Credentials &	Discover LLM Hallucinations					External Harms	
Search Victim-Owned Websites &	Publish Hallucinated Entities				RAG Poisoning		LLM Prompt Obfuscation		Discover LLM System Information					Spamming AI System with Chaff Data	
	Publish Poisoned AI Agent Tool						LLM Trusted Output Components Manipulation		Process Discovery &						
	Publish Poisoned Datasets						Manipulate User LLM Chat History								
	Publish Poisoned Models						Masquerading &								
	Retrieval Content Crafting						Modify AI Agent Configuration								
	Stage Capabilities &						Virtualization/Sandbox Evasion &								

# ATLAS Matrix

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. <sup>&</sup> 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 <sup>&amp;</sup> 8 techniques	Resource Development <sup>&amp;</sup> 13 techniques	Initial Access <sup>&amp;</sup> 7 techniques	AI Model Access 4 techniques	Execution <sup>&amp;</sup> 5 techniques	Persistence <sup>&amp;</sup> 7 techniques	Privilege Escalation <sup>&amp;</sup> 4 techniques	Defense Evasion <sup>&amp;</sup> 13 techniques	Credential Access <sup>&amp;</sup> 6 techniques	Discovery <sup>&amp;</sup> 8 techniques	Lateral Movement <sup>&amp;</sup> 2 techniques	Collection <sup>&amp;</sup> 4 techniques	AI Attack Staging 6 techniques	Command and Control <sup>&amp;</sup> 3 techniques	Exfiltration <sup>&amp;</sup> 4 techniques	Impact <sup>&amp;</sup> 6 techniques
Active Scanning <sup>&amp;</sup>	Acquire Infrastructure	AI Supply Chain Compromise	AI Model Inference API Access	AI Agent Tool Invocation	AI Agent Context Poisoning	AI Agent Tool Invocation	Corrupt AI Model	AI Agent Tool Credential Harvesting	Cloud Service Discovery <sup>&amp;</sup>	Phishing <sup>&amp;</sup> Use Alternate Authentication Material <sup>&amp;</sup>	AI Artifact Collection	Craft Adversarial Data	AI Agent	Exfiltration via AI Agent Tool Invocation	Data Destruction via AI Agent Tool Invocation
Gather RAG-Indexed Targets	Acquire Public AI Artifacts	Drive-by Compromise <sup>&amp;</sup>	AI-Enabled Product or Service	Command and Scripting Interpreter <sup>&amp;</sup>	LLM Prompt Self-Replication	Escape to Host <sup>&amp;</sup>	Delay Execution of LLM Instructions	Credentials from AI Agent Configuration	Discover AI Agent Configuration		Data from AI Services	Create Proxy AI Model	AI Service API	Exfiltration via Cyber Means	Denial of AI Service
Gather Victim Identity Information <sup>&amp;</sup>	Develop Capabilities <sup>&amp;</sup>	Evade AI Model	Full AI Model Access	Deploy AI Agent	Manipulate AI Model	LLM Jailbreak	Evade AI Model	Exploitation for Credential Access <sup>&amp;</sup>	Discover AI Artifacts	Data from Information Repositories <sup>&amp;</sup>	Generate Deepfakes	Reverse Shell	LLM Data Leakage	Erode AI Model Integrity	
Search Application Repositories	Establish Accounts <sup>&amp;</sup>	Exploit Public-Facing Application <sup>&amp;</sup>	Physical Environment Access	LLM Prompt Injection	Modify AI Agent Configuration	Valid Accounts <sup>&amp;</sup>	Exploitation for Defense Evasion <sup>&amp;</sup>	OS Credential Dumping <sup>&amp;</sup>	Discover AI Model Ontology	Data from Local System <sup>&amp;</sup>	Generate Malicious Commands		LLM Response Rendering	Erode Dataset Integrity	
Search Open AI Vulnerability Analysis	LLM Prompt Crafting	Phishing <sup>&amp;</sup>		User Execution <sup>&amp;</sup>	Poison Training Data		False RAG Entry Injection	RAG Credential Harvesting	Discover AI Model Outputs		Manipulate AI Model			Erode Dataset Integrity	
Search Open Technical Databases <sup>&amp;</sup>	Obtain Capabilities <sup>&amp;</sup>	Prompt Infiltration via Public-Facing Application			Prompt Infiltration via Public-Facing Application		Impersonation <sup>&amp;</sup>	Unsecured Credentials <sup>&amp;</sup>	Discover AI LLM Hallucinations		Verify Attack			Evade AI Model	
Search Open Websites/Domains <sup>&amp;</sup>	Poison Training Data	Valid Accounts <sup>&amp;</sup>			RAG Poisoning		LLM Jailbreak		Discover LLM Hallucinations					External Harms	
Search Victim-Owned Websites <sup>&amp;</sup>	Publish Hallucinated Entities						LLM Prompt Obfuscation		Discover LLM System Information						
	Publish Poisoned AI Agent Tool						LLM Trusted Output Components Manipulation		Process Discovery <sup>&amp;</sup>						
	Publish Poisoned Datasets						Manipulate User LLM Chat History								
	Publish Poisoned Models						Masquerading <sup>&amp;</sup>								
	Retrieval Content Crafting						Modify AI Agent Configuration								
	Stage Capabilities <sup>&amp;</sup>						Virtualization/Sandbox Evasion <sup>&amp;</sup>								



# LLM Prompt Injection

An adversary may craft malicious prompts as inputs to an LLM that cause the LLM to act in unintended ways. These "prompt injections" are often designed to cause the model to ignore aspects of its original instructions and follow the adversary's instructions instead.

Prompt Injections can be an initial access vector to the LLM that provides the adversary with a foothold to carry out other steps in their operation. They may be designed to bypass defenses in the LLM, or allow the adversary to issue privileged commands. The effects of a prompt injection can persist throughout an interactive session with an LLM.

Malicious prompts may be injected directly by the adversary ([Direct](#)) either to leverage the LLM to generate harmful content or to gain a foothold on the system and lead to further effects. Prompts may also be injected indirectly when as part of its normal operation the LLM ingests the malicious prompt from another data source ([Indirect](#)). This type of injection can be used by the adversary to a foothold on the system or to target the user of the LLM. Malicious prompts may also be [Triggered](#) user actions or system events.

**ID:** AML.T0051

**Subtechniques:** [LLM Prompt Injection: Direct](#), [LLM Prompt Injection: Indirect](#), [LLM Prompt Injection: Triggered](#)

**Tactic:** [Execution](#)

**Maturity:** Realized

**Number Of Case Studies:** 1

**Number Of Mitigations:** 6

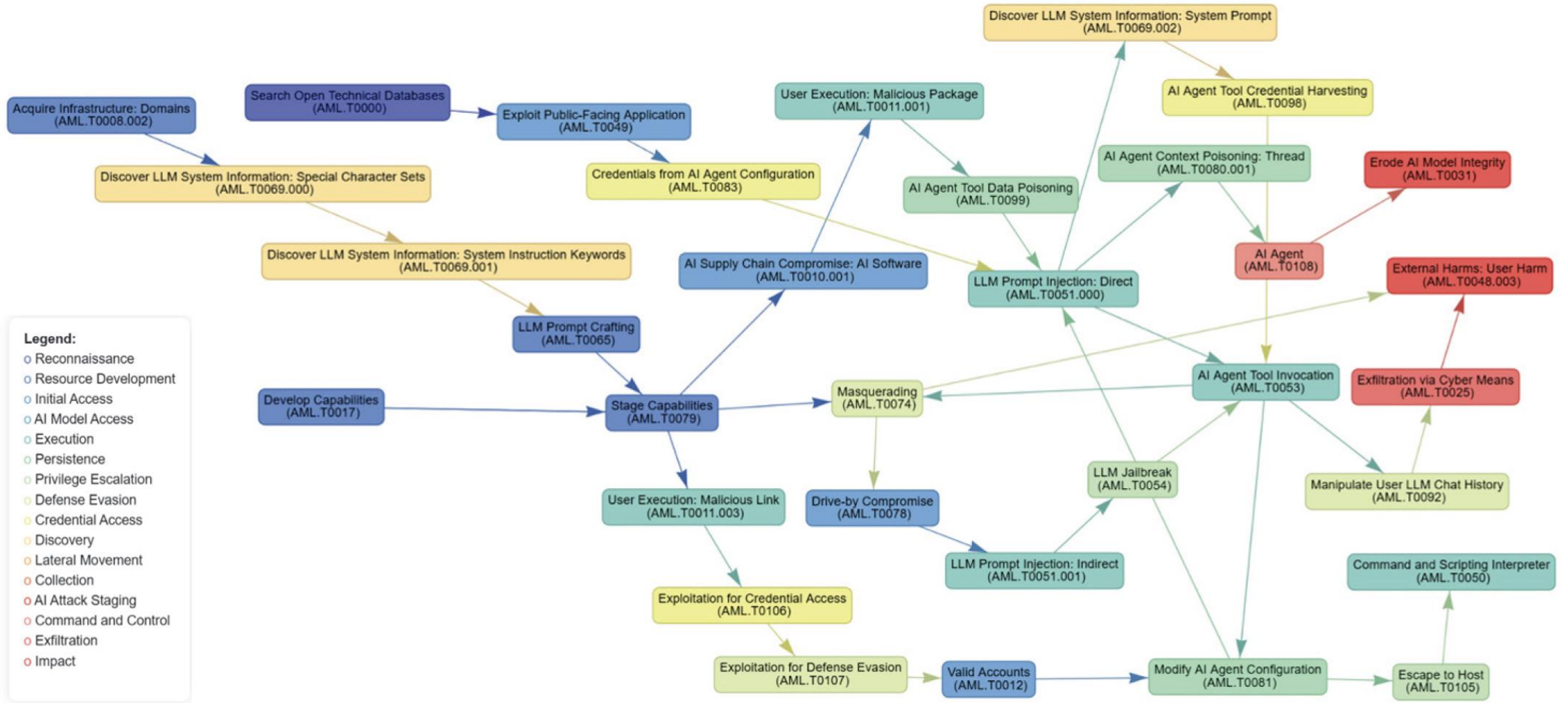
**Created:** 25 October 2023

**Last Modified:** 05 November 2025

# MITRE ATLAS OPENCLAW INVESTIGATION

<https://www.mitre.org/news-insights/publication/mitre-atlas-openclaw-investigation>





OPENCLAW ATTACK GRAPH

# Taxonomy of Failure Mode in Agentic AI Systems



<https://cdn-dynmedia-1.microsoft.com/is/content/microsoftcorp/microsoft/final/en-us/microsoft-brand/documents/Taxonomy-of-Failure-Mode-in-Agentic-AI-Systems-Whitepaper.pdf>

	Safety	Security
Novel	<ul style="list-style-type: none"> <li>• Intra-agent Responsible AI (RAI) issues</li> <li>• Harms of allocation in multi-user scenarios</li> <li>• Organizational knowledge loss</li> <li>• Prioritization leading to user safety issues</li> </ul>	<ul style="list-style-type: none"> <li>• Agent compromise</li> <li>• Agent injection</li> <li>• Agent impersonation</li> <li>• Agent flow manipulation</li> <li>• Agent provisioning poisoning</li> <li>• Multi-agent jailbreaks</li> </ul>
Existing	<ul style="list-style-type: none"> <li>• Insufficient transparency and accountability</li> <li>• Parasocial relationships</li> <li>• Bias amplification</li> <li>• User impersonation</li> <li>• Insufficient intelligibility for meaningful consent</li> <li>• Hallucinations</li> <li>• Misinterpretation of instructions</li> </ul>	<ul style="list-style-type: none"> <li>• Memory poisoning and theft</li> <li>• Targeted knowledge base poisoning</li> <li>• XPIA</li> <li>• Human-in-the-loop bypass</li> <li>• Function compromise and malicious functions</li> <li>• Incorrect permissions</li> <li>• Resource exhaustion</li> <li>• Insufficient isolation</li> <li>• Excessive agency</li> <li>• Loss of data provenance</li> </ul>



Image: Microsoft Designer

Build

Test

Deploy

Operate

### Modelscan

Scans downloaded model weights.

<https://github.com/provectai/modelscan>

### ZenML

Tracking lineage of data and models.

<https://github.com/zenml-io/zenml>



Image:  
Microsoft Designer

### Nvidia NeMo Guardrails

Rules for restricting interactions.

<https://github.com/NVIDIA/NeMo-Guardrails>

### META Llama Guard

Specialised LLM as security classifier.

<https://huggingface.co/meta-llama/Llama-Guard-4-12B>

### Langfuse

Visual traces of LLM and agent workflows.

<https://github.com/langfuse/langfuse>

### Arize Phoenix

Context poisoning monitoring/detector.

<https://phoenix.arize.com/>

---

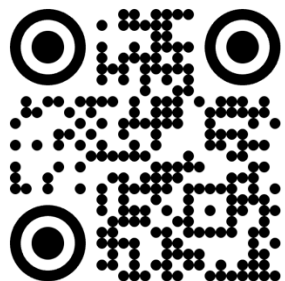
## NVIDIA Garak

The automated vulnerability scanner (Nmap for LLMs).

Fast, broad-strokes probing for known CVEs and basic jailbreaks.

Security Auditors establishing initial safety baselines.

<https://garak.ai/>



## Microsoft PyRIT

The advanced persistent threat simulator.

Multi-turn, compounding conversational attacks and multi-modal testing.

Dedicated Red Teams testing complex AI agents.

<https://github.com/Azure/PyRIT>



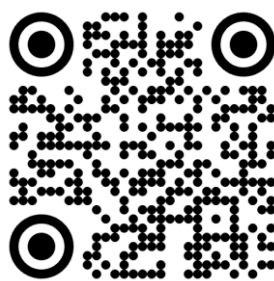
## Promptfoo

The DevSecOps pipeline defender.

Lightweight, YAML-based testing that runs automatically on every code push.

Developers & QA catching regressions during daily builds.

<https://promptfoo.dev/>



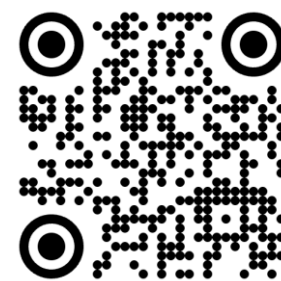
## DeepTeam

The out-of-the-box compliance checker.

Pre-loaded with 40+ vulnerabilities; automatically maps to OWASP and NIST.

Governance & Risk (GRC) teams needing quick compliance reports.

<https://trydeepteam.com/>



**BRAND NEW**  
**SUV**

**NOW WITH 0% APR!**

**SHOP NOW**

... ..



Image: ChatGPT 5.2

Master di Secondo Livello in  
Cybersecurity e Compliance  
Aziendale Integrata

Diritto della  
sicurezza  
11 CFU

Sicurezza  
delle reti  
7 CFU

Sicurezza  
dei Sistemi  
di Elab.  
Tradizionali  
13 CFU

Sicurezza  
dell'AI  
6 CFU

Elementi di Ingegneria dell'Informazione  
10 CFU

Federico Cerutti

[federico.cerutti@unibs.it](mailto:federico.cerutti@unibs.it)