

Threat Intel Report on

Mitre Atlas Framework

BY

Team SafeNet

Mayur Pawar - 2033

Rahul Gouda - 2069

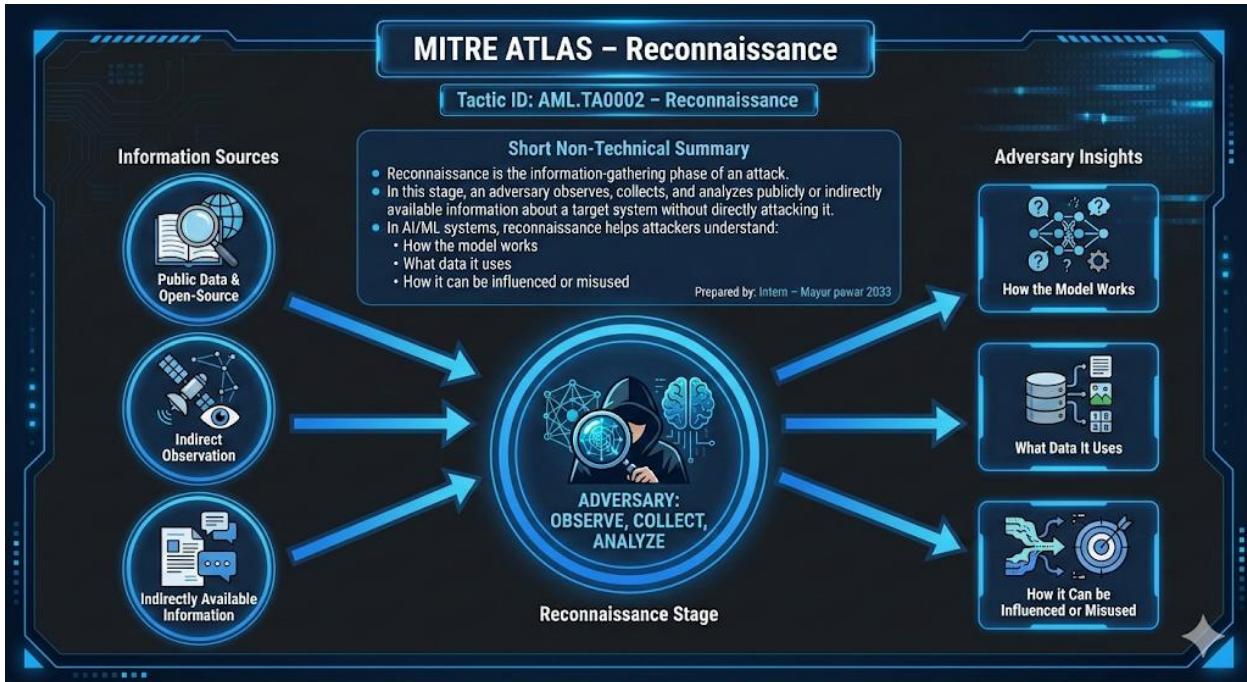
Janhavi Pandey – 2032

Sanika Jankar – 2038

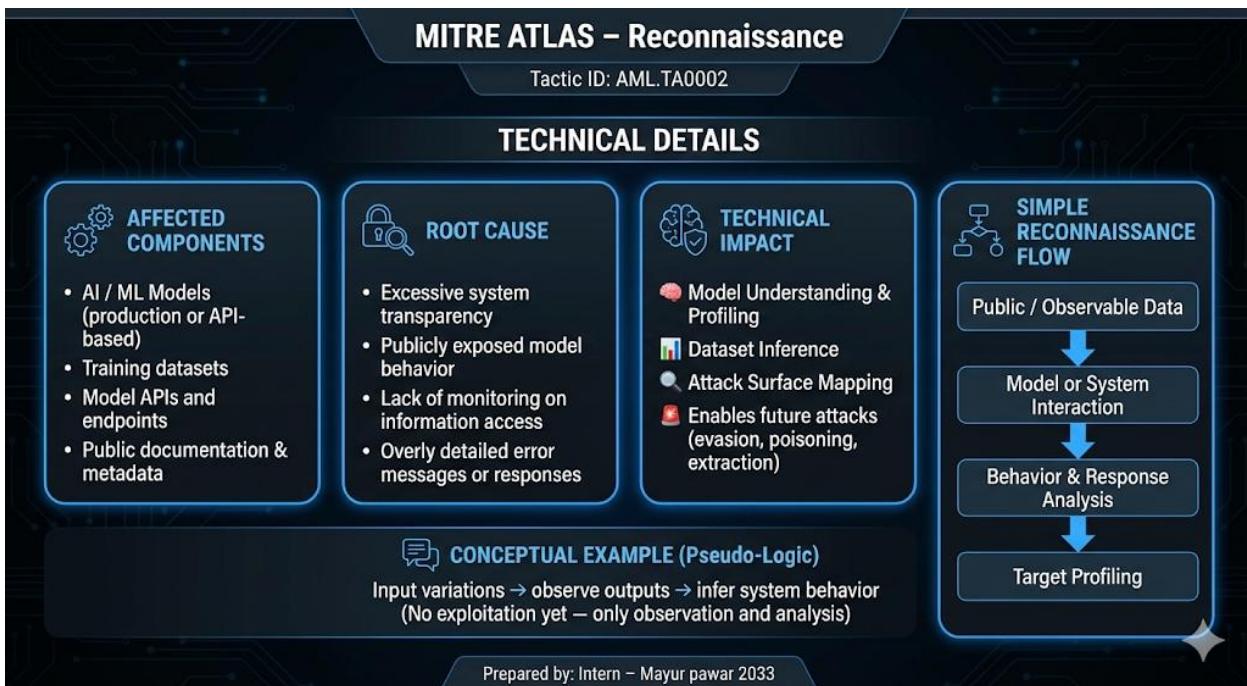
Shruti Rane - 2066

1. Reconnaissance

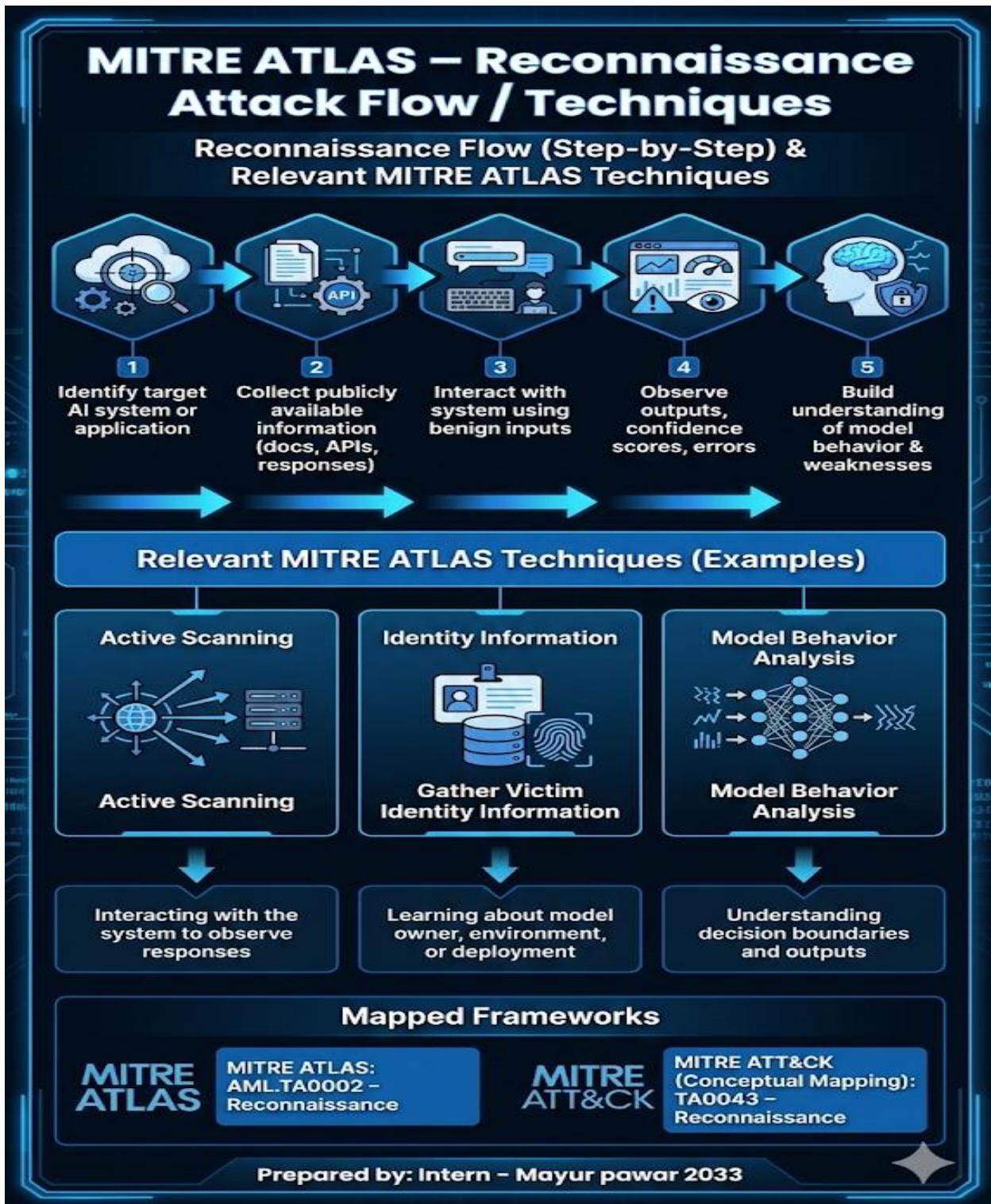
Overview



Technical Details:

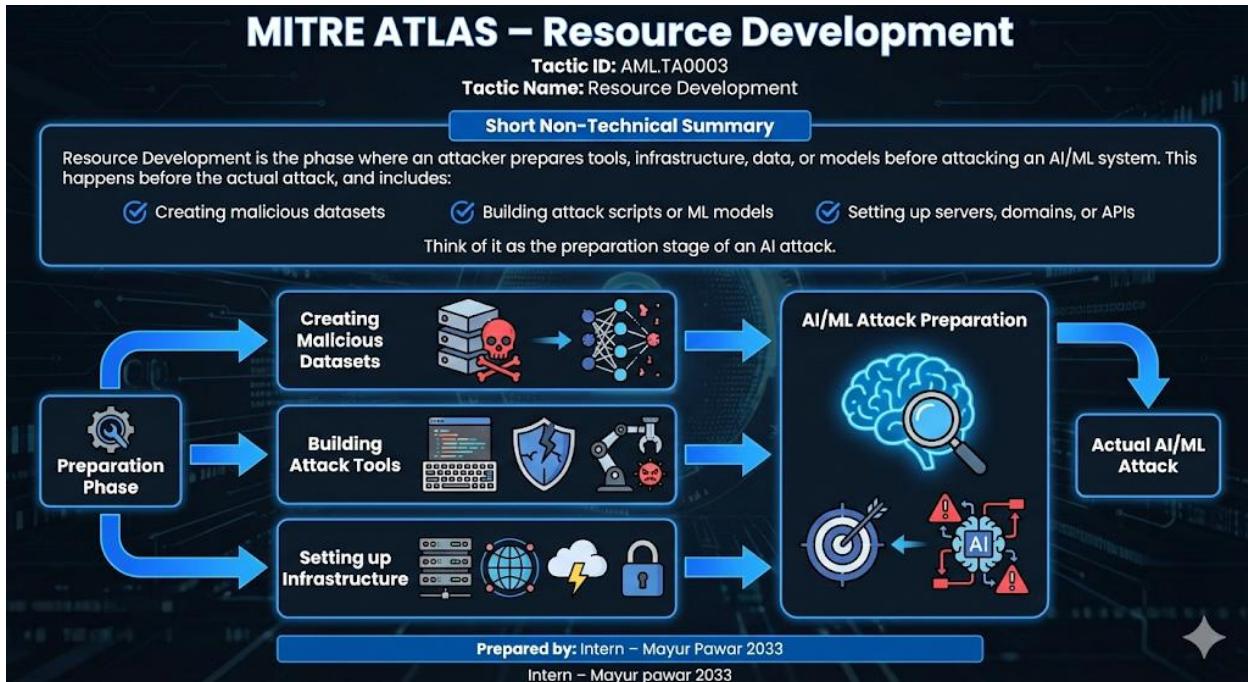


Attack Flow

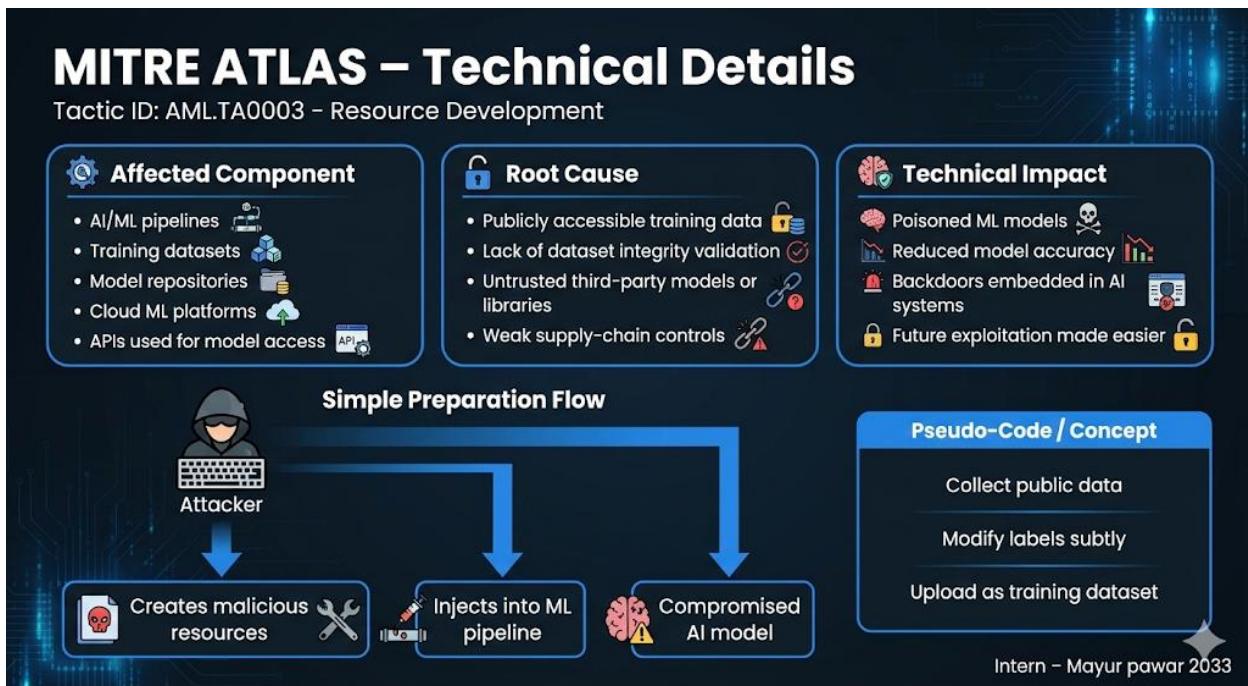


2. Resource Development

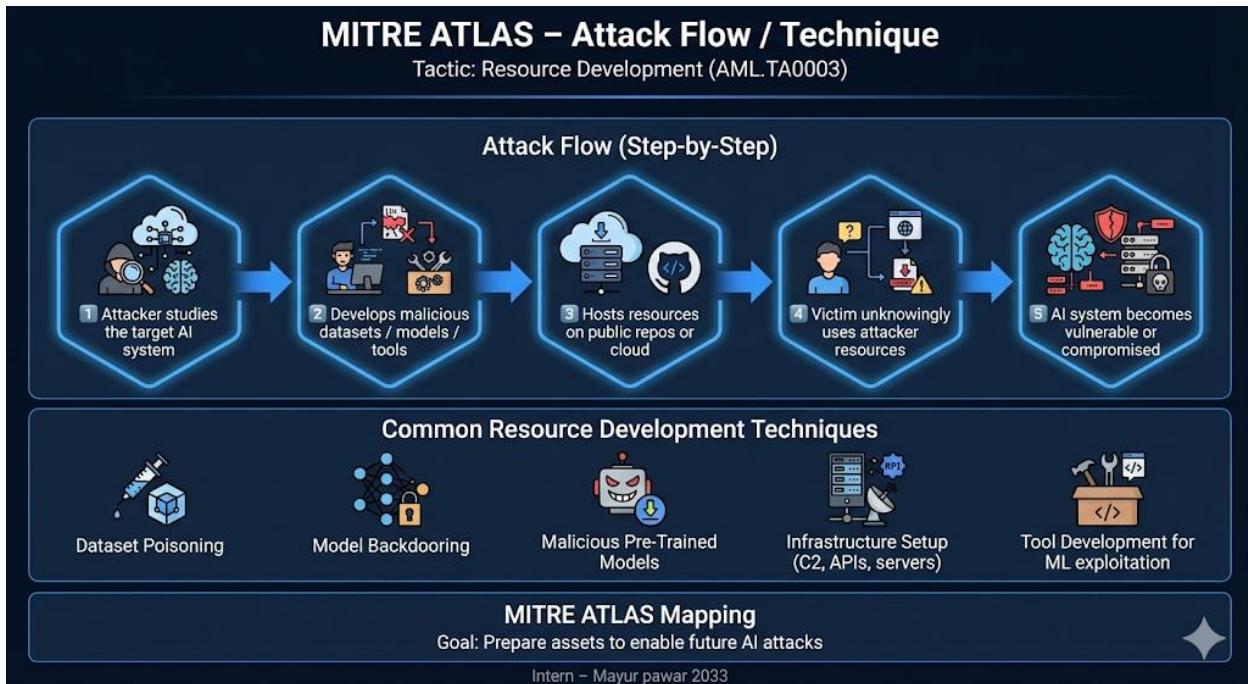
Overview



Technical Details:

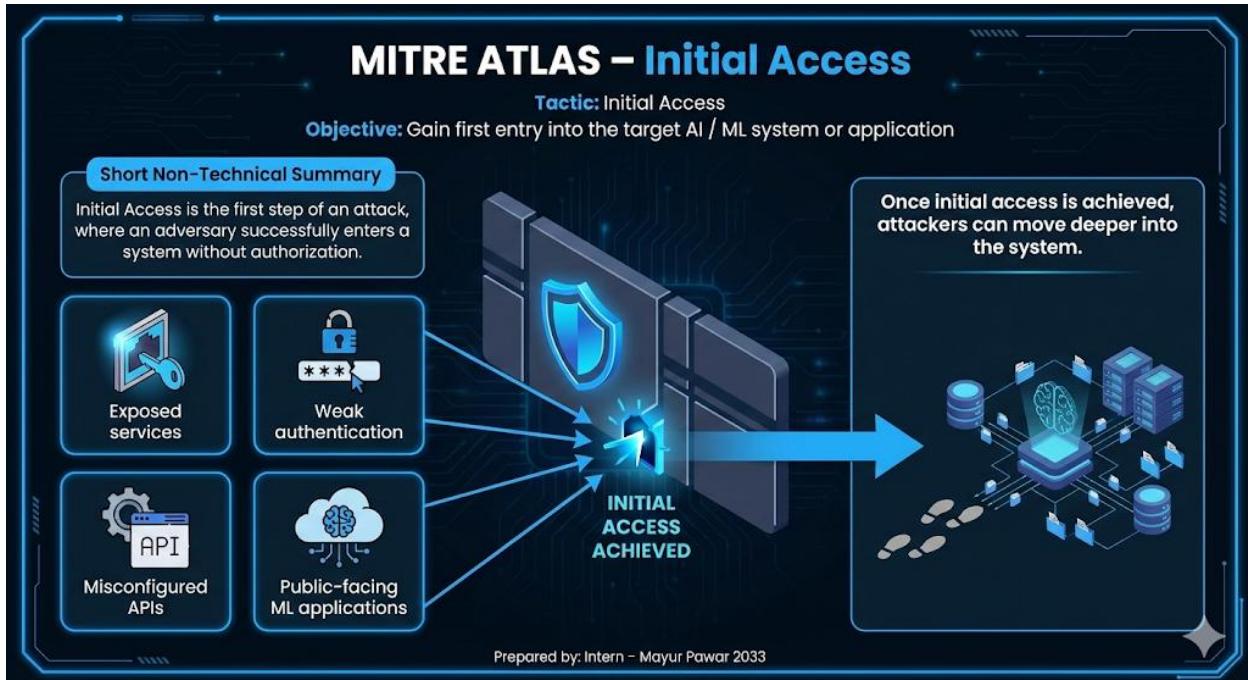


Attack Flow:

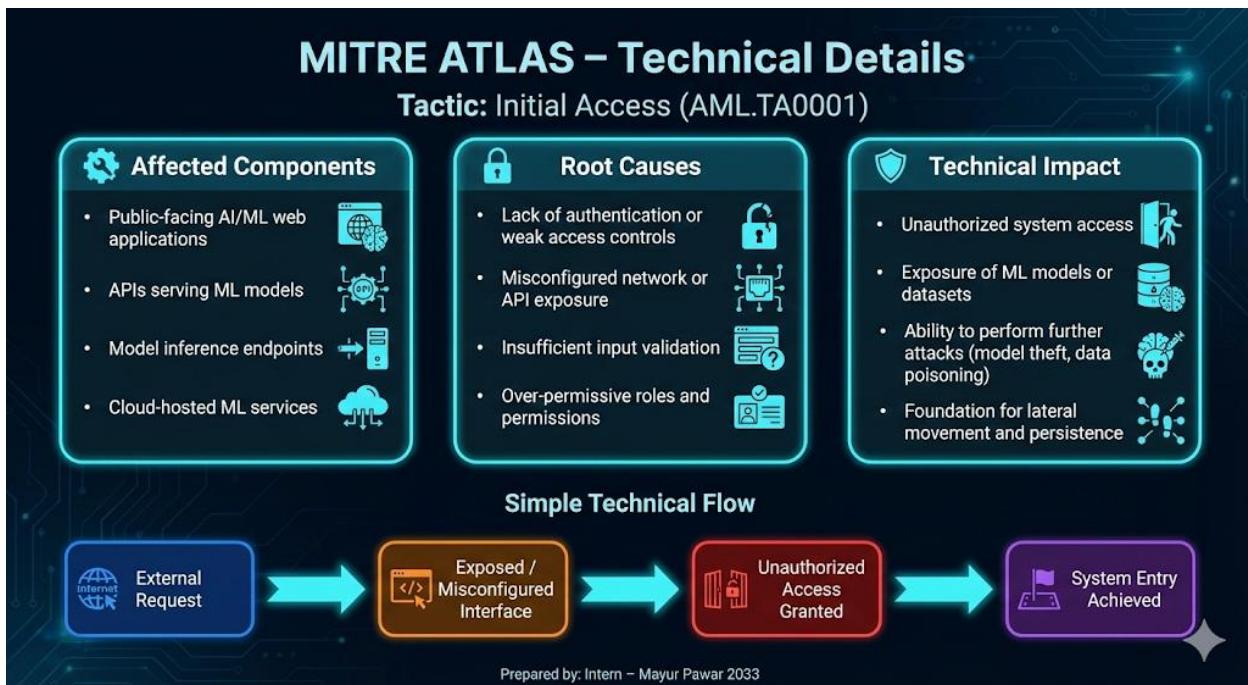


3. Initial Access

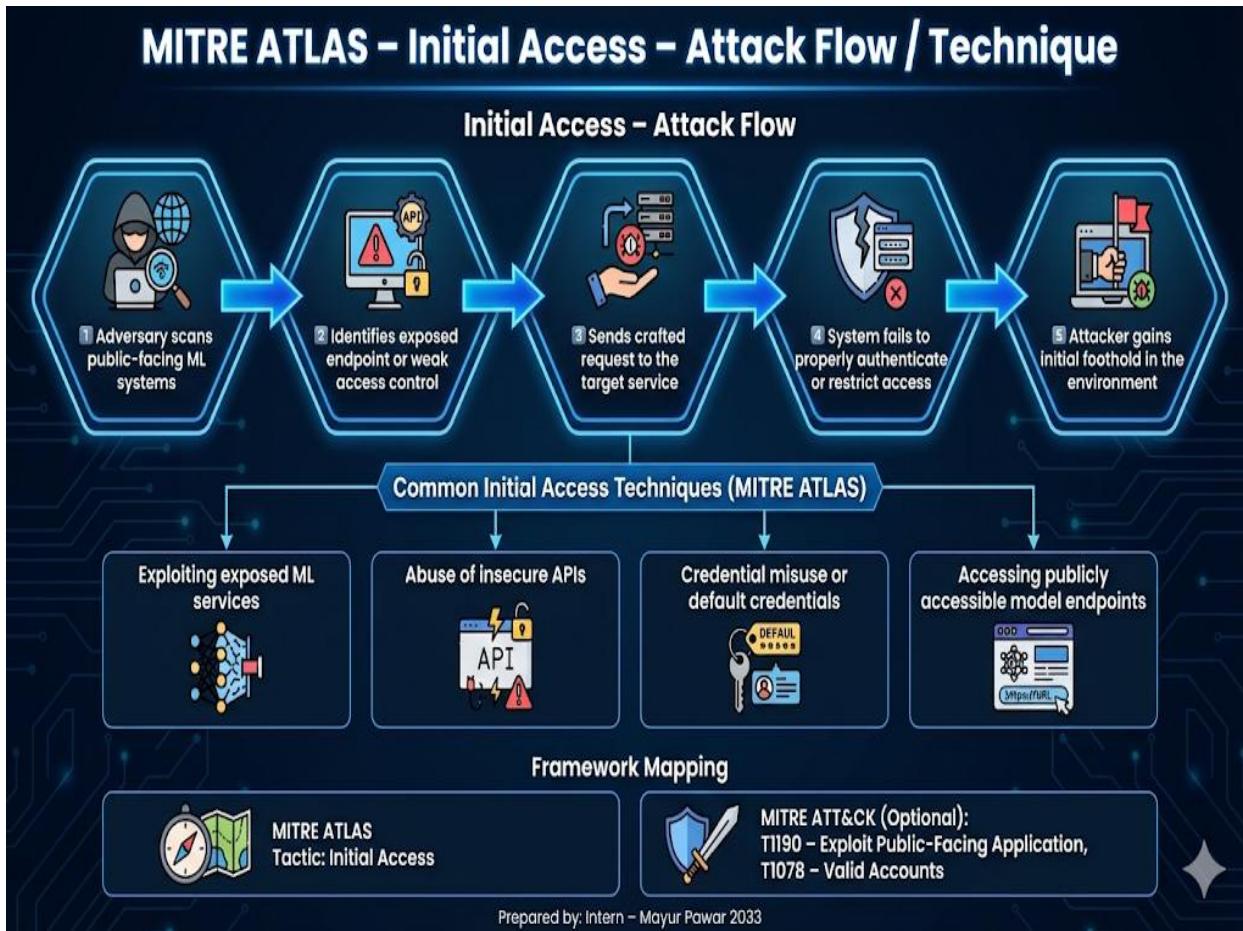
Overview



Technical Details:



Attack Flow:



4. AI Model Access

Overview

MITRE ATLAS - AI Model Access

Tactic Name: AI Model Access

AI MTRE Access is the stage where an attacker gains the ability to interact with AI model.

- Directly, through AI inference APIs
- Indirectly, through AI-enabled applications like chatbots or services

By accessing the model, attackers can:

- Observe how AI behaves
- Identify weaknesses or patterns for future harmful attacks
- This tactic commonly acts as an entry point for attacks targeting AI systems.

```
graph LR; 1[1 Attacker Observation] --> 2[2 Interact with AI-enabled App/API]; 2 --> 3[3 Access AI Model]; 3 --> 4[4 Observe AI Behavior & Patterns]; 4 --> 5[5 Prepare for Advanced Attacks]
```

Technical Details

MITRE ATLAS - AI MODEL ACCESS

TECHNICAL DETAILS

Affected Component

- AI inference APIs
- AI-enabled web or mobile applications

- Cloud-hosted AI services
- Public or unauthenticated model endpoints

Root Cause

- Missing or weak authentication
- Excessive trust in public access
- Lack of rate limiting and monitoring

Technical Impact

- Model behavior reconstruction
- Data exposure via AI responses
- Model misuse or abuse
- Gateway AI models attract AI attacks

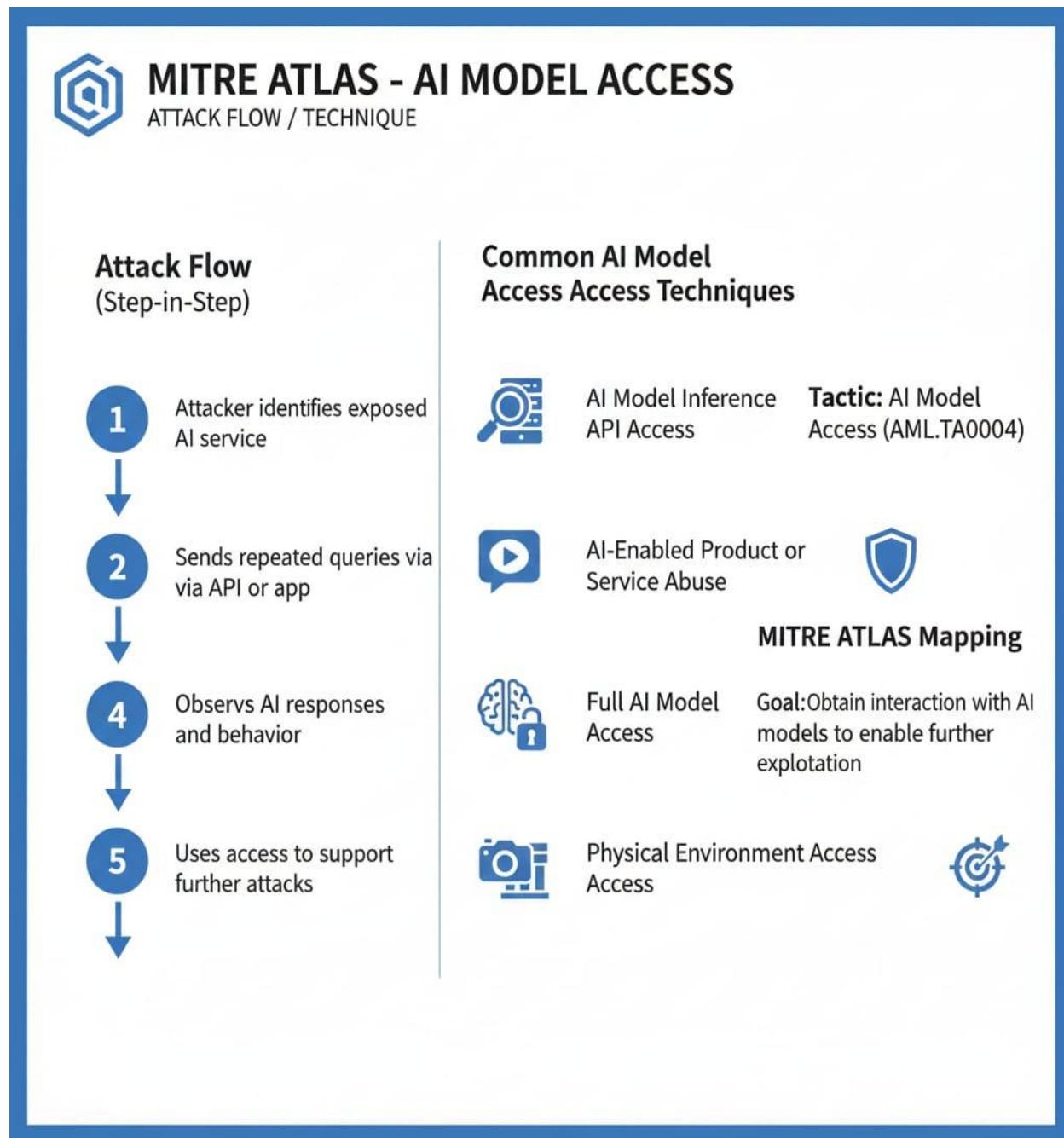
Simple Access Flow

```
graph LR; Attacker --> PAI[Public AI Interface / API]; PAI --> AIModel[AI Model]; AIModel --> UAO[Unrestricted AI Output]
```

Concept Pseudo Flow

- Send repeated prompts
 - Observe responses
 - Identify weaknesses

Attack Flow



5. Execution

Overview

MITRE ATLAS – Execution

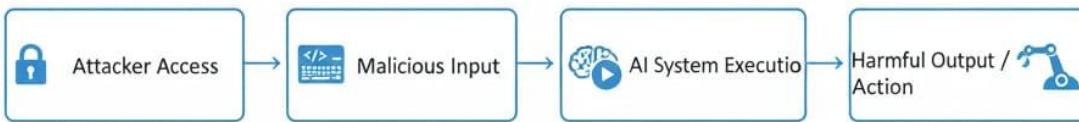
Tactic Name: Execution

Tactic non-technical summary

Execution refers in the stage in an AI attack where adversary actions, or carried by this system. In this tactic, attackers make AI models, agents or AI-enabled software execute unintended instructions, such as generating or malicious, running commands, running commands, or triggering automated actions. This tactic represents active misuse of AI capabilities and its progress toward data theft, system abuse, or data theft, specifically. This tactic represents where AI system actively participates in the attack.

Key Actions	Technical Impact	MITRE ATLAS Mapping
 Issue Malicious Prompts	 Harmful Content Generation	 Tactic: Execution (AML.TA0005)
 Inject Malicious Data	 System Command Execution	 Goal: Cause AI system to perform malicious actions
 Trigger Unauthorized Commands	 System Command Execution	
 Automate Harmful Tasks	 Data Leakage	
	 Service Disruption	

Execution Flow



Technical Details

MITRE ATLAS – Execution (Technical Details)					
Affected Component	Root Cause	Technical Impact			
<ul style="list-style-type: none">AI ModelsAI agentsAI-enabled applicationsOur AI servicesCloud AI servicesPrompt processing logic	<ul style="list-style-type: none">Input ValidationOver-privileged agentsExecution boundariesBlind trustPoor isolation	<ul style="list-style-type: none">Malicious instructionsUnprivileged actionsHarmful outputsAutomated attacker goalsData leakage			
Root Cause					
<h3>Simple Execution Flow</h3> <pre>graph LR; Attacker[Attacker] --> Prompt[Malicious Prompt / Instruction]; Prompt --> Agent[Agent]; Agent --> Execution[Untended Execution]; Execution --> Outcome[Malicious Outcome]</pre>					
<h3>Concept Flow</h3> <table><tbody><tr><td><ul style="list-style-type: none">Crafting promptBypassing safeguardsBlind trust</td><td><ul style="list-style-type: none">Triggers safeguardsTrigger executionObserving harmful behavior</td><td></td></tr></tbody></table>			<ul style="list-style-type: none">Crafting promptBypassing safeguardsBlind trust	<ul style="list-style-type: none">Triggers safeguardsTrigger executionObserving harmful behavior	
<ul style="list-style-type: none">Crafting promptBypassing safeguardsBlind trust	<ul style="list-style-type: none">Triggers safeguardsTrigger executionObserving harmful behavior				

Attack Flow

MITRE ATLAS – Execution (Attack Flow)

Attack Flow (Step-by-Step)



Attacker gains access to an AI system or application



Attacker crafts a malicious prompt or instruction



Prompt bypasses AI safety controls or guardrails



AI model or agent executes unintended actions



Malicious output or action is produced

MITRE ATLAS Mapping



Tactic: Execution (AML.TA0005)



Goal: Cause AI systems to perform unintended or malicious actions

Common Execution Techniques



- LLM Prompt Injection
- Guardrail Bypass
- Tool or Plugin Abuse
- AI Agent Command Execution
- Indirect Prompt Injection



6. Persistence

Overview

MITRE ATLAS – Persistence

Tactic Name: Persistence

- Persistence refers to the stage where attackers continue to influence an AI system for long periods of time.
- In this tactic, the attacker modifies AI configurations, data sources, or system behavior so that malicious effects remain active even after restarts, updates, or normal use.
- Persistence allows attackers to repeatedly exploit or manipulate AI systems without access each time. This tactic focuses on maintaining long-term presence within AI-enabled environments.

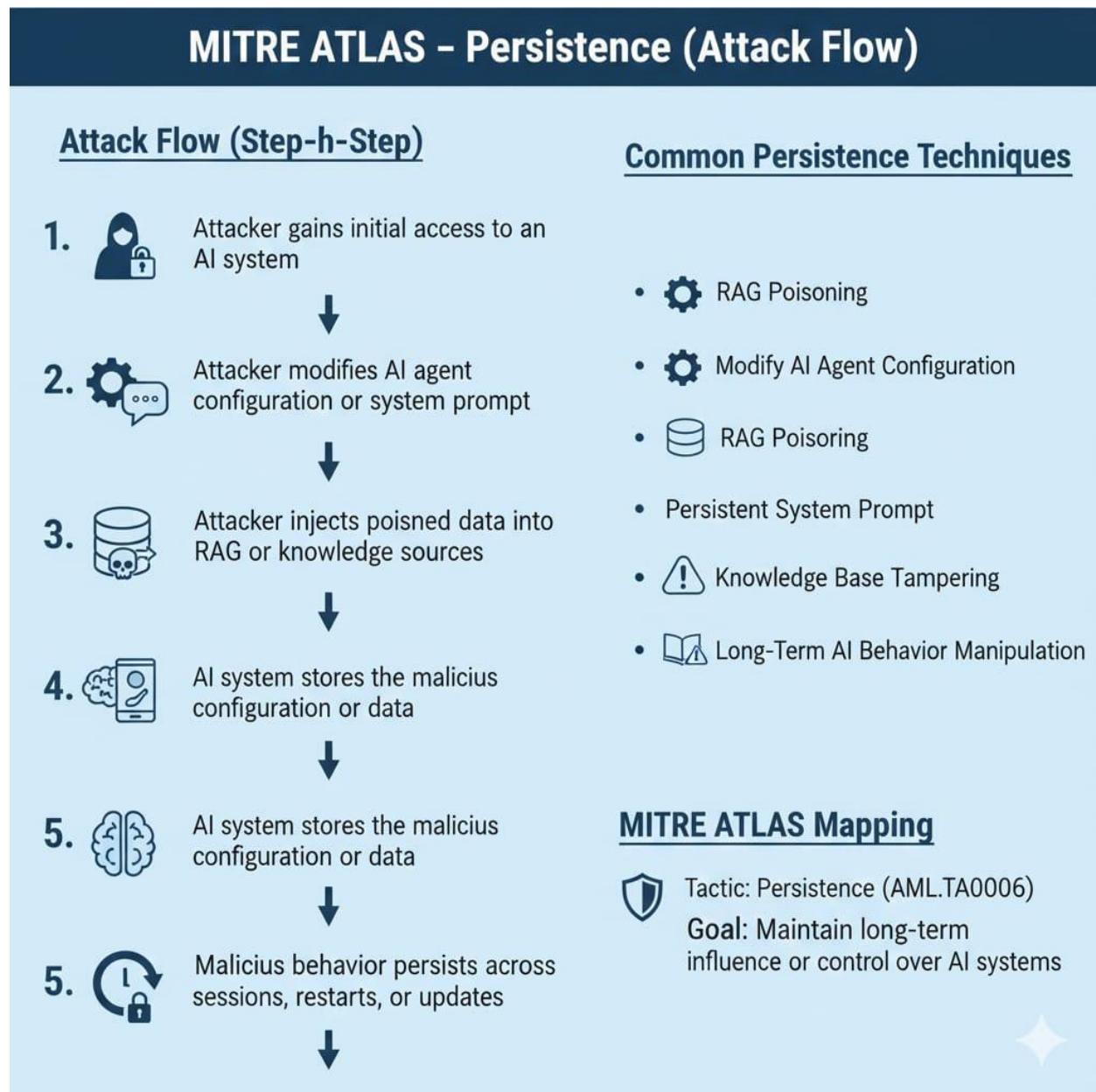
Key Actions	Technical Impact	MITRE ATLAS Mapping
 Modify AI Configurations	 Undetected Manipulation	 Tactic: Persistence (AML.TA0006)
 Backdoor Training Data	 Recurring Malicious Outputs	 Goal: Maintain long-term control over AI systems
 Hijack Automation Pipelines	 Persistent Access	
 Schedule Malicious Tasks	 Long-Term System Compromise	

Persistence Flow

```
graph LR; A[Attacker Access] --> B[Malicious Configuration / Data]; B --> C[AI System Modification]; C --> D[Persistent Control]
```



Technical Details



Attack Flow

MITRE ATLAS – Persistence (Technical Details)

Affected Component	Root Cause	Technical Impact
AI agent configurations and system prompts	Lack of integrity checks on AI configurations	Persistent manipulation of AI outputs
Retrieval-Augmented Generation (RAG) data sources	Unrestricted modification or system prompts	Long-term misdirection or bias in responses
AI knowledge bases and metadatabase bases and document stores	Unvalidated or untrusted data sources in RAG	Embedded malicious logic in AI agents
Model configuration files and policies	Poor access control or management on AI management interfaces	Repeated exploitation without re-access
AI workflow automation settings	Absence of monitoring for long-AI behavior changes	Increased difficulty in detection and recovery

Simple Persistence Flow

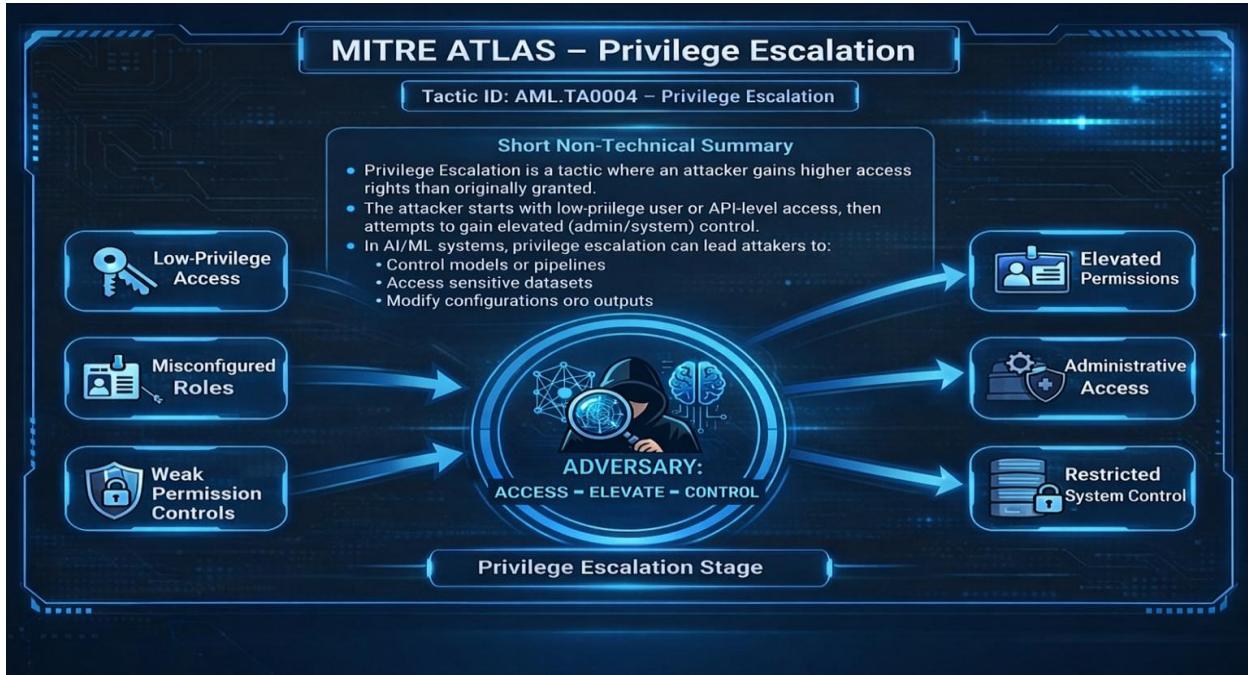
```
graph LR; A[Attacker] --> B[Modifies AI Configuration / Data Source]; B --> C[AI System Stores Changes]; C --> D[Changes Persist Across Sessions]; D --> E[Changed Persist Across Behavior]
```

Concept Flow

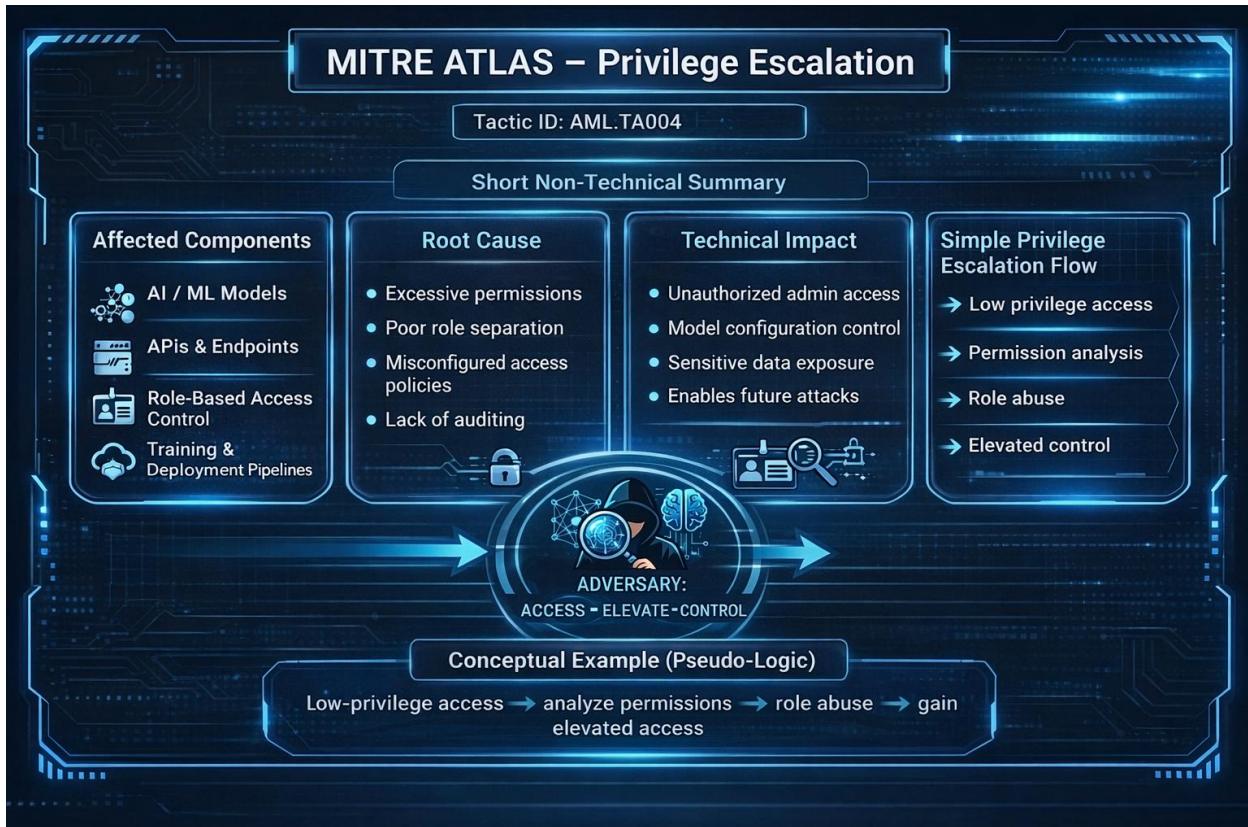
```
graph LR; A[Alter AI configuration or prompt] --> B[Changes saved by the system]; B --> C[Malicious behavior persists over time]
```

7. Privilege Escalation

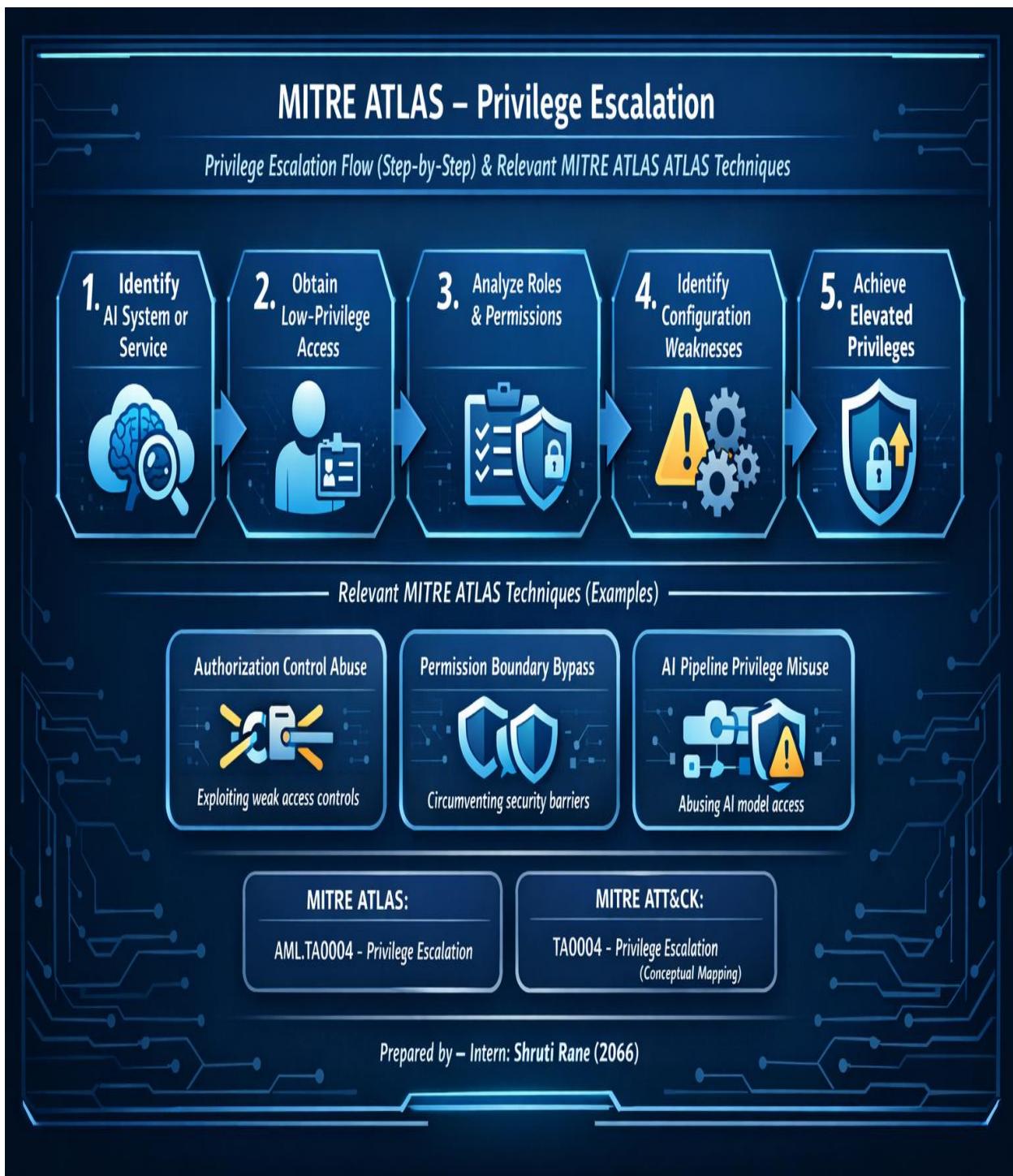
Overview



Technical Details



Attack Flow

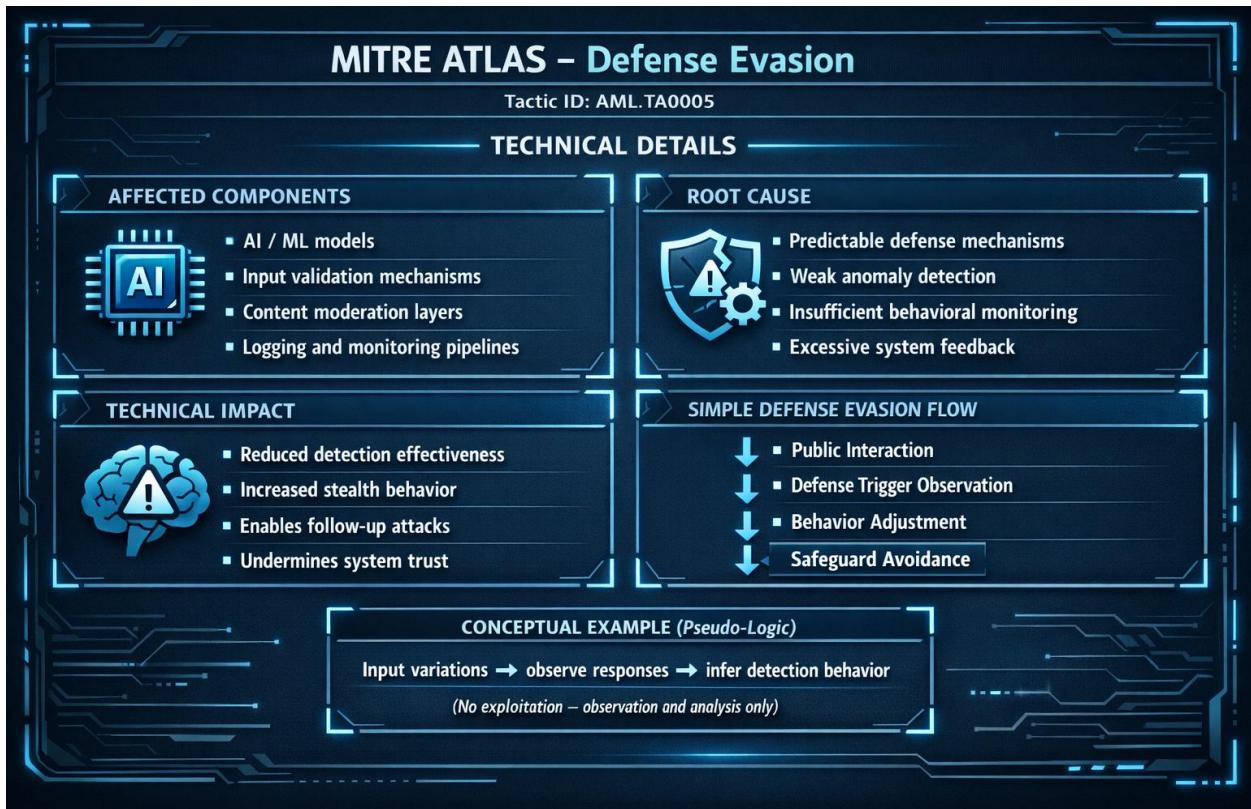


8. Defense Evasion

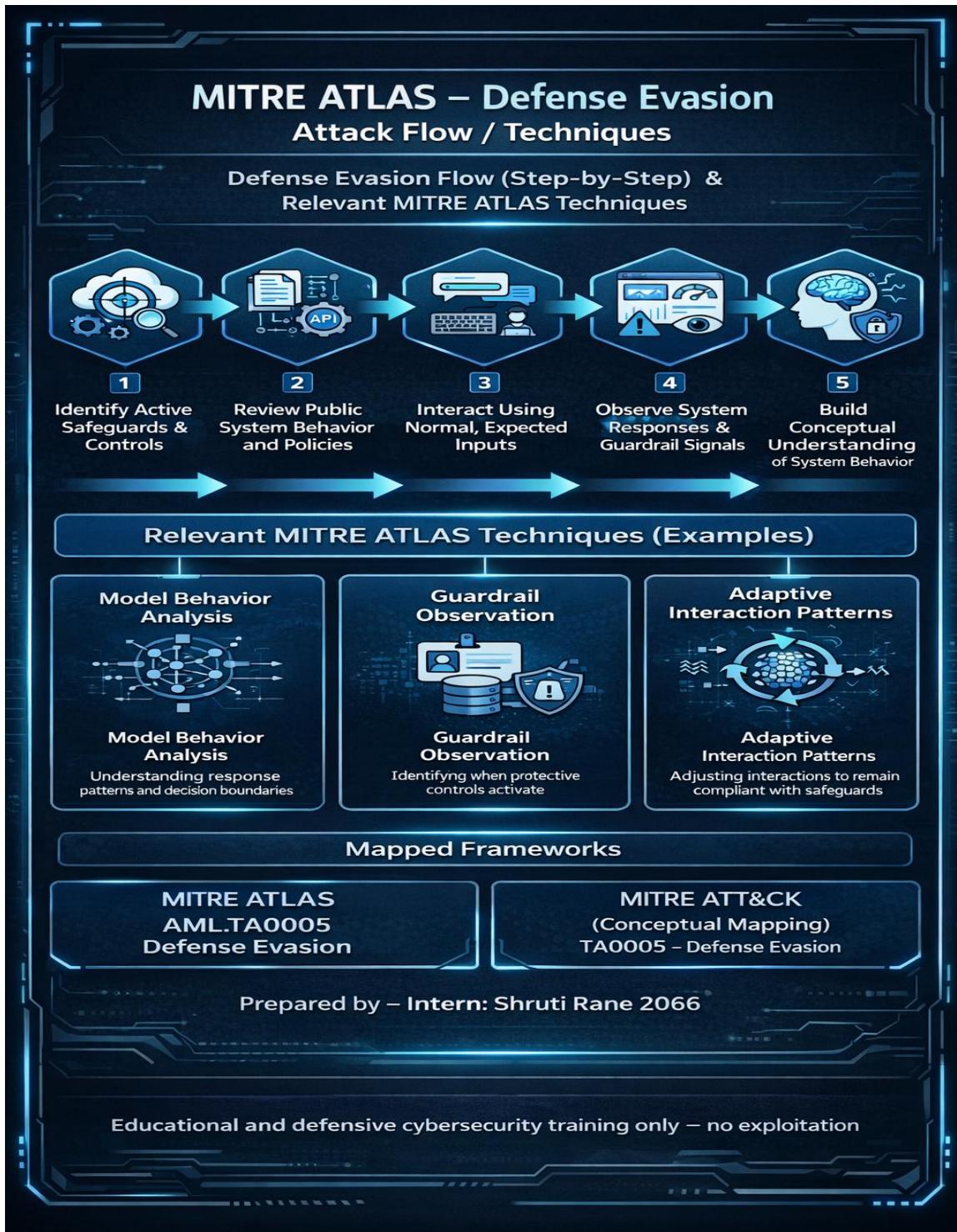
Overview



Technical Details:



Attack Flow:

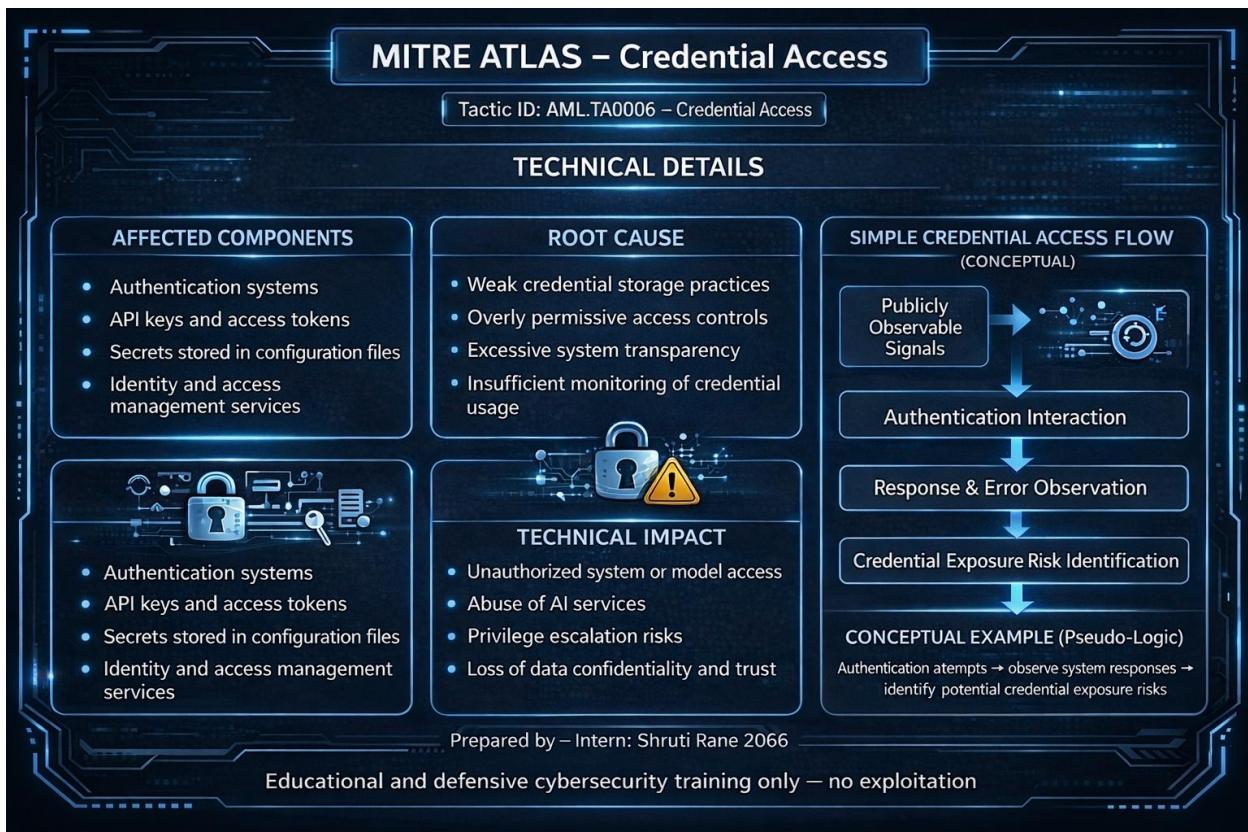


9. Credential Access

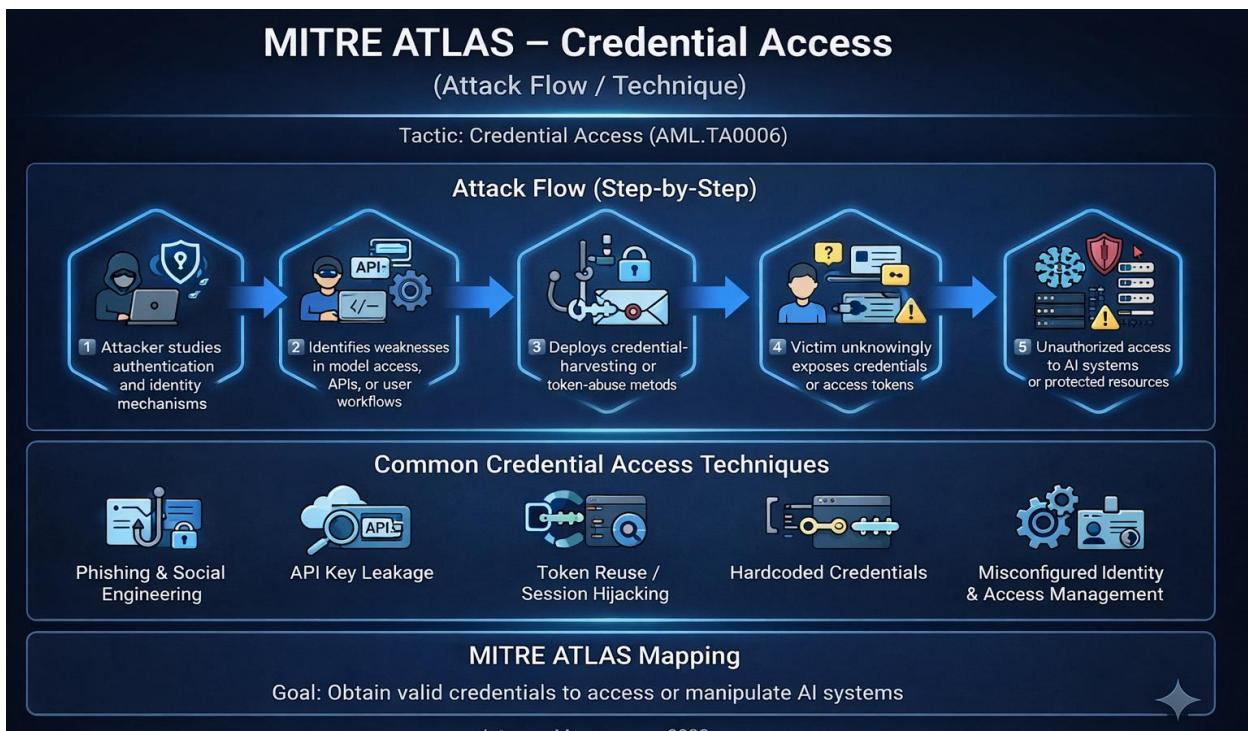
Overview



Technical Details:

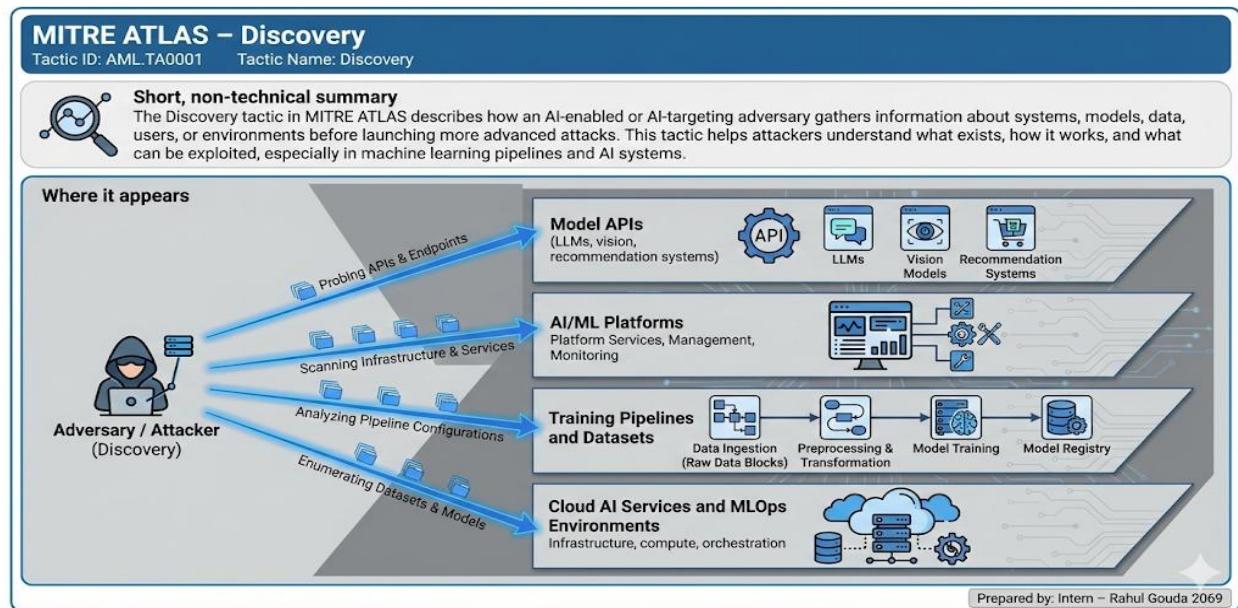


Attack Flow

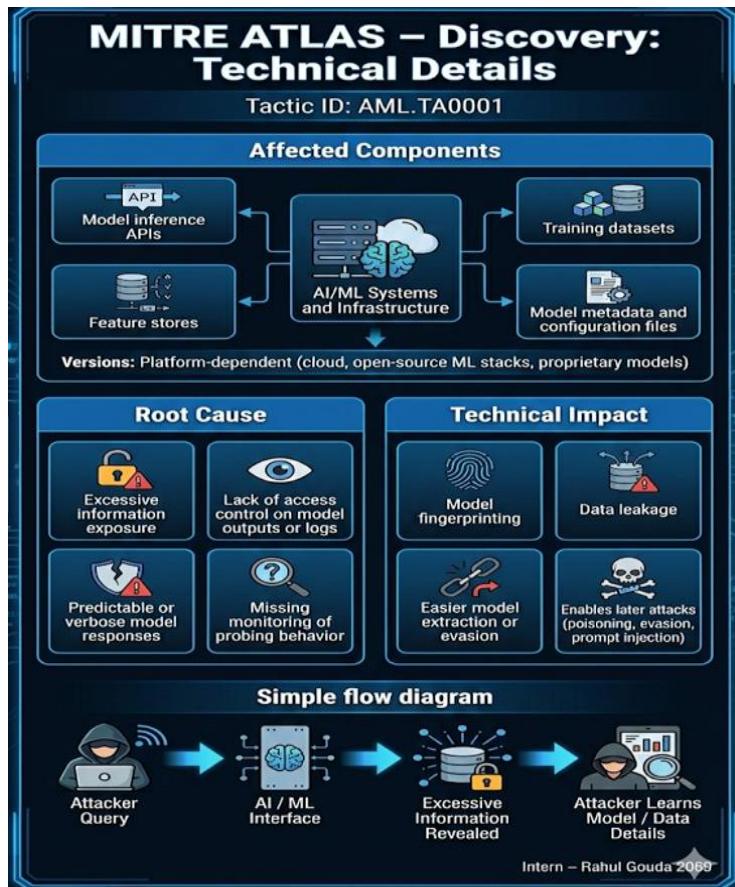


10. Discovery

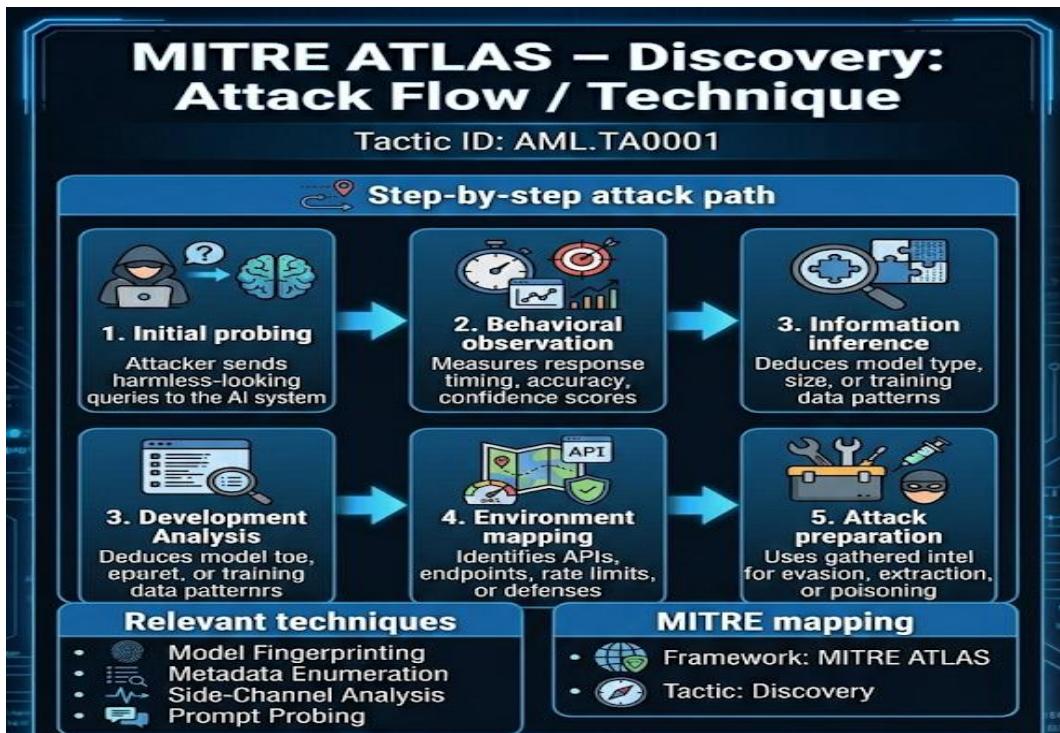
Overview



Technical Details:

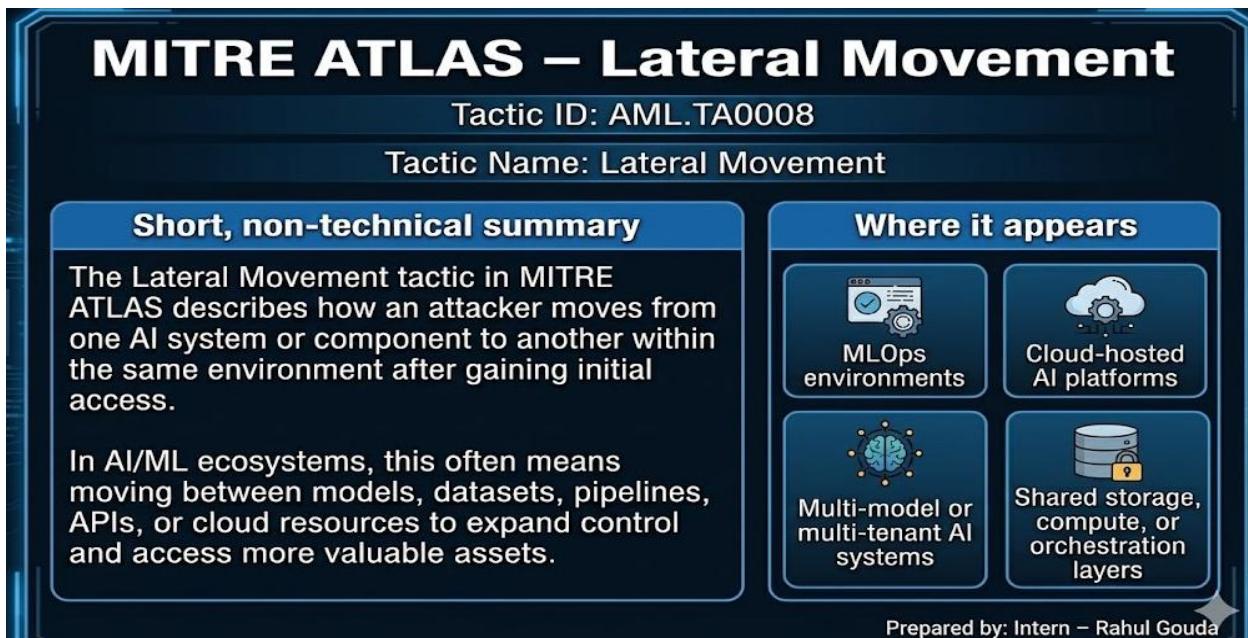


Attack Flow:

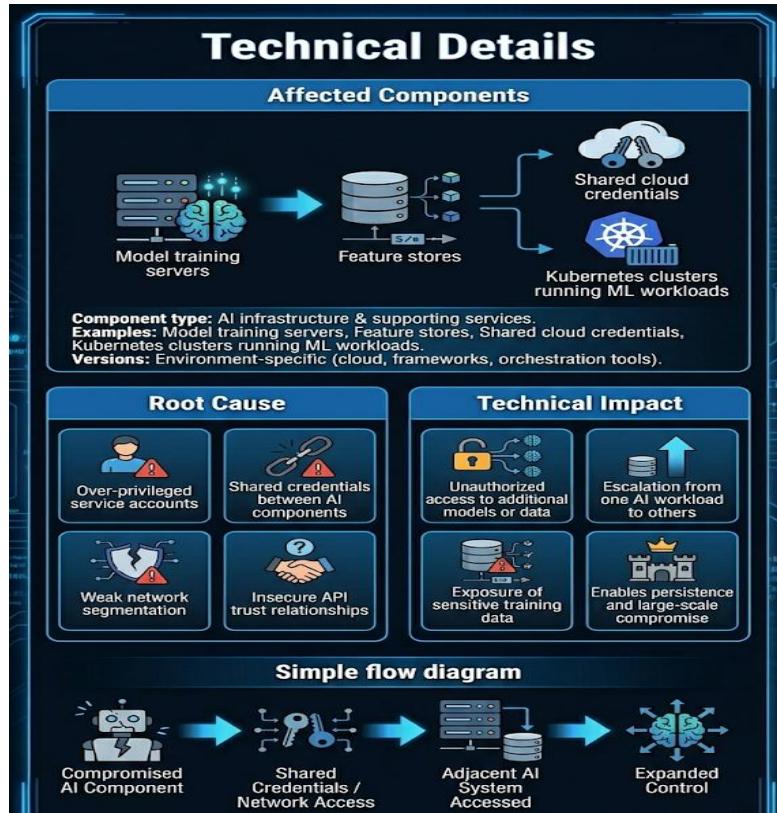


11. Lateral Movement

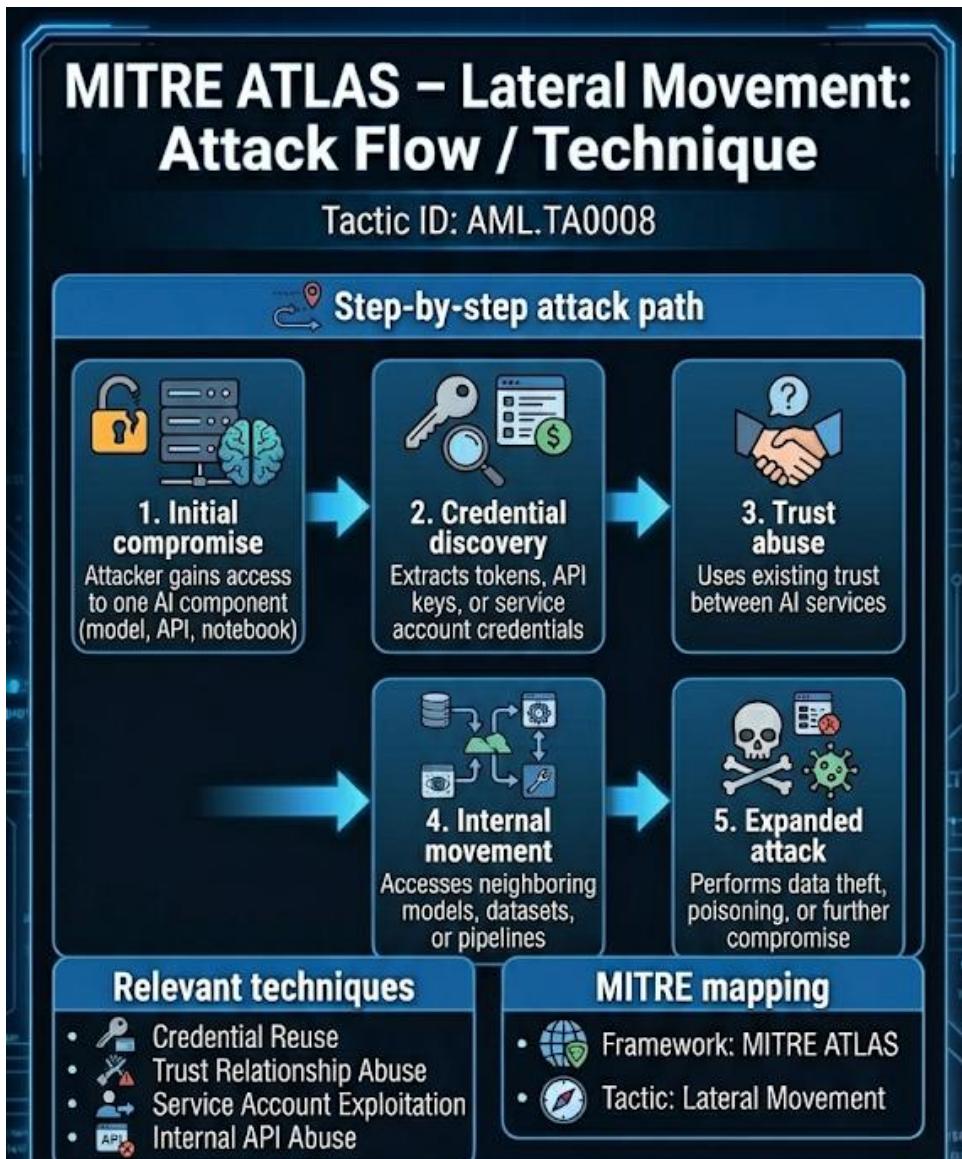
Overview



Technical Details:

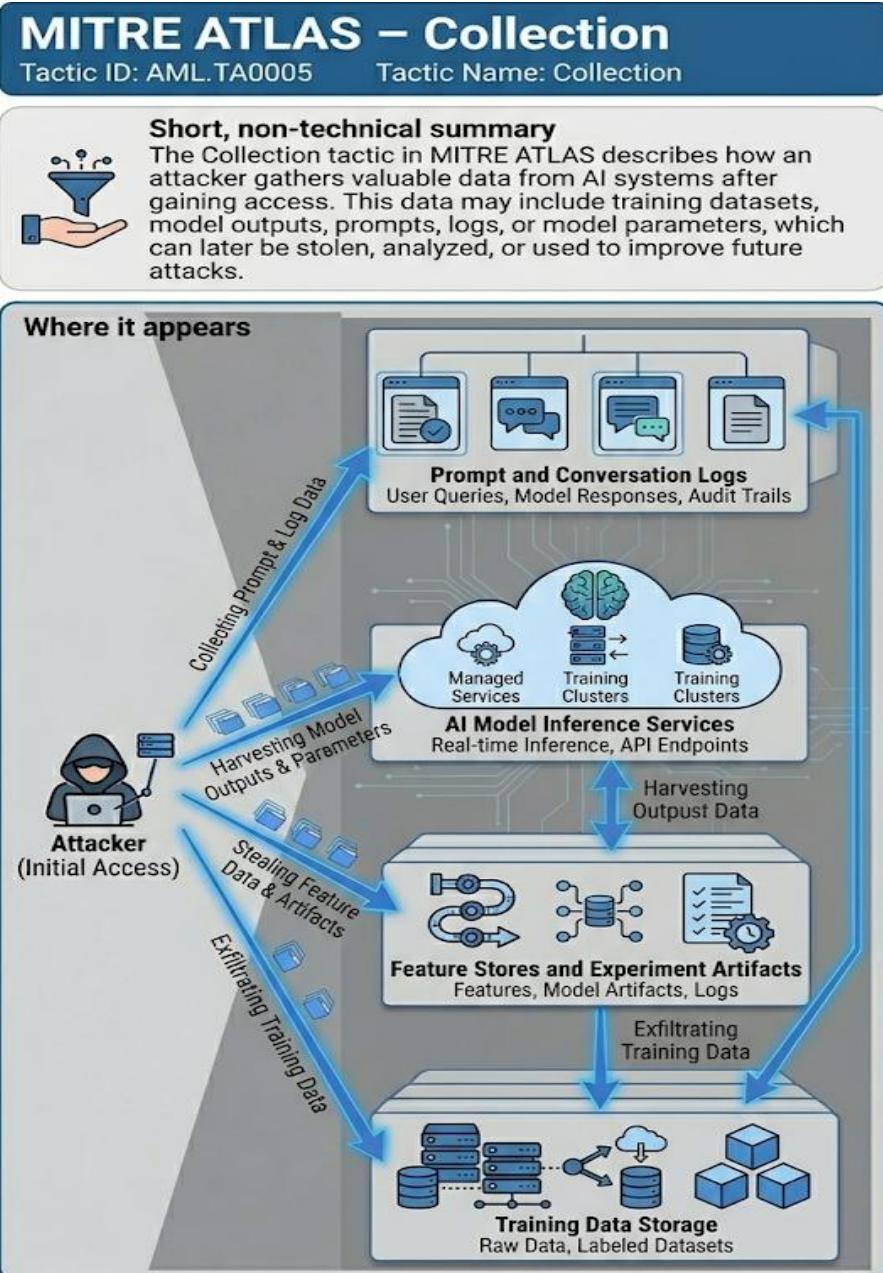


Attack Flow:

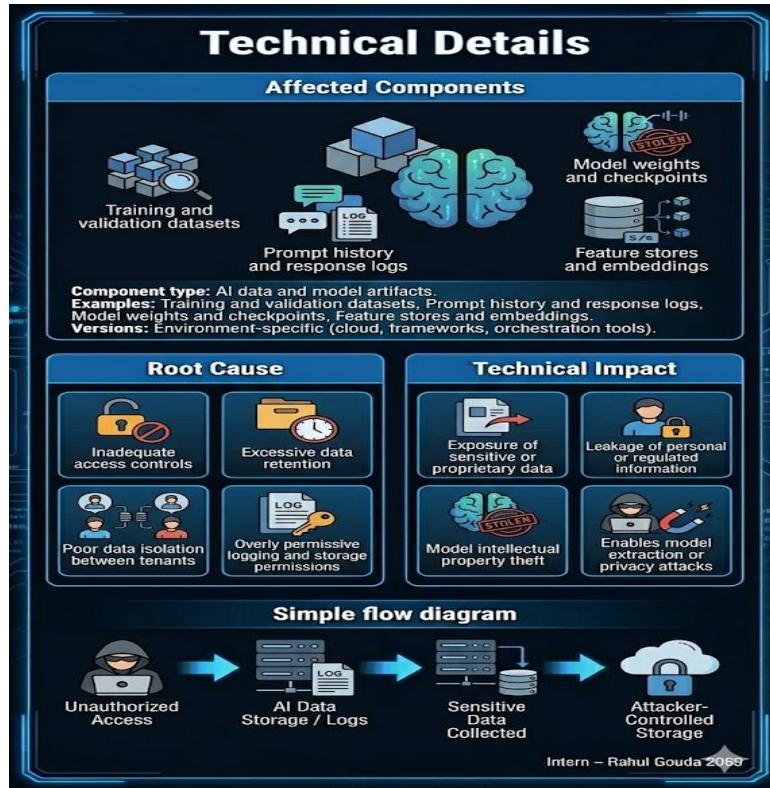


12. Collection

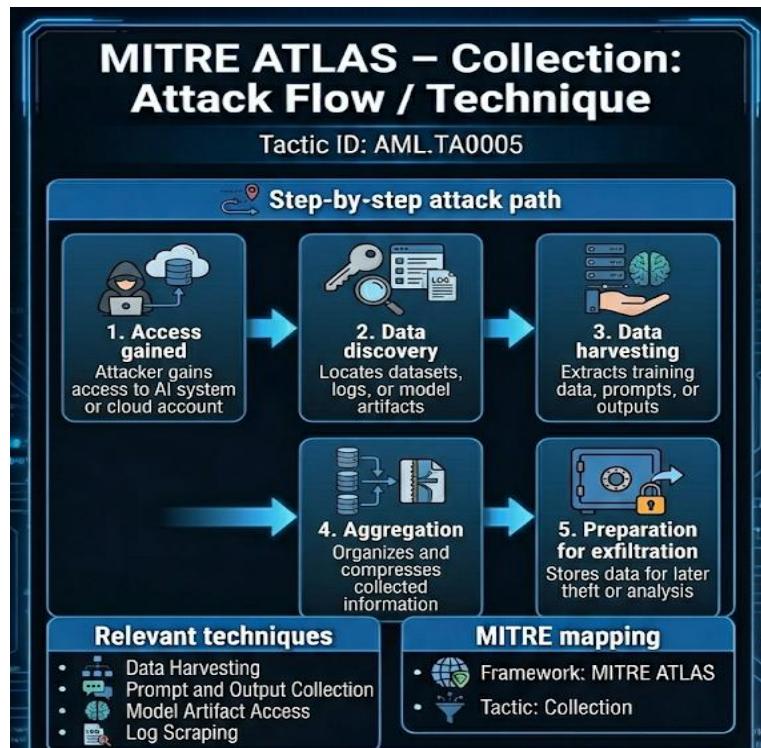
Overview



Technical Details:



Attack Flow:



13. Command and Control

Overview

MITRE ATLAS
COMMAND AND CONTROL

Tactic ID: AML.TA0011
Tactic Name: Command and control

Short Non-Technical Summary

Command and control is the phase where an attacker prepares tools, infrastructure, data, or models before attacking an AI/ML system.

This happens **before the actual attack**, and includes:

- Creating malicious datasets
- Building attack scripts or ML models
- Setting up servers, domains, or APIs

Think of it as the **preparation stage** of an AI attack.



• Creating malicious datasets	Building attack scripts or ML models	Seryer & Domain Setup
-------------------------------	--------------------------------------	-----------------------

Think of it as the **preparation stage** of an AI attack.

Prepared by: Intern – Sanika Jankar

Technical detail:

MITRE ATLAS

TECHNICAL DETAILS

COMMAND AND CONTROL

Affected Components



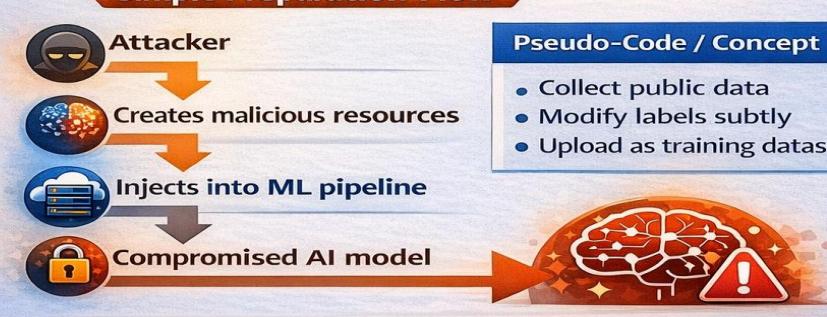
- AI/ML pipelines
- Training datasets
- Model repositories
- Cloud ML platforms
- APIs used for model access

Root Cause



- Publicly accessible **training data**
- Lack of dataset integrity validation
- Untrusted third-party models or libraries
- Weak supply-chain controls

Simple Preparation Flow



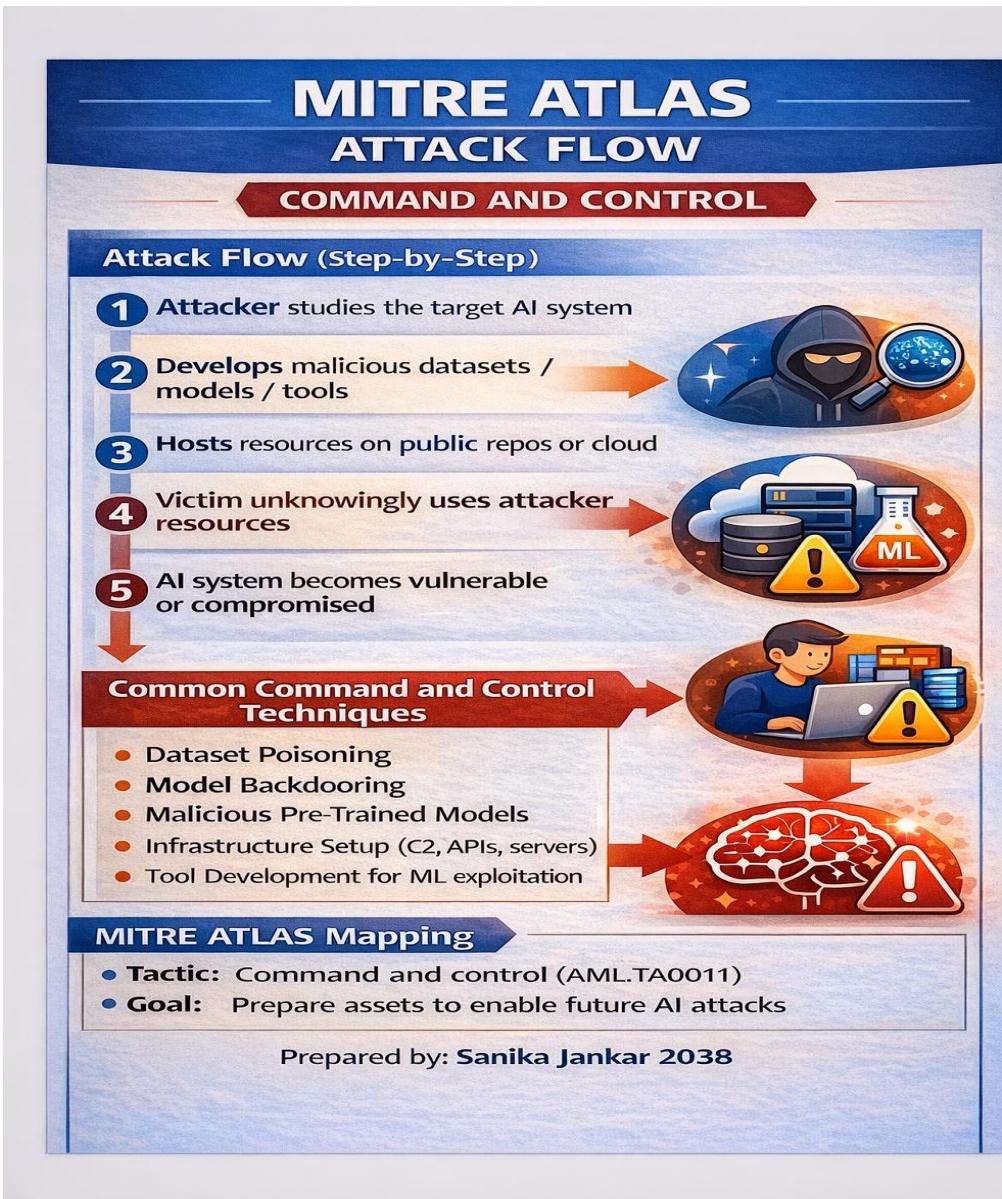
```
graph TD; Attacker((Attacker)) --> Resources((Creates malicious resources)); Resources --> Pipeline((Injects into ML pipeline)); Pipeline --> Model((Compromised AI model));
```

Pseudo-Code / Concept

- Collect public data
- Modify labels subtly
- Upload as training dataset

Prepared by: Intern – Sanika Jankar

Attack Flow:



14. Exfiltration

Overview

MITRE ATLAS
Exfiltration

Tactic ID: AML.TA0010
Tactic Name: Exfiltration

Short Non-Technical Summary

Exfiltration is the phase where an attacker prepares tools, infrastructure, data, or **models** before attacking an AI/ML system. This happens **before the actual attack**, and includes:

- Stealing model training data
- Extracting proprietary ML models
- Exporting sensitive datasets via APIs

Think of it as the **preparation stage** of an AI attack.



Think of it as the **preparation stage** of an AI attack.

Q: Short Non-Technical Summary

• Exfiltration	• Model Backdooring
• Stealing model training data	• Exporting sensitive datasets via APIs

Prepared by: Sanika Jankar 2038

Technical detail:

MITRE ATLAS
TECHNICAL DETAILS

EXFILTRATION
Tactic ID: AML.TA0010

Affected Components



- AI/ML pipelines
- Training datasets
- Model repositories
- Cloud ML platforms
- APIs used for model access

Root Cause

- Publicly accessible training data
- Lack of dataset integrity validation
- Untrusted third-party models or libraries
- Weak supply-chain controls



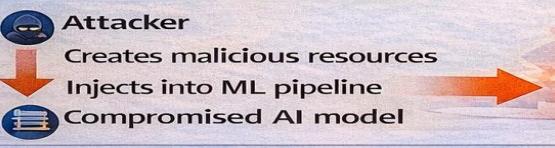
Technical Impact

- 🎯 Poisoned ML models
- 📉 Reduced model accuracy
- 🔒 Backdoors embedded in AI systems
- 👤 Future exploitation made easier

Pseudo-Code / Concept

- Collect public data
- Modify labels subtly
- Upload as training dataset

Simple Preparation Flow



- Attacker
Creates malicious resources
- Injects into ML pipeline
- Compromised AI model



Prepared by: Intern – Sanika Jankar 2038

Attack flow:



15. Impact

Overview

MITRE ATLAS – Impact
Tactic ID: AML.TA0040

Short Non-Technical Summary
the attacker phase where attacker tools tools infrastructure data or before attacking AI/ML system.

This happens before actual attack, and includes:

- Creating malicious datasets
- Building attack scripts or ML models
- Building attack scripts or ML models
- Infrastructure domains, or APIs



Prepared By: Intern – Sanika Jankar 2038

Think of ais as preparation stage of an AI attack.



Technical detail:

MITRE ATLAS – Impact

Tactic ID: AML.TA0040

Affected Component

AI/ML Compelines

- AI/ML pipelines training data
- Building attack scripts or ML models
- Model repositories
- Lack of dataset integrity validation
- APIs used of model access

Root Cause

- Publicly accessible training data
- Lack of dataset integrity validation
- Untrusted third-party models or libraries
- Weak supply-chain controls

Prepared By: Intern – Sanika Jankar 2038

Technical Impact

- Poisoned ML models
- Reduced model accuracy
- Backdoors embedded AI systems
- Future exploitation made easier

Simple Preparation Flow

Attacker

- ↓
Creates public data
Injects into ML pipeline
Compromised AI model

Pseudo-Code / Concept

Modify labels subtly
Upload as training dataset

Attack flow:

MITRE ATLAS: Impact (AML.TA0040) - Attack Flow



Common Impact Techniques

- Dataset Poisoning
- Model Backdooring
- Malicious Pre-Trained Models
- C2 Infrastructure

MITRE ATLAS Mapping

Tactic: Impact (AML.TA0040)

Goal: Strategic preparation of assets to ensure the success a future AI attack

Prepared by: Intern – Sanika Jankar 2038

