

Microsoft Al Tour

In partnership with **NIVIDIA**.







Securing Generative Al Applications

Rod Trent Senior Program Manager, Microsoft



Agenda



Introduction

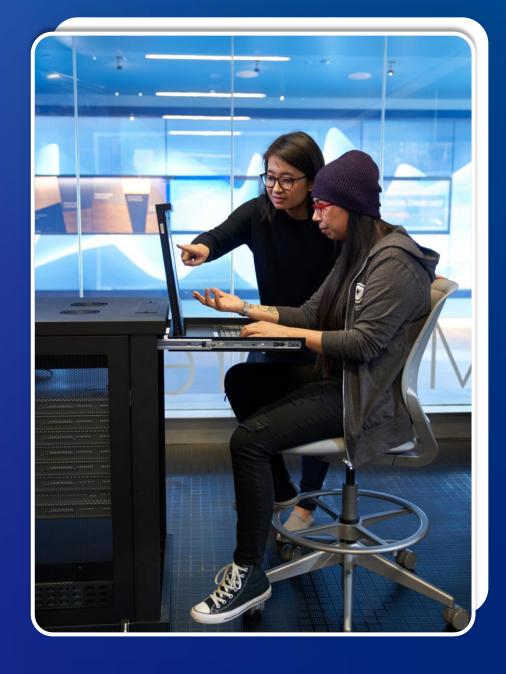
Understanding AI architecture in a security context

The AI security threat landscape

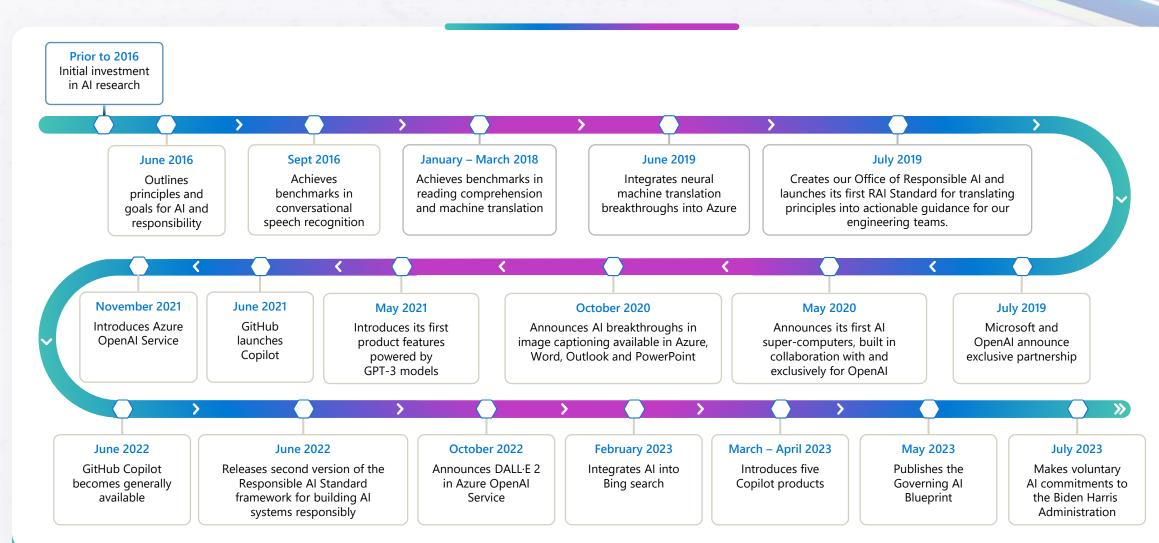
How Microsoft secures AI platforms

Security controls for developers building AI-enabled applications

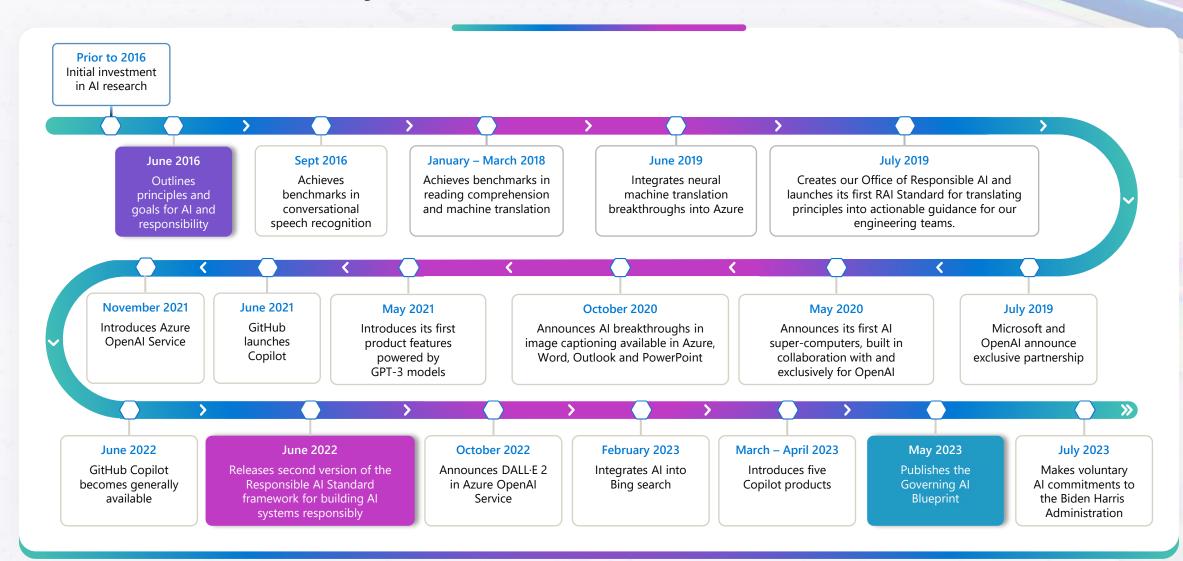
Introduction



Microsoft's Journey in Al



Microsoft's Journey in Al



Microsoft's Responsible Al Framework













>>

Fairness

Al systems should treat all people fairly.

Reliability & Safety

Al systems should perform reliably and safely.

Privacy & Security

Al systems should be secure and respect privacy.

Inclusiveness

Al systems should empower everyone and engage people. Transparency

Al systems should be understandable.

Accountability

People should be accountable for Al systems.

Learn More

https://www.microsoft.com/en-us/ai/responsible-ai

Microsoft's Responsible Al Framework













>>

Fairness

Al systems should treat all people fairly.

Reliability & Safety

Al systems should perform reliably and safely.

Privacy & Security

Al systems should be secure and respect privacy.

Inclusiveness

Al systems should empower everyone and engage people.

Transparency

Al systems should be understandable.

Accountability

People should be accountable for Al systems.

Learn More

https://www.microsoft.com/en-us/ai/responsible-ai

Understanding Al architecture in a security context



Al Architecture Overview

AI Usage This layer focuses on the user interaction with the AI interface, either standalone or built into existing application UI.

Current Examples:

- Public Access (Bing Chat/ChatGPT/Bard)
- Code Development (GitHub Copilot)
- Microsoft Copilots (Office, Viva, Windows)

AI Application Development of an Al-integrated application through secured software development practices (SDL). Additional plugins are enabled to provide specific functions and controls.

Current Examples:

- Microsoft development of Copilots
- Customers developing their own products
- 3rd party solution offerings running on Azure

AI Platform Access to traditional and generative Al models to build your own Al-integrated solutions, running in your compliance boundary on secure and reliable cloud infrastructure.

Current Examples:

- Azure ML Model Catalogue (Hugging Face, Llama 2)
- Azure OpenAl Service (GPT4, ChatGPT, DALL-E)

Detailed AI architecture breakdown

AI Usage	User Prompt	Prompt Response			
AI Application	User Interface	Content Filter	Grounding	Semantics	Plugins
	Orchestration	Prompt Engineering	Memory	Audit/Controls	
AI Platform	API	Orchestration	Plugin Management	Deep Safety System	Generative Al Model
	Al Infrastructure	Audit/Controls			

Al Shared Responsibility Model

Illustrates which responsibilities are typically performed by an organization and application developer and which are performed by their Al provider (such as Microsoft)

IaaS PaaS SaaS

		(BYO Model)	(Azure AI)	(Copilot)
AI Usage	User training and accountability			
	Usage policy, admin controls			
	Identity, device, and access management			
	Data governance			
AI Application	Al plugins and data connections			
	Application design and implementation			
	Application infrastructure			
	Application safety systems			
AI Platform	Model safety and security systems			
	Model accountability			
	Model tuning			
	Model design and implementation			
	Model Training Data Governance			
	Al Compute Infrastructure			

Al Pain Points











Integration

New technologies and design decisions introduce new risks and vulnerabilities.

User training needs to be adapted to the new capabilities of the Al solutions selected for use by the organization.

Data & Privacy

Sensitive data access and processing via Al systems creates new risks.

Transparency and control needs to be established and maintained through out the lifecycle.

Al Supply Chain

Increased focus on potentially vulnerable or malicious code or 3rd party components.

Lack of compliance standards and rapidly developing best practices.

Trusted AI

Very similar to the early days of BYOD:
Employees likely already using GenAl to achieve their tasks.

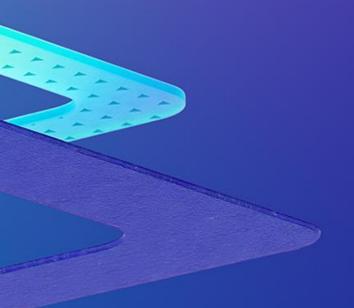
Leaders must establish a trusted pathway to GenAl integrated applications to protect the organization.

The Unknowns

GenAl is new and brings unique challenges such as Al Hallucinations.

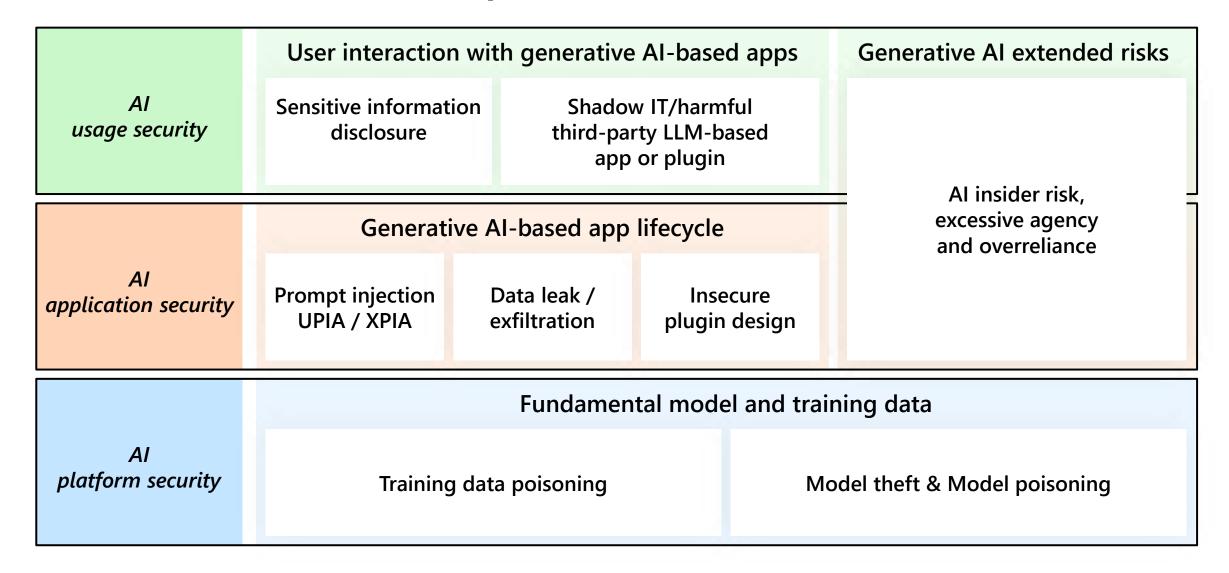
Exciting prospects of the potential, the ROI is not yet proven in real-world scenarios.

The Al security threat landscape





Generative AI threat map



The cybersecurity bell curve

Basic security hygiene still protects against 98% of attacks¹



Enable multifactor authentication

Make it harder for bad actors to utilize stolen or phished credentials by enabling multifactor authentication.
Always authenticate and authorize based on all available data points, including user identity, location, device health, service or workload, data classification, and anomalies.

Apply least privilege access

Prevent attackers from spreading across the network by applying least privilege access principles, which limits user access with just-in-time and just-enough-access (JIT/JEA), risk-based adaptive polices, and data protection to help secure both data and productivity.

Keep up to date

Mitigate the risk of software vulnerabilities by ensuring your organization's devices, infrastructure, and applications are kept up to date and correctly configured. Endpoint management solutions allow policies to be pushed to machines for correct configuration and ensure systems are running the latest versions.

Utilize antimalware

Stop malware attacks from executing by installing and enabling antimalware solutions on endpoints and devices.

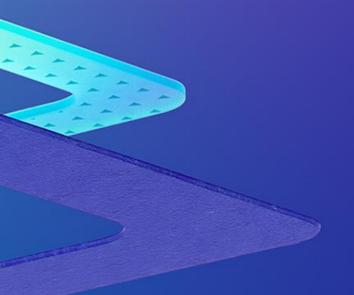
Utilize cloud-connected antimalware services for the most current and accurate detection capabilities.

Know where your sensitive data is stored and who has access. Implement information protection best practices such as applying sensitivity labels and data loss prevention policies. If a breach does occur, it's critical that security teams know where the most sensitive

data is stored and accessed.

Protect data

How Microsoft secures Al platforms





Building new principles for AI security

Ensure your data is *your* data



Microsoft will not use customers data to train the foundational Al models, without explicit consent. Data governance is a shared responsibility.

Customer:

Protect your data as a top priority. Ensure it remains private and controlled, end to end.

Be secure by design



The Microsoft AI stack is designed and built on decades of secure software practices, and a strong supply chain of partners, following mature SDL tools and processes.

Customer:

Ask for transparency in every Al system you connect to your data, for the whole Al supply chain.

Secure by intention and in practice

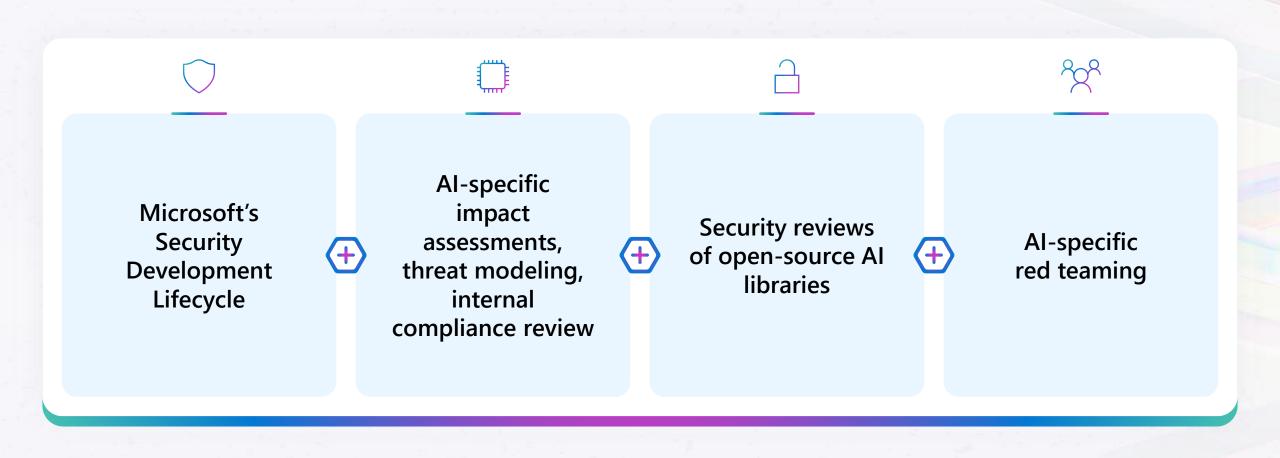


Microsoft grounds its efforts to advance AI in strong principles that govern the technology's implementation and use: Positioning people at the center of the equation.

Customer:

Strong Zero Trust, and Data Governance programs will matter more than ever.

How does Microsoft address risks from attackers?



How does Microsoft address risks from the use or misuse of AI?



Example: Overreliance

Using AI to justify a viewpoint or action

Assuming the AI must fair or accurate

Al doing something that the user can't meaningfully check.

User is simply too busy to check it carefully



Microsoft Approach

Ground on authoritative data sources

Provide greater transparency and explainability

Design UX interfaces to mitigate overreliance

Bing is powered by AI, so surprises and mistakes are possible.

How does Microsoft address risks from the use or misuse of AI?



Example: Hostile Misuse

Hostile misuse involves using AI system to intentionally to cause harm including circumventing safeguards.

Generating malicious code

Asking instructions for harmful purposes

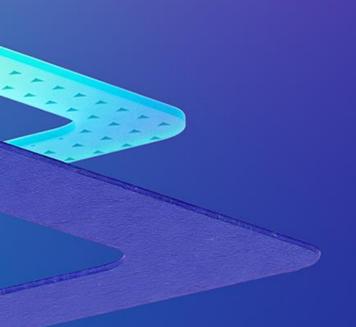


Microsoft Approach

Microsoft has invested heavily in a defense in depth approach including a deep safety layer that provides security by default to disallow the AI to perform tasks that are harmful or dangerous to the user, intentionally deceptive, or likely to adversely affect the public interest

Acceptable Use Policy governs Al usage

Microsoft's Al Red Teaming Approach





Security Community view of Red Teaming



Double Blind

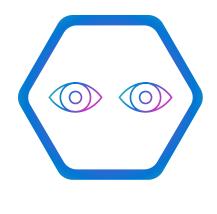


Emulate real world adversaries



Mature toolkit and processes

RAI Community view of Red Teaming





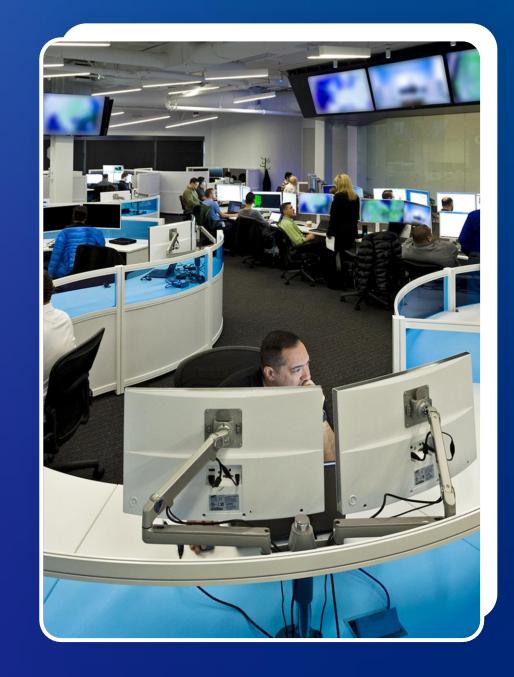


Adversarial and Benign

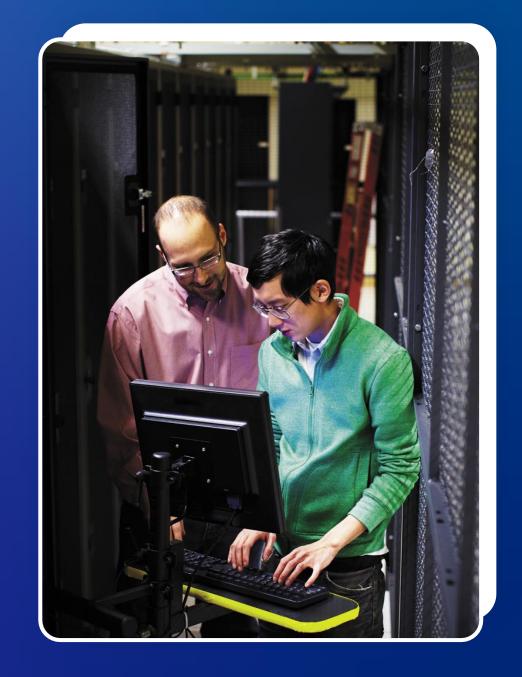


Rapidly Evolving
Tools and processes

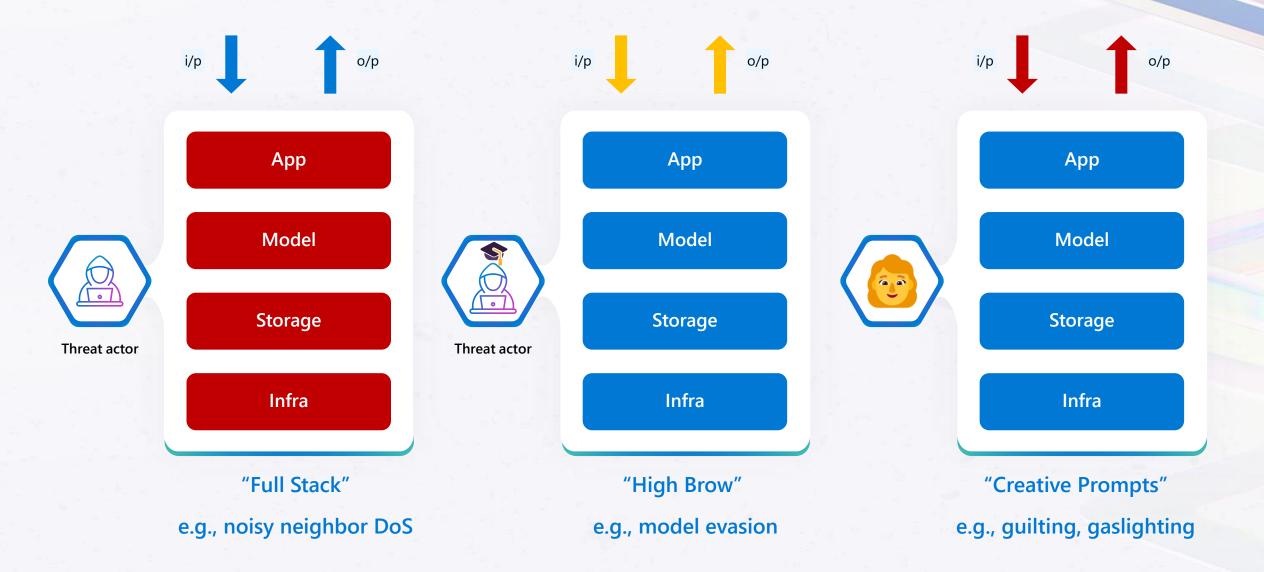
Al Red Teaming combines the best of both worlds



Al Red Team = Probing for Security + Responsible Al Harms



Three Flavors of AI Red Teaming



Three Flavors of AI Red Teaming



Focusing on the entire AI stack

Leveraging Traditional Security skills



High Brow

Focus only on the i/p and o/p

Leveraging Adversarial ML skills



Creative Prompt

Focuses on the i/p and o/p

Leverages a broad skillset to cause failures



Al-specific red teaming hardens the effectiveness of security protections



Expands the definition and scope for Al



Focuses on failures from both malicious and benign personas

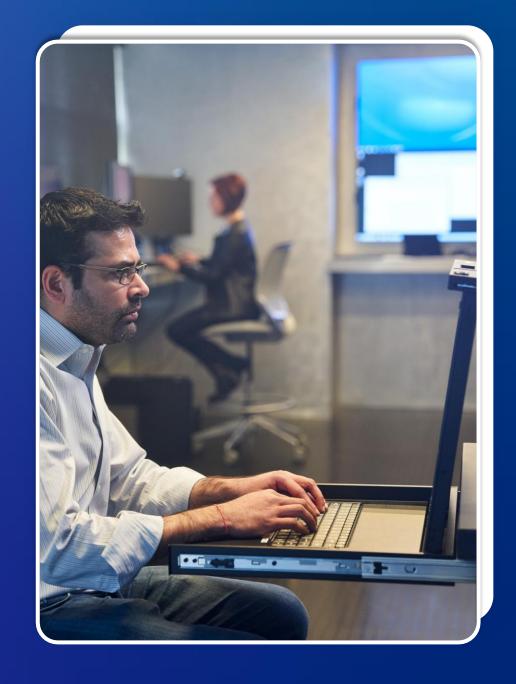


Recognizes that
Al systems are
constantly evolving

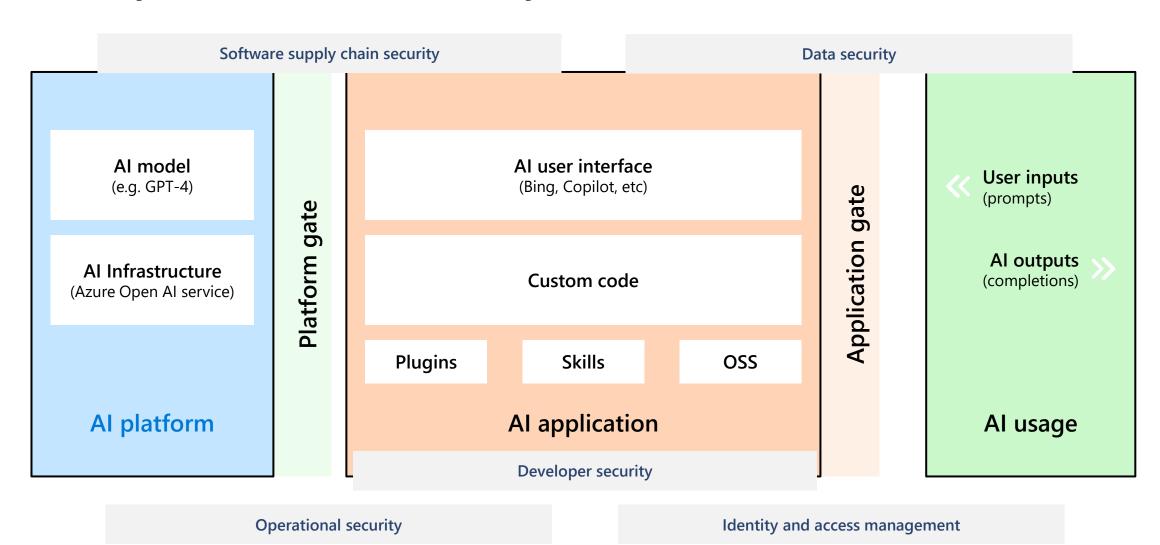
Learn More

https://www.microsoft.com/en-us/security/blog/2023/08/07/microsoft-ai-red-team-building-future-of-safer-ai/

Security controls for developers building Al-enabled applications



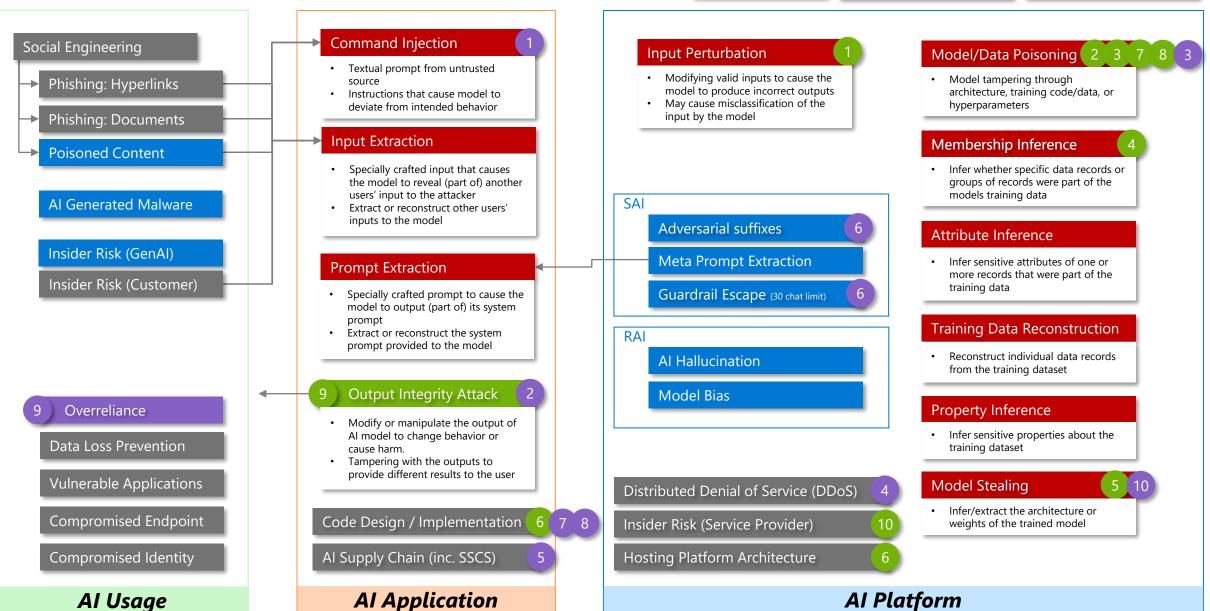
Security controls within AI systems



Threat Modelling Scenarios

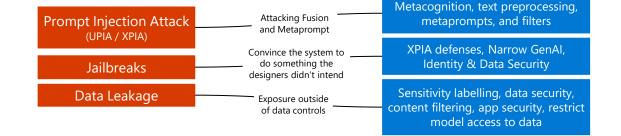
MSRC AI Bug Bar OWASP Top 10 for ML Generative AI Specific

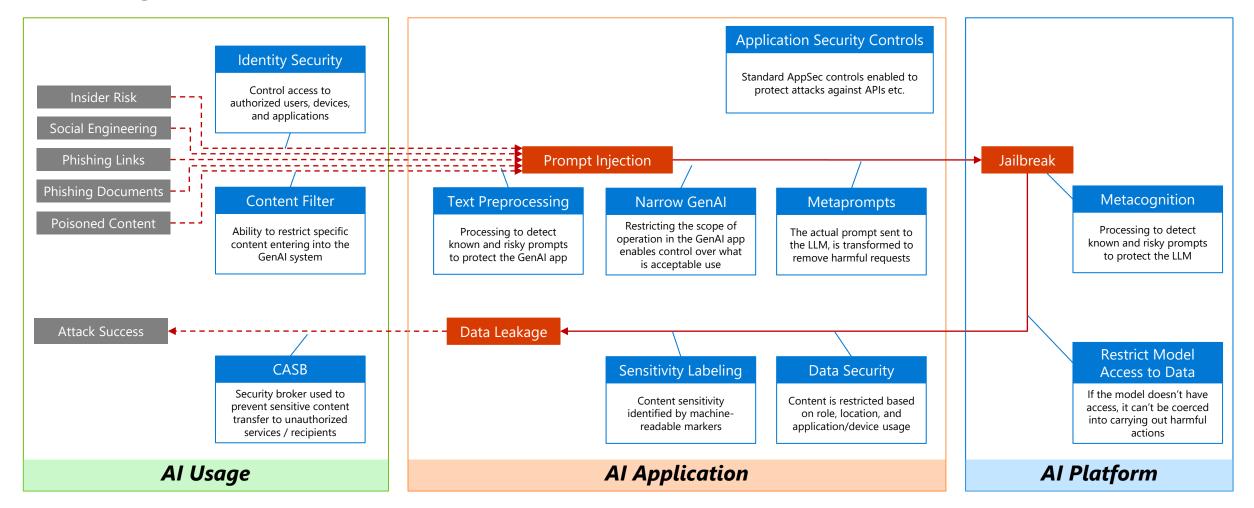
MITRE ATLAS OWASP Top 10 for LLM Common Cyber Threats



Threat Mapping Template

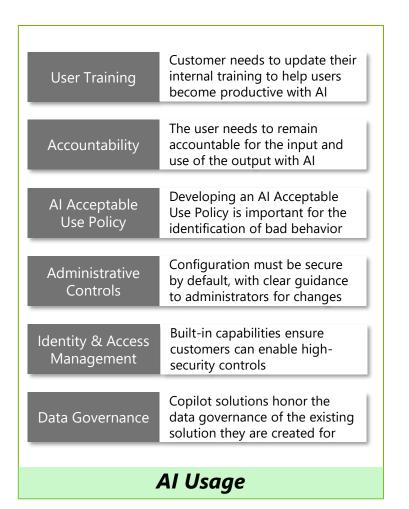
This framework provides a repeatable method of articulating both the vulnerabilities (red) and the mitigations (blue)

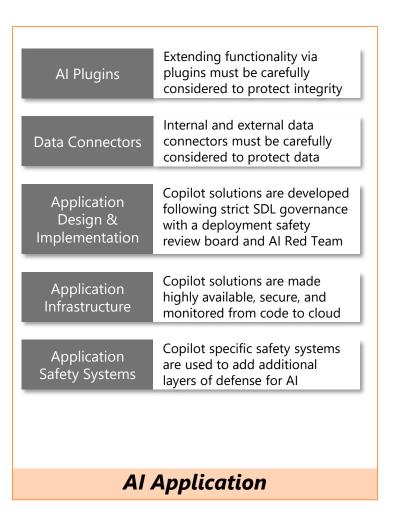


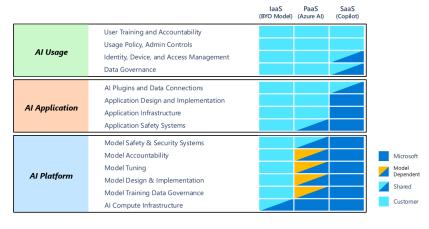


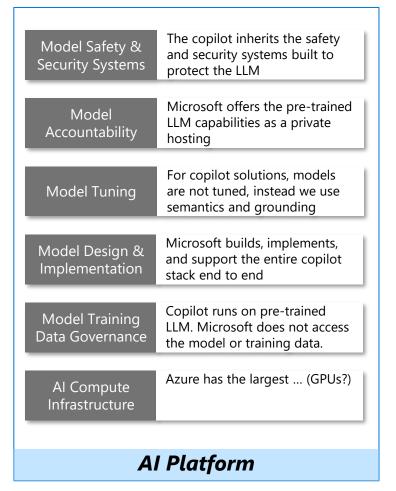
End-to-End Secure AI

SaaS (Copilot)

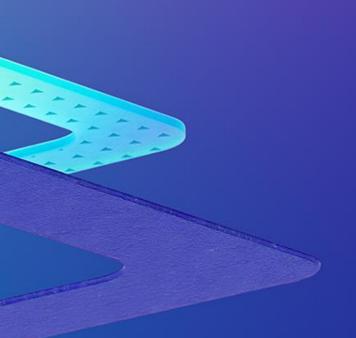


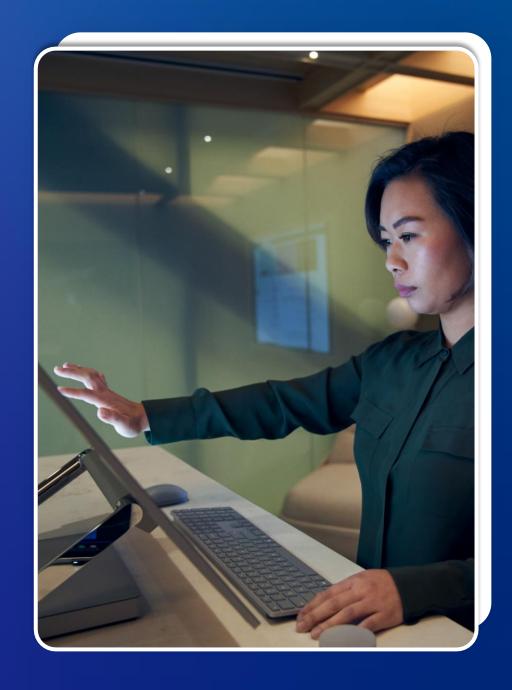






Wrap up





Security for AI is an ever-evolving process









Develop

Continually evolving secure by design and by default requirements

Test

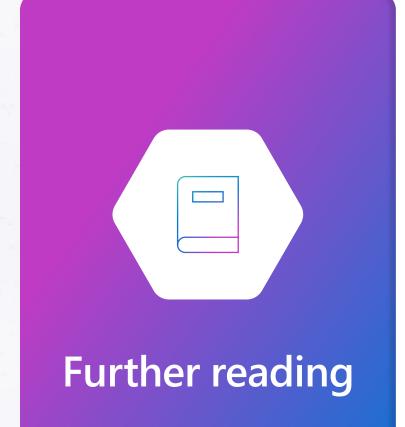
Continual assessment and AI-specific red teaming

Monitor/Respond

Analyze 65T+ threat signals, respond to incidents

Evolve

Learn, share, collaborate





Microsoft Security Copilot documentation | Microsoft Learn



Al shared responsibility model – Microsoft Azure | Microsoft Learn



Best practices for AI security risk management | Microsoft Security Blog



Microsoft Cybersecurity Reference Architectures
(MCRA) – Security documentation | Microsoft Learn



