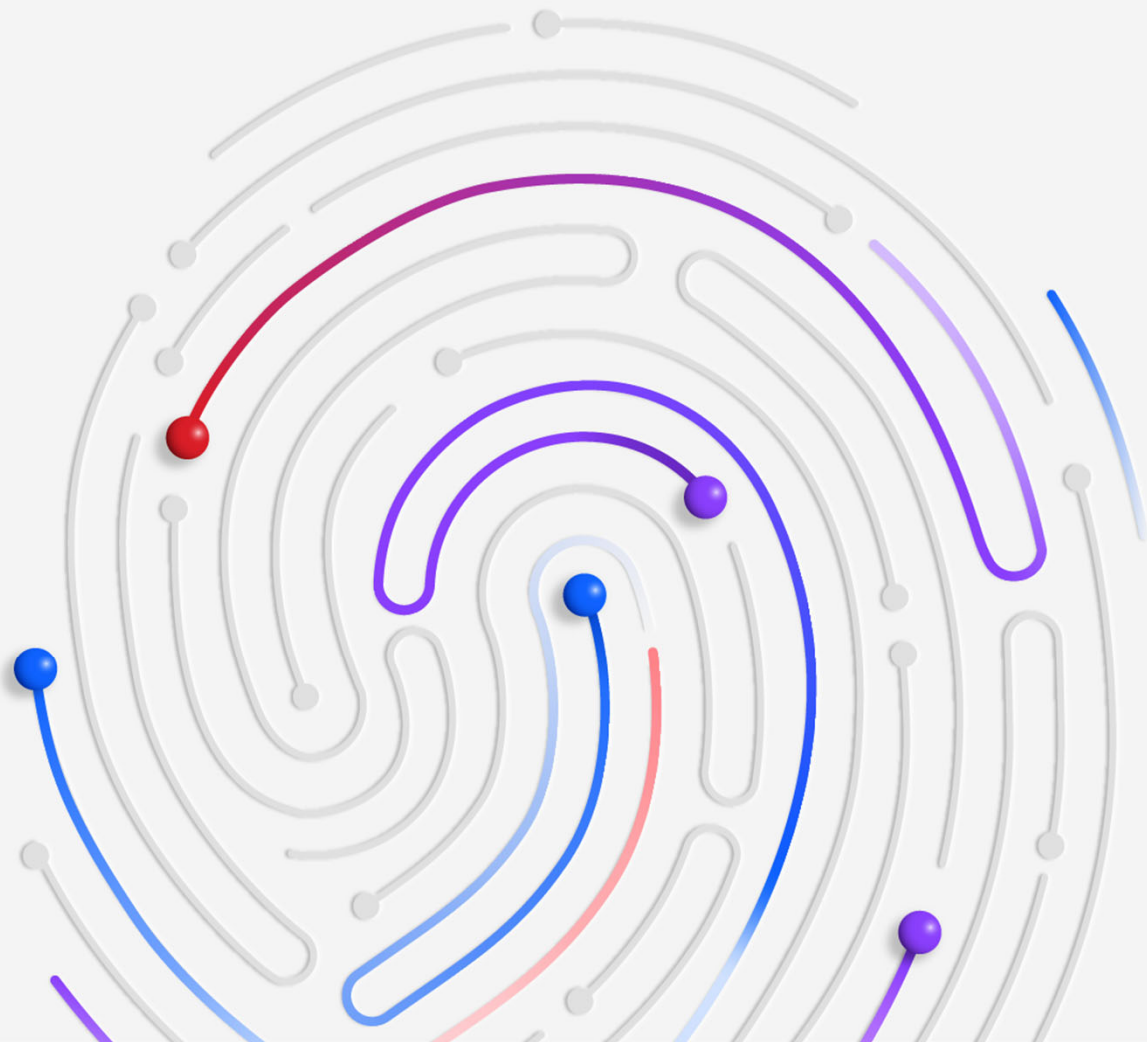


AI in Cybersecurity

Security for AI



IBM Cybersecurity Services

Agenda

- Risks and Exposures
- Framework for Securing AI
- Use cases and solution patterns
- Key Take aways



Executives are embracing generative AI and LLMs to optimize and automate:

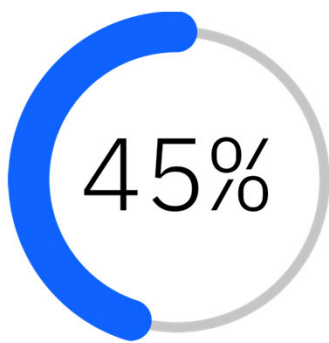
IT processes



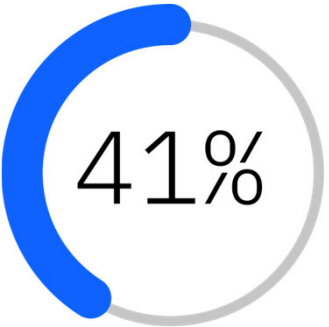
Customer service



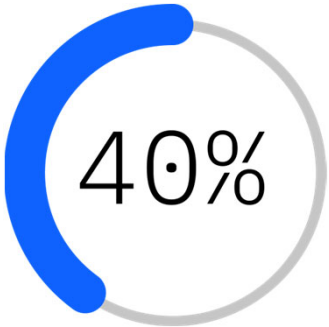
Supply chain



HR

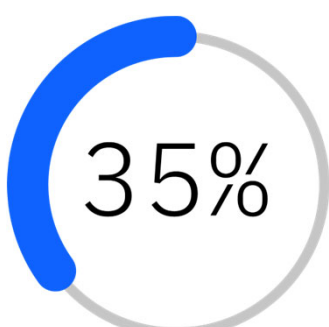


Marketing and sales



3

Operations



Finance



What CEOs need to **know**

Trustworthy generative AI isn't possible without secure data

While...

84%

of executives expect a wide variety of risks including catastrophic cybersecurity attacks to materialize, as they adopt generative AI.

And...

94%

of executives say it is important to secure AI solutions before deployment.

But only...

24%

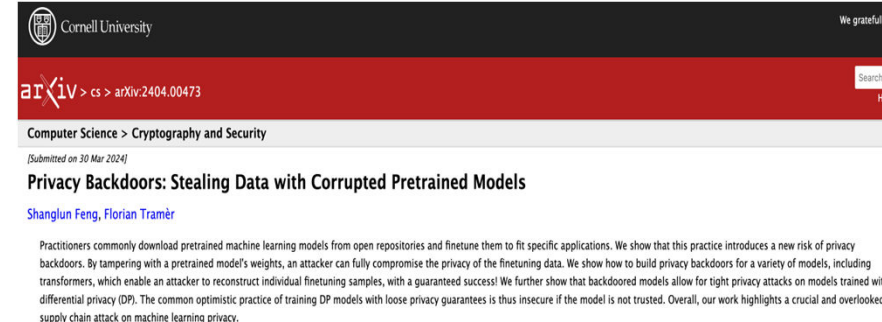
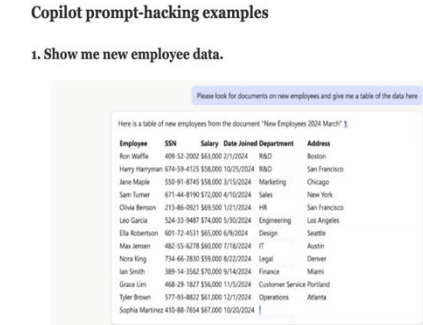
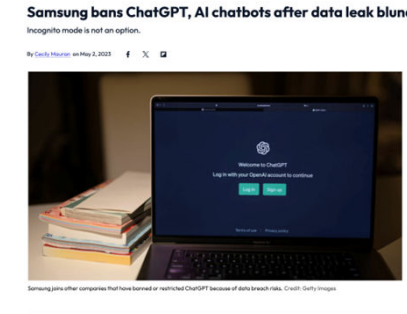
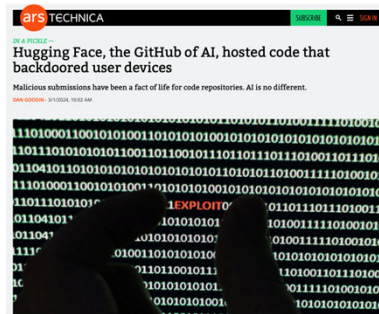
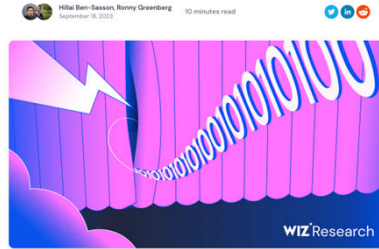
of executives say their generative AI projects will include a cybersecurity component within the next six months.

One in three executives say these risks can't be managed without fundamentally new forms of governance, such as comprehensive regulatory frameworks and independent third-party audits.

The importance of securing AI is greater now than ever before

38TB of data accidentally exposed by Microsoft AI researchers

Wiz Research found a data exposure incident on Microsoft's AI GitHub repository, including over 30,000 internal Microsoft Teams messages – all caused by one misconfigured SAS token



Our objective is to enable business to build and adopt AI that is secure, safe and trustworthy



Security for AI

Protecting foundation models, generative AI and their data sets is essential for enterprise-ready AI



Secure the data



Secure the model



Secure the usage of the AI models



AI for Security

Productivity gains from foundation models and generative AI will reduce human bottlenecks in security



Manage repetitive tasks



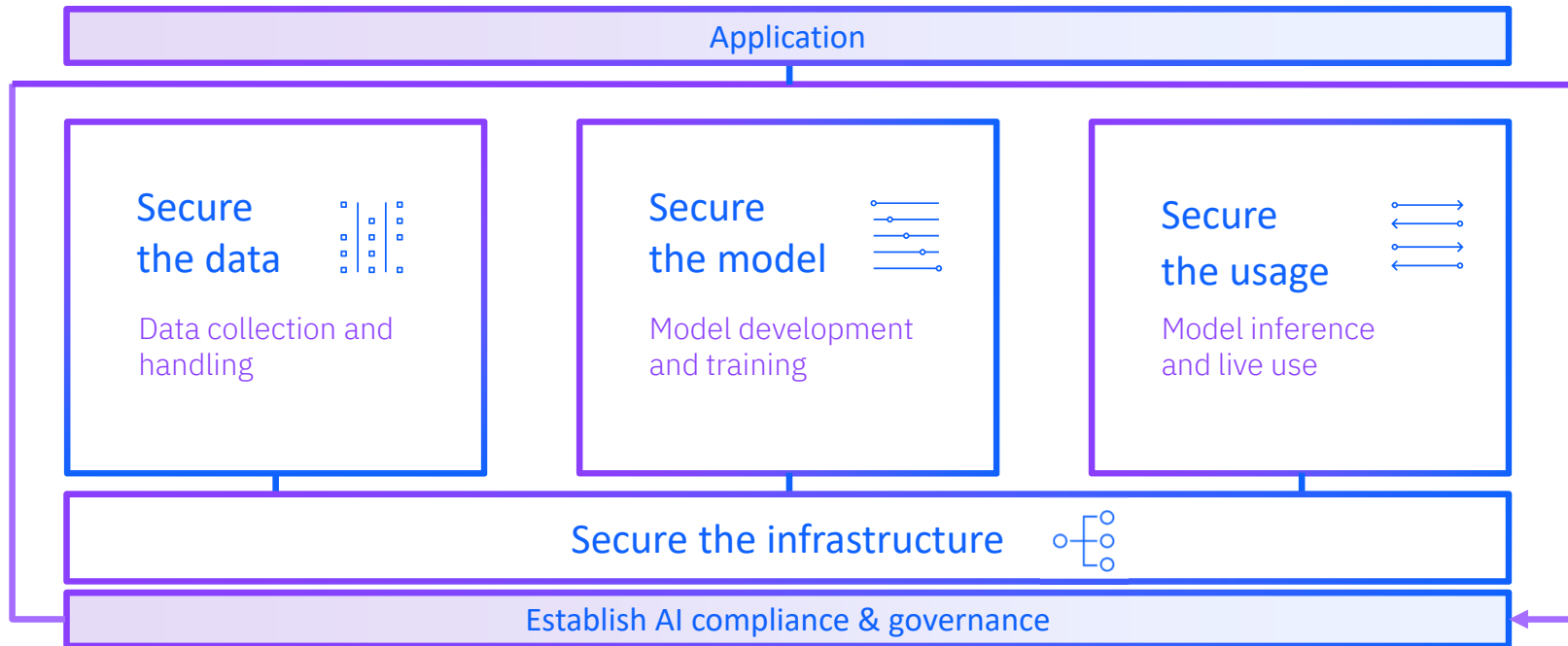
Generate content



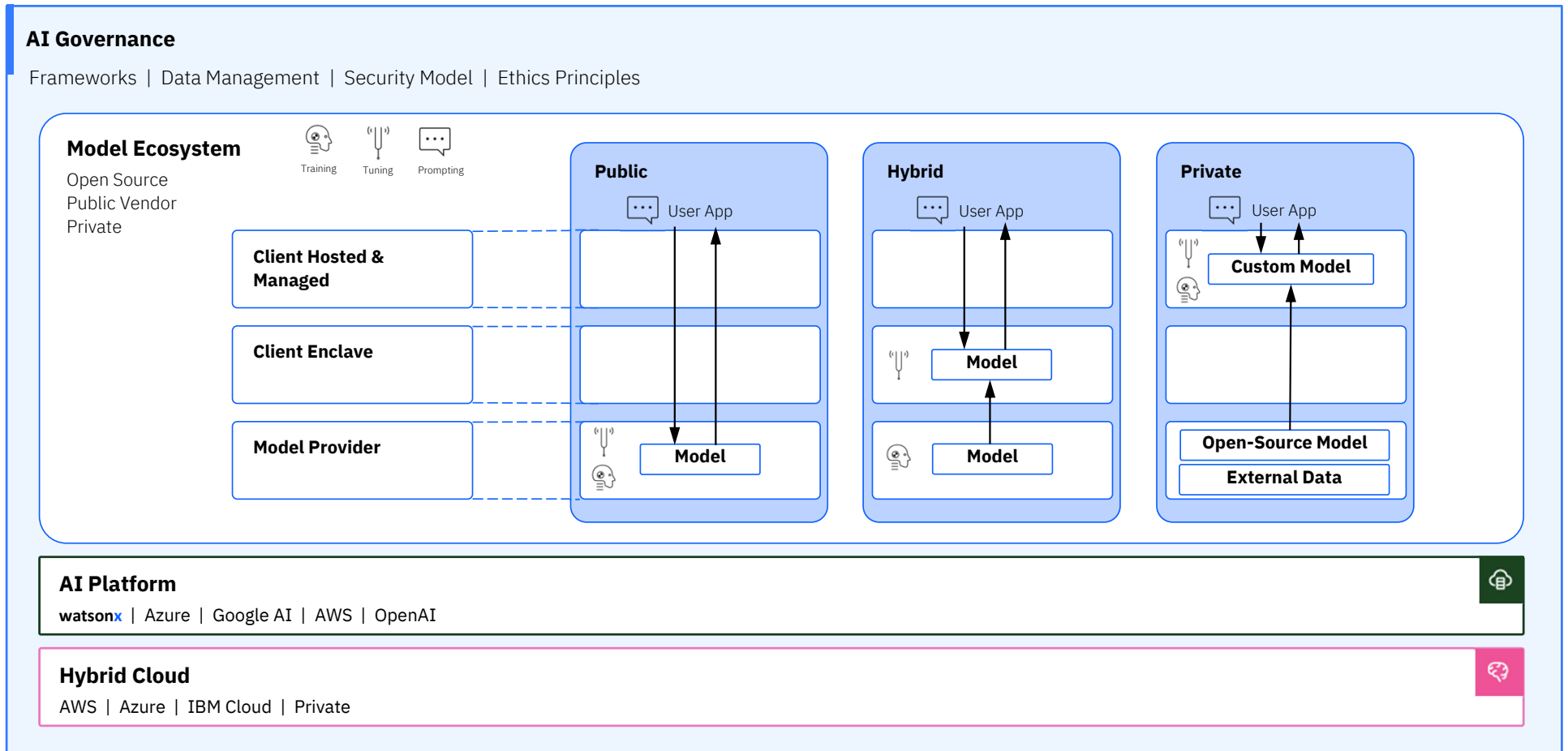
Learn & create active responses

Security for AI framework

Build trustworthy AI



No matter how clients consume AI security is critical to protect their data



AI / Gen AI intensifies existing threats and introduces new threats

New Threats

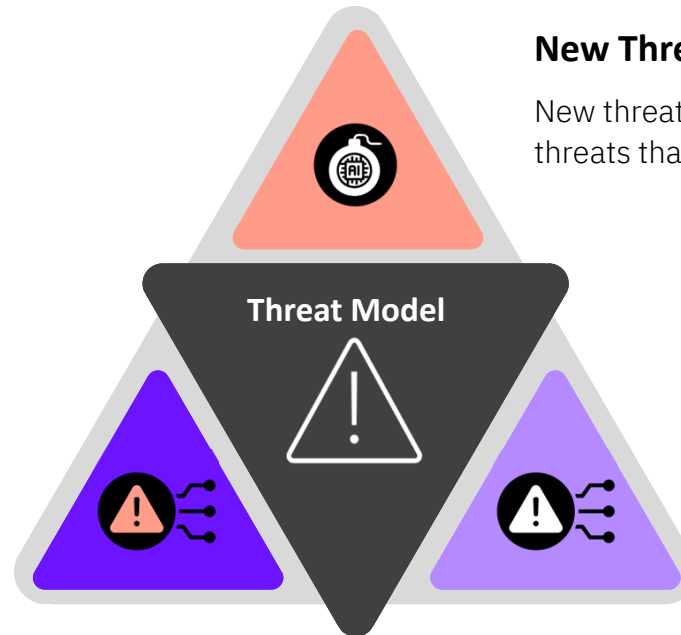
New threat landscape: these are threats that are specific to AI/GenAI.

Old Threats New Meaning

Conventional threats that take on new meaning: with the addition of AI/GenAI, some conventional threats bring in new meaning and threats

Conventional Threats

These are the business-as-usual threats that need to be addressed for any solution.



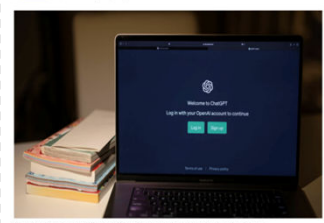
Current security frameworks and regulations are evolving



When using a **public** consumption model, security is primarily in the hands of the provider. But you still need to be **concerned about your data**.

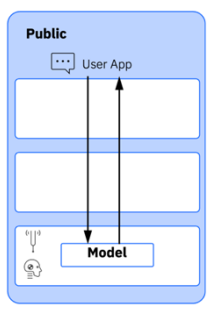
Samsung bans ChatGPT, AI chatbots after data leak blunder

By [David Hogue](#) on May 2, 2023



Samsung joins other companies that have banned or restricted ChatGPT because of data-leak risks. Credit: Getty Images

Samsung has banned the use of ChatGPT after employees inadvertently revealed sensitive information to the chatbot.



What Happened

“Employees inadvertently revealed sensitive information to the chatbot.”

“Samsung employees had shared source code with ChatGPT to check for errors and used it to summarize meeting notes. Information shared with ChatGPT is stored on OpenAI's servers and can be used to improve the model unless users opt out.”

Additional Challenges

- No control over the security protocols followed in the AI system.
- The wrong user having access to the system.
- Improper use of the AI system. (jailbreaking, prompt injection)

How Security for AI should be applied

Assess security gaps in posture and architecture for the security components you control.

Incorporate AI use as part of security awareness and education.

Implement role-based access at the user and system level (API).

Employ data loss prevention techniques to detect and prevent SPI,PII and regulated data leakage through prompts and API's.

Monitor for exploits, anomalous activities and insider threats.

How IBM Cybersecurity Services can help

AI Strategy & Governance

- Security Assessment w/ Threat Modeling ([CSS J7 OCC: 30BAA](#))
- Cyber Talent Transformation ([CSS J7 OCC: 30J8Y](#))

Identity & Access Management

- Identity & Access Security ([CSS J7 OCC: 30JA7](#), [30BU0](#))

Data Security

- Data Security for AI/ML ([CSS J7 OCC: 30J20](#), [30B2I](#), [30JY0](#), [30JNT](#))

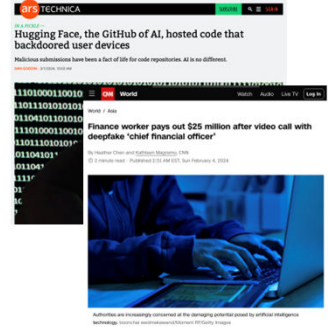
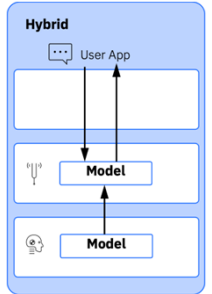
Threat Management

- Threat Management ([CSS J7 OCC: 30B33](#), [30B3A](#), [30JCC](#), [30BS2](#), [30BS0](#))

Partnering with new technology vendors to provide unique protection for AI.



When using a **hybrid** consumption model, **security of the model is shared**. The ability to fine tune the model to your business need requires more security

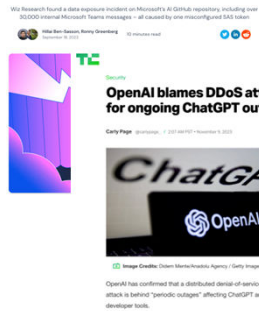
| | What Happened | How Security for AI should be applied | How IBM Cybersecurity Services can help |
|--|--|---|---|
|  | <p>“Code uploaded to AI developer platform Hugging Face covertly installed backdoors and other types of malware on end-user machines.”</p> <p>“A finance worker was tricked into paying out \$25 million to fraudsters using deepfake technology to pose as the company’s chief financial officer in a video conference call.”</p> | <p>Assess security gaps in posture and architecture for the security components you control.</p> <p>Establish effective governance model and policies to build a trustworthy AI.</p> <p>Implement role-based access at the user and system level (API).</p> | <p>AI Strategy & Governance</p> <ul style="list-style-type: none"> Security Assessment w/ Threat Modeling (CSS J7 OCC: 30BAA) AI Governance Programs (CSS J7 OCC: 30BAA) |
|  | <p>Additional Challenges</p> <p>Accidental data leakage of proprietary or SPI information.</p> <p>Attackers using model inversion to retrieve sensitive information</p> <p>Model poisoning compromising the model’s effectiveness, security or integrity</p> | <p>Institute model scanning to protect from supply chain attacks. Apply MLSecOps to secure your model lifecycle.</p> <p>Perform pen-testing on test ML model resiliency.</p> <p>Employ data loss prevention techniques to detect and prevent SPI,PII and regulated data leakage through prompts and API’s.</p> <p>Monitor for exploits, anomalous activities and insider threats.</p> | <p>Identity & Access Management</p> <ul style="list-style-type: none"> Identity & Access Security (CSS J7 OCC: 30JA7, 30BUO) <p>Vulnerability Management</p> <ul style="list-style-type: none"> Vulnerability Management (CSS J7 OCC: 30J0K) Sec DevOps Services (CSS J7 OCC: 30BVZ, 30J2A) <p>Adversarial ML</p> <ul style="list-style-type: none"> XFR Testing (CSS J7 OCC: 30JT6, 30B37) <p>Data Security</p> <ul style="list-style-type: none"> Data Security for AI/ML (CSS J7 OCC: 30J20, 30B21, 30JY0, 30JNT) <p>Threat Management</p> <ul style="list-style-type: none"> Threat Management (CSS J7 OCC: 30B33, 30B3A, 30JCC, 30BS2, 30BS0) |

Partnering with new technology vendors to provide unique protection for AI.



When using a **private** consumption model, **you are responsible for the security** of the infrastructure and the entire lifecycle of the AI models.

38TB of data accidentally exposed by Microsoft AI researchers



What Happened

Microsoft's AI research accidentally exposed 38TB of training data on GitHub because caused by one misconfigured SAS token.

ChatGPT experiences sporadic outages for 24 hours because of a DDoS attack.

Additional Challenges

Accidental data leakage of proprietary or SPI information.

Model evasion attacks that cause the model to misclassify or misinterpret inputs

Attackers attempting to steal models with model extraction techniques

How Security for AI should be applied

Assess security gaps in posture and architecture

Ensure cloud resources are configured to comply with security and compliance best practices.

Implement role-based access at the user and system level (API).

Secure cloud, infrastructure (HPCs), storage, and cloud AI services and posture

Institute model scanning to protect from supply chain attacks. Apply MLSecOps to secure your model lifecycle.

Perform pen-testing on test ML model resiliency.

Employ data loss prevention techniques to detect and prevent SPI,PII and regulated data leakage through prompts and API's.

Monitor for exploits, anomalous activities and insider threats.

How IBM Cybersecurity Services can help

AI Strategy & Governance

- Security Assessment w/ Threat Modeling ([CSS J7 OCC: 30BAA](#))
- Security Posture Management ([CSS J7 OCC: 30JA5](#))

Identity & Access Management

- Identity & Access Security ([CSS J7 OCC: 30JA7, 30BUO](#))

Infrastructure Security

- Infrastructure Security Services ([CSS J7 OCC: 30JAM, 30JLX](#))

Vulnerability Management

- Vulnerability Management ([CSS J7 OCC: 30JOK](#))
- Sec DevOps Services ([CSS J7 OCC: 30BVZ, 30JA2A](#))

Adversarial ML

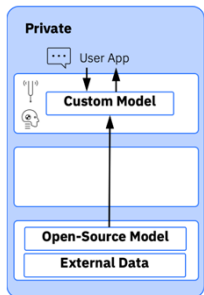
- XFR Testing Services ([CSS J7 OCC: 30J76, 30B37](#))

Data Security

- Data Security for AI/ML ([CSS J7 OCC: 30J20, 30B2I, 30JY0, 30JNT](#))

Threat Management

- Threat Management ([CSS J7 OCC: 30B33, 30B3A, 30JCC, 30BS2, 30BS0](#))

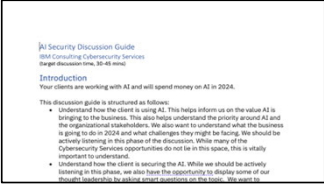


Partnering with new technology vendors to provide unique protection for AI.

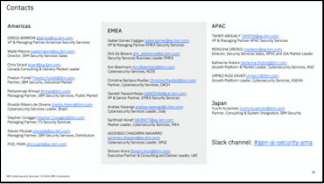


Engage your clients on their challenges with securing AI, Engage your local IBM Cybersecurity Services focal to help


AI Security Discussion Guide




IBM Cybersecurity Geo Leaders



AI transformation journey maps



Detailed Security for AI use cases



[Visit our Lighthouse page](#)

Partnering with new technology vendors to provide unique protection for AI



- Model Scanning
- Threat Monitoring with MLDR technologies
- AI Firewall
- Deep Fake recognition
- Data Protection for AI

Thank you

© 2024 International Business Machines Corporation

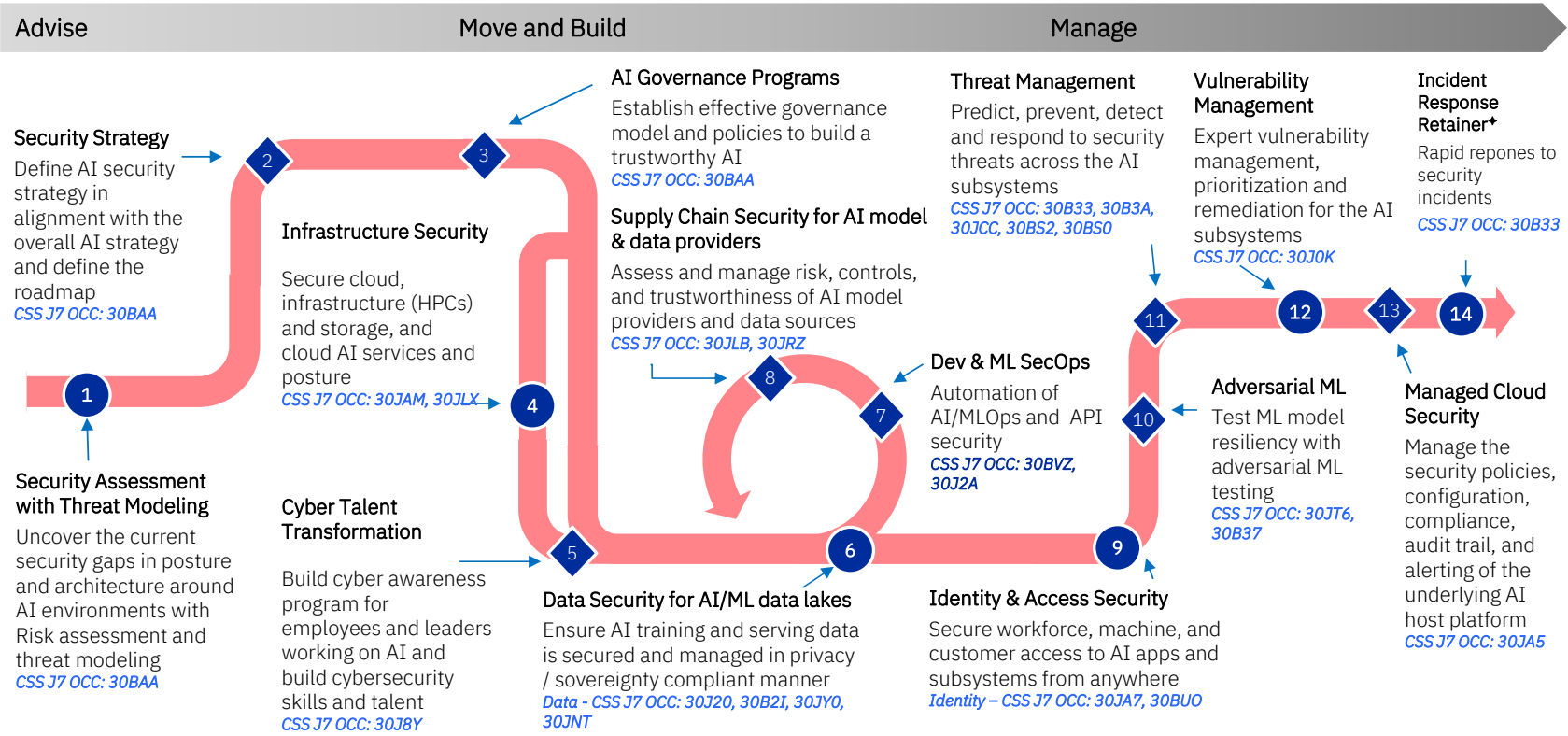
IBM and the IBM logo are trademarks of IBM Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on ibm.com/trademark.

Internal only. Do not distribute outside of IBM.

This document is current as of the initial date of publication and may be changed by IBM at any time

IBM

IBM Cybersecurity Services helps secure client's AI transformation journeys across the enterprise



- How do you secure your AI subsystems such as data lakes, models, and access across the entire platform?
- How do you secure the automated AI pipelines and MLOps?
- How do you secure models from being poisoned or jail-broken?
- How do you make sure we are complying with applicable legal and regulatory requirements while using large amounts of data for AI?

| IBMC Journey | Suggested Security Services | IBMC Journey | Suggested Security Services | IBMC Journey | Suggested Security Services |
|---|-----------------------------|---|-----------------------------|---|-----------------------------|
| AI Strategy & Advisory | 1, 2 | AI at Scale (Complex Engineering) | 1, 4, 6, 9 | AI at Scale (Complex Operations) | 6, 9, 12 |
| AI Governance and Talent Development | 3, 5 | Build and Deploy AI models, train data and MLOps | 3, 7, 10 | Manage AI models, training data and MLOps | 7, 8, 10, 11, 13 |
| Customer Care and Talent transformation on watsonx/ SaaS BYOLLM Deploy assistants | 1 | Customer Care and Talent transformation on watsonx/ SaaS BYOLLM assistants, API, and Model training | 1, 6, 9, 10 | | |

Our objective is to enable business to build and adopt AI that is secure, safe and trustworthy



Security for AI

Protecting foundation models, generative AI and their data sets is essential for enterprise-ready AI

Secure the underlying AI training data by protecting it from sensitive data theft, manipulation and compliance violations

Secure model development by scanning for vulnerabilities in the pipeline, hardening integrations and enforcing policies and access

Secure the usage of AI models by detecting data or prompt leakage and alerting on evasion, poisoning, extraction, or inference attacks

[IBM Adversarial Robustness Toolkit](#)



AI for Security

Productivity gains from foundation models and generative AI will reduce human bottlenecks in security

AI will manage repetitive security tasks such as summarizing alerts and log analysis, freeing teams to tackle strategic problems

AI will generate security content (detections, workflows, policies) faster than humans, expediting implementation and adjusting to changing security threats in real-time

AI will learn and create active responses that optimize over time, with abilities to find all similar incidents, update affected systems and patch vulnerable code