# AI & CYBERSECURITY
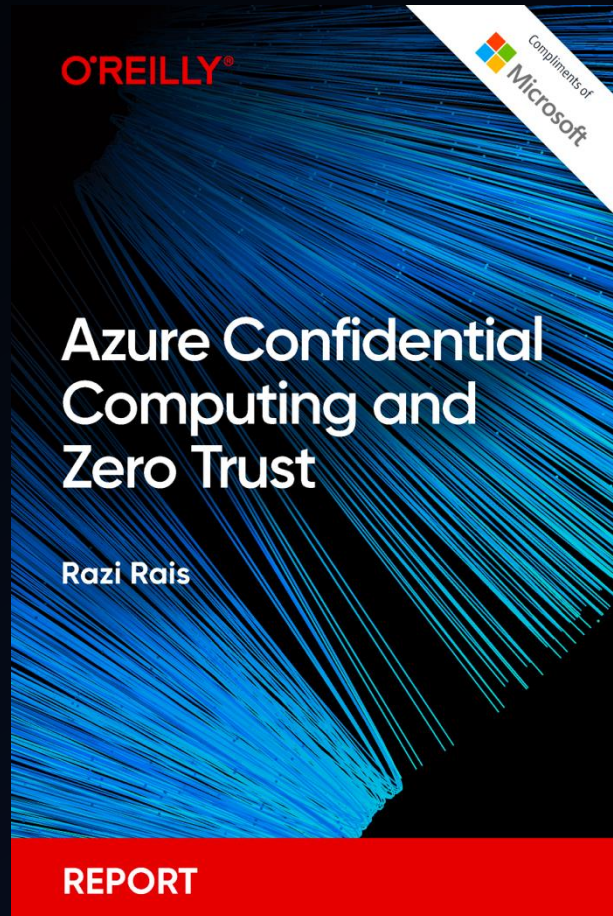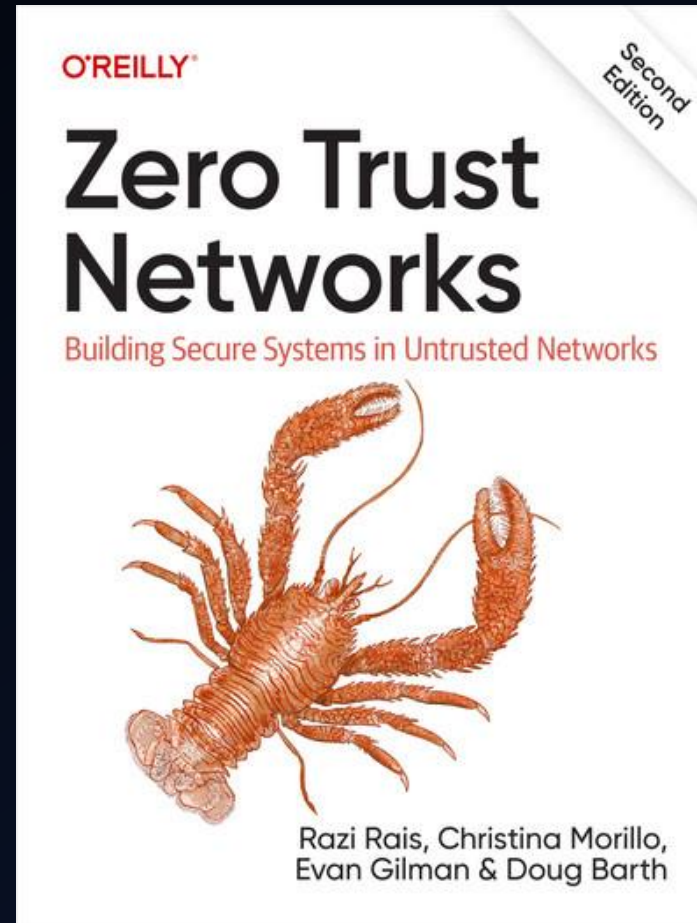# THE INTERSECTION OF DEFENSE, THREATS, AND PROTECTION

**Razi Rais**

# About

❖ **20+ years of experience** in the software development, architecture, and product management.

❖ **10+ years at Microsoft** in various teams as a software engineer, product manager, and architect.

❖ Currently working as a senior technical product manager at Microsoft helping businesses secure their digital identities at cloud scale.

❖ **Published author, speaker, and trainer**.

# Cybersecurity books



**Azure Confidential Computing and Zero Trust**
O'REILLY® — Compliments of Microsoft
Razi Rais
REPORT



**Zero Trust Networks**
O'REILLY® — Second Edition
Building Secure Systems in Untrusted Networks
Razi Rais, Christina Morillo, Evan Gilman & Doug Barth

[Read Online]                    [Read Online]

⚡ **Let's Connect!**

Looking for mentorship in AI and cybersecurity? Need an expert speaker for your next event? Working on an AI-powered cybersecurity project and need strategic guidance? Let's collaborate—reach out today on LinkedIn!
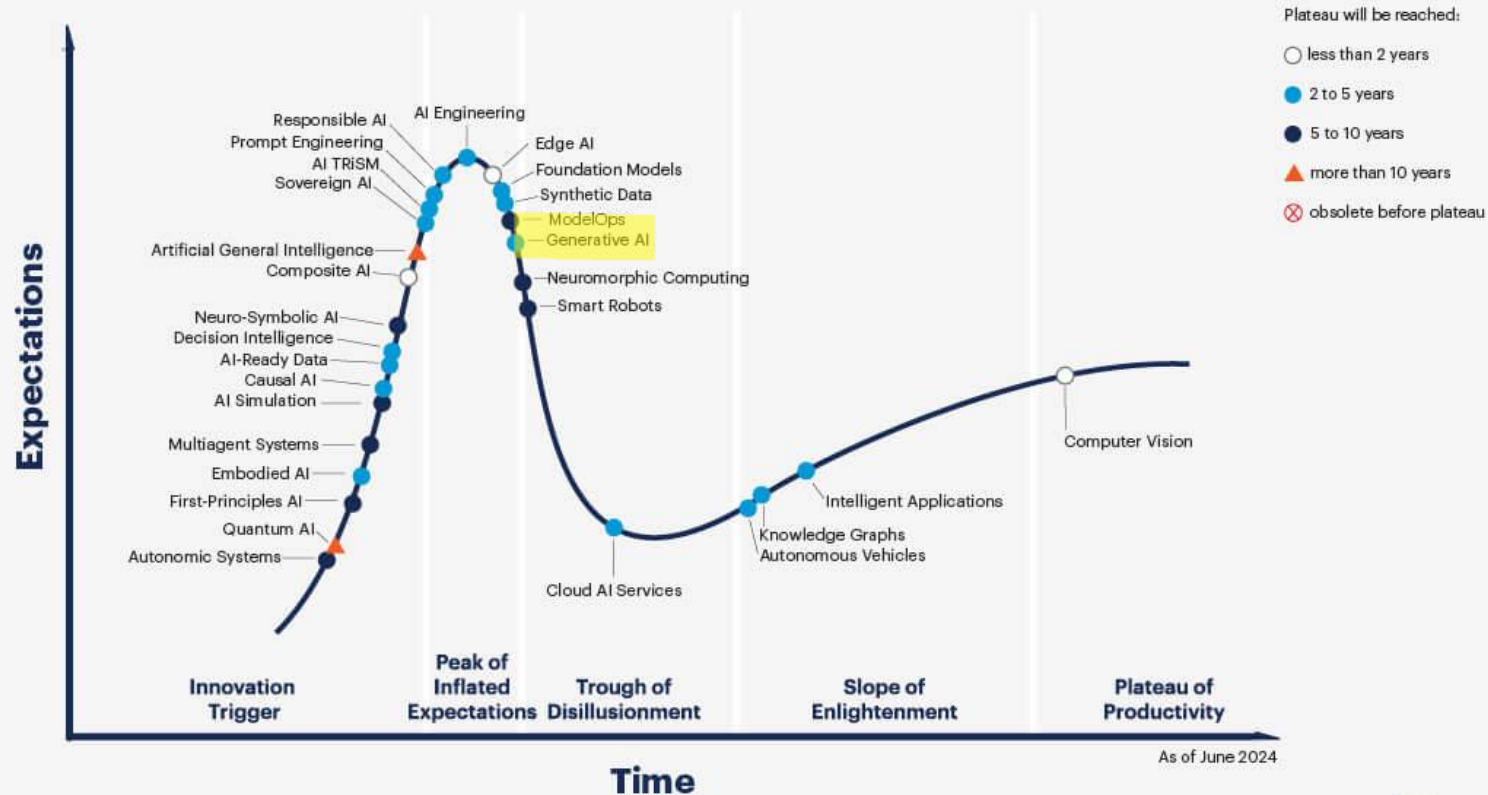
https://www.linkedin.com/in/razirais/

# Agenda

- ✓ Role of AI in Cybersecurity: Security of AI + AI in Security

- ✓ Learning Resources

- ✓ Discussion

# Gartner – Hyper Cycle for AI 2024

# Artificial Intelligence  (AI)

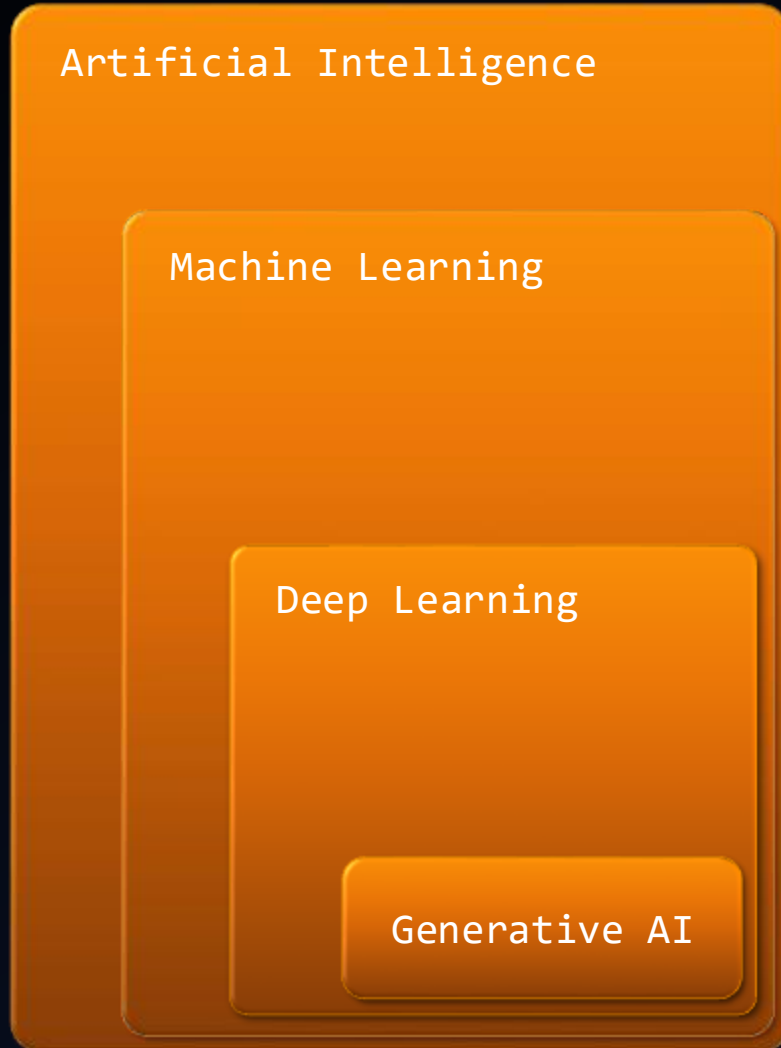**Artificial Intelligence**

**Machine Learning**

**Deep Learning**

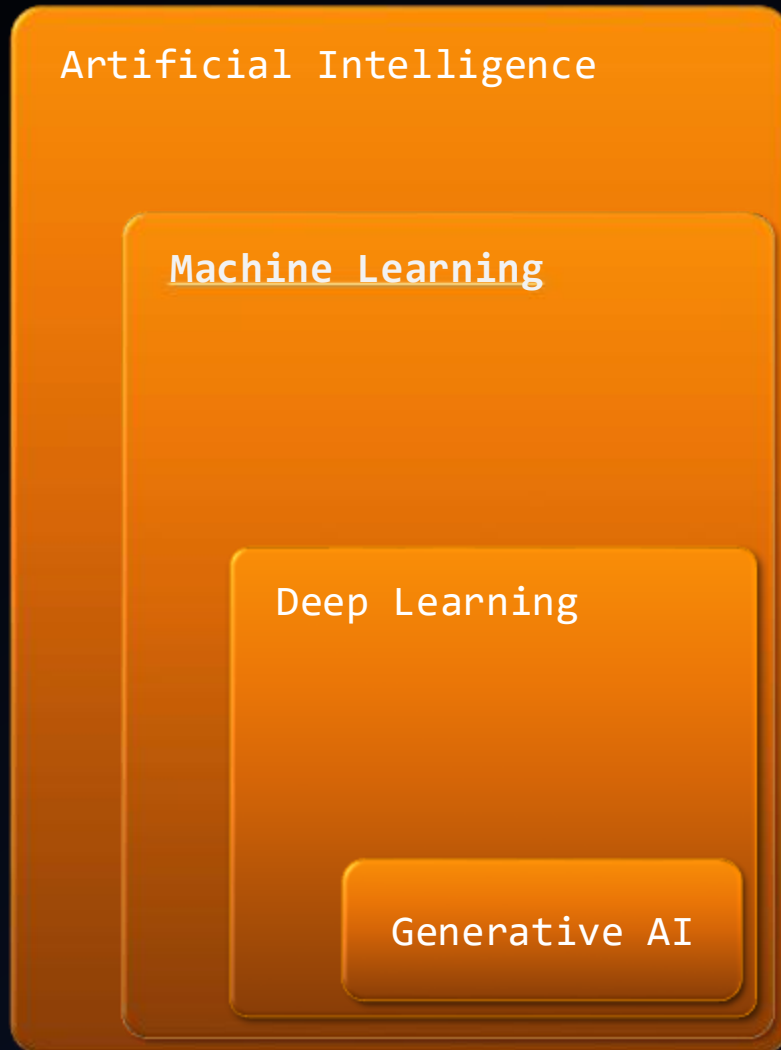**Generative AI**

- **Artificial Intelligence**: The field of computer science that seeks to create intelligent machines that can replicate or exceed human intelligence

- **Machine Learning**: Subset of AI that enables machines to learn from existing data and improve upon that data to make decisions or predictions

- **Deep Learning**: A machine learning technique in which layers of neural networks are used to process data and make decisions

- **Generative AI:** Create new written, visual, and auditory content given prompts or existing data

# Artificial Intelligence  (AI) Cont.

**Artificial Intelligence**

**Machine Learning**
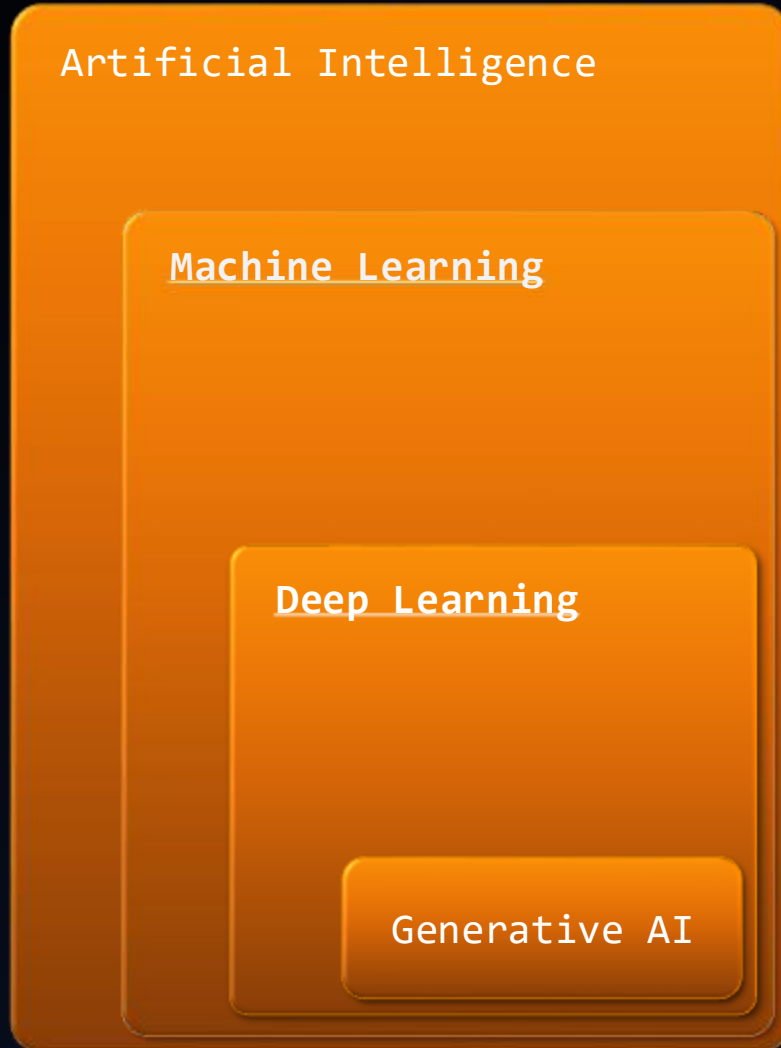
Deep Learning

Generative AI

- **Machine Learning**: The first decade of the 2000s marked the rapid advance of various machine learning techniques that could analyze massive amounts of online data to draw conclusions – or "learn" – from the results. Since then, companies have viewed machine learning as an incredibly powerful field of AI for analyzing data, finding patterns, generating insights, making predictions and automating tasks at a pace and on a scale that was previously impossible.
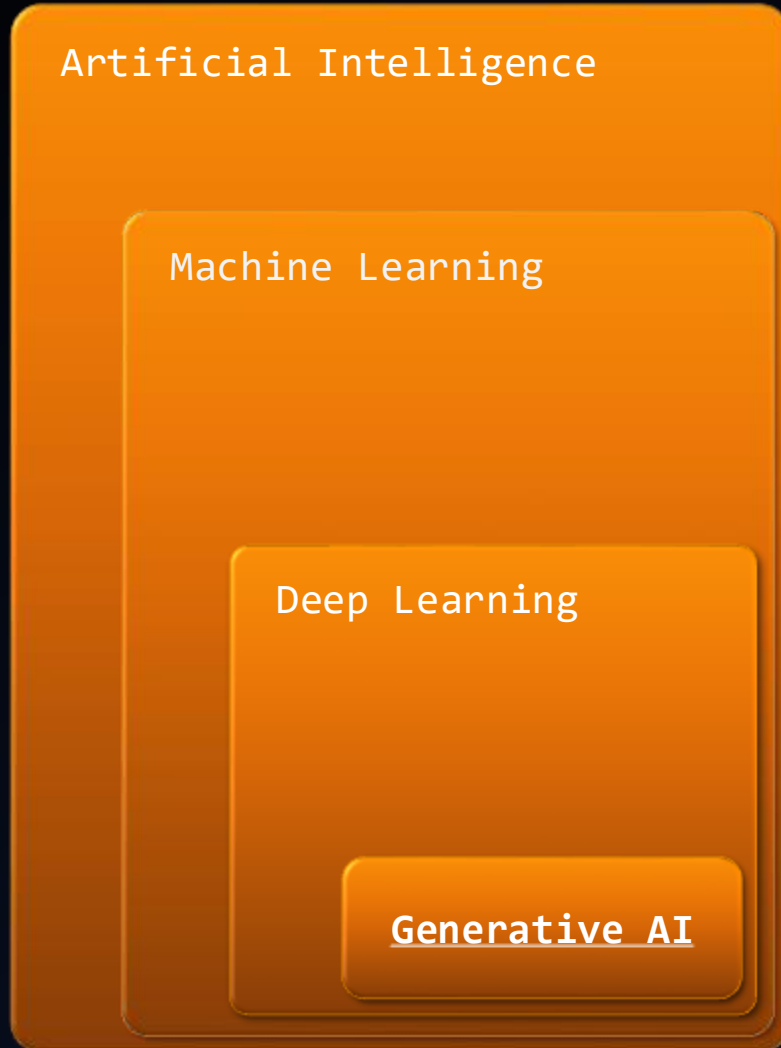
# Artificial Intelligence (AI) Cont.

Artificial Intelligence

Machine Learning

Deep Learning

Generative AI

- **Deep Learning:** The 2010s produced advances in AI's perception capabilities in the field of machine learning called deep learning. Breakthroughs in deep learning enable the computer vision that search engines and self-driving cars use to classify and detect objects, as well as the voice recognition that allows popular AI speech assistants to respond to users in a natural way.

# Artificial Intelligence (AI) Cont.

Artificial Intelligence

Machine Learning

Deep Learning

**Generative AI**

- **Generative AI:** Building on exponential increases in the size and capabilities of deep learning models, the 2020s will be about language mastery. The GPT-4 language model, developed by OpenAI, marks the beginning of a new phase in the abilities of language-based AI applications. Models such as this will have far-reaching consequences for business, since language permeates everything, an organization does day to day—its institutional knowledge, communication and processes.

# Terminology

❖ Model → Patterns or relationships in data

❖ Large Language Models (LLM)

  ▪ Large: More data than can be manually labeled

  ▪ Language: Match context and words (e.g., word prediction, creative writing)

  ▪ Model: Semi-supervised learning

# ChatGPT

❖ChatGPT: AI Chatbot, developed by [OpenAI](), trained to perform conversational tasks and creative tasks

❖Conversation-in and message-out

❖Trained over 175 billion machine learning parameters

❖GPT-4 and above are multimodal (e.g., images + text)

# Domain Specific LLMs

❖ While generic LLMs (like ChatGPT) are great fit for general queries but they cannot understand the specific context beyond the massive datasets they are trained with.

❖ If you didn't train a language model with domain specific data, the results may be less than ideal.

❖ So, we need custom models with a better language understanding of a specific domain (Finance, Healthcare, Cybersecurity etc.)

# Domain Specific LLMs (Cont.)

- ❖ BlackrockGPT
- ❖ GoldmanSachsGPT
- ❖ StripeGPT
- ❖ MorningstarGPT
- ❖ RobinhoodGPT
- ❖ VanguardGPT
- ❖ SoFiGPT

- ❖ ChubbGPT
- ❖ RevolutGPT
- ❖ ChatLAW
- ❖ KAI-GPT
- ❖ FinGPT
- ❖ ClimateBERT

# AI in Security

**Using AI to enhance cybersecurity, such as preventing cyberattacks, optimizing security processes, and improving security resilience.**

- Intrusion Detection and Prevention
- User Behavior Analytics
- Vulnerability Assessment
- Cyber threat intelligence
- Phishing Protection

# SOC Use case: Security posture management

- Assist with evaluating whether an organization is vulnerable to known vulnerabilities and exploits.

- Assist in risk prioritization

- Assist in resolving weaknesses by making specific advice.

# SOC Use case: Incident response

- Assist in identification of an ongoing attack

- Assist in assessing scale of an attack

- Assist in providing  guidance around remediation

# SOC Use case: Security reporting

- Assist in easily summarizing an event, incident, or threat

- Assist with the preparation of information into shareable and customizable reports

# Directionally where are we heading?



**Dialing in the human-agent ratio**

As leaders assemble human-agent teams, they'll need to get the balance right for each role, function, or project to ensure optimal performance on both sides of the equation.

**Too few agents per person**
Underutilizes both agentic and human resources, leaving potential efficiencies on the table

**Too many agents per person**
Overwhelms the human capacity for applying judgment and decision making, introducing business risk and potential employee burnout

**Optimal balance**
Agents enhance productivity and innovation while humans provide robust guidance and oversight

# Security of AI

**AI can be used by bad actors with malicious intent such as criminals, terrorists, and hostile nation-states.**

- Deep fakes

- Disinformation campaigns

- Misuse of military robots

- Autonomous weapon systems

- Social engineering

- Hacking and cyber attacks

# Adversaria use of AI



## Adversarial use of AI in influence operations

| Capability | 🇨🇳 China | 🇷🇺 Russia | ☪ Iran & proxies |
|---|---|---|---|
| Text | MEDIUM / LOW | MEDIUM / LOW | LOW |
| Image | HIGH | HIGH | MEDIUM / LOW |
| Audio/video | HIGH | HIGH | LOW |
| Example | May 2024: Bespoke Taizi Flood AI-generated cartoon | June 2024: AI-generated audio of Elon Musk narrating fabricated documentary | April 2024: Likely AI-generated video leading up to Iranian military operation |

# Nation State Actors & Targeted Sectors

## Russia

### Nation state threat actor activity

#### Targeting by region

| | Sector | Percentage |
|---|---|---|
| 1 | Europe & Central Asia | 68% |
| 2 | North America | 20% |
| 3 | Middle East & North Africa | 5% |
| 4 | East Asia & Pacific | 3% |
| 5 | Latin America & Caribbean | 3% |
| 6 | South Asia | 1% |
| 7 | Sub-Saharan Africa | 1% |

Approximately 75% of targets were in Ukraine or a NATO member state, as Moscow seeks to collect intelligence on the West's policies on the war. Ukraine remains the country most targeted by Russian actors.

#### Most targeted sectors

| | Sector | Percentage |
|---|---|---|
| 1 | Government | 33% |
| 2 | IT | 15% |
| 3 | Think tanks and NGOs | 15% |
| 4 | Education and Research | 9% |
| 5 | Inter-governmental organization | 4% |
| 6 | Defense Industry | 4% |
| 7 | Transportation | 3% |
| 8 | Energy | 2% |
| 9 | Media | 2% |
| 10 | All others | 13% |

Russian actors focused their targeting against European and North American government agencies and think tanks, likely for intelligence collection related to the war in Ukraine. Actors like Midnight Blizzard also targeted the IT sector, suggesting it was in part planning supply-chain attacks to gain access to these companies' client's networks for follow-on operations.

## China

### Nation state threat actor activity

#### Targeting by region

| | Sector | Percentage |
|---|---|---|
| 1 | East Asia & Pacific | 39% |
| 2 | North America | 33% |
| 3 | Europe & Central Asia | 12% |
| 4 | Latin America & Caribbean | 8% |
| 5 | South Asia | 4% |
| 6 | Middle East & North Africa | 2% |
| 7 | Sub-Saharan Africa | 2% |

Chinese threat actors' targeting efforts remain similar to the last few years in terms of geographies targeted and intensity of targeting per location. While numerous threat actors target the United States across a wide variety of sectors, targeting in Taiwan is largely limited to one threat actor, Flax Typhoon.

#### Most targeted sectors

| | Sector | Percentage |
|---|---|---|
| 1 | IT | 24% |
| 2 | Education and Research | 22% |
| 3 | Government | 20% |
| 4 | Think tanks and NGOs | 10% |
| 5 | Manufacturing | 4% |
| 6 | Defense Industry | 3% |
| 7 | Communications | 3% |
| 8 | Finance | 3% |
| 9 | Transportation | 2% |
| 10 | All others | 9% |

Most Chinese threat activity is for intelligence collection purposes and was especially prevalent in ASEAN countries around the South China Sea. Granite Typhoon and Raspberry Typhoon were the most active in the region, while Nylon Typhoon continued to target government and foreign affairs entities globally.

# OWASP: Top 10 for LLM

**LLM01**

## Prompt Injection

This manipulates a large language model (LLM) through crafty inputs, causing unintended actions by the LLM. Direct injections overwrite system prompts, while indirect ones manipulate inputs from external sources.

**LLM02**

## Insecure Output Handling

This vulnerability occurs when an LLM output is accepted without scrutiny, exposing backend systems. Misuse may lead to severe consequences like XSS, CSRF, SSRF, privilege escalation, or remote code execution.

**LLM03**

## Training Data Poisoning

Training data poisoning refers to manipulating the data or fine-tuning process to introduce vulnerabilities, backdoors or biases that could compromise the model's security, effectiveness or ethical behavior.

**LLM04**

## Model Denial of Service

Attackers cause resource-heavy operations on LLMs, leading to service degradation or high costs. The vulnerability is magnified due to the resource-intensive nature of LLMs and unpredictability of user inputs.

**LLM05**

## Supply Chain Vulnerabilities

LLM application lifecycle can be compromised by vulnerable components or services, leading to security attacks. Using third-party datasets, pre- trained models, and plugins add vulnerabilities.

**LLM06**

## Sensitive Information Disclosure

LLMs may inadvertently reveal confidential data in its responses, leading to unauthorized data access, privacy violations, and security breaches. Implement data sanitization and strict user policies to mitigate this.

**LLM07**

## Insecure Plugin Design

LLM plugins can have insecure inputs and insufficient access control due to lack of application control. Attackers can exploit these vulnerabilities, resulting in severe consequences like remote code execution.

**LLM08**

## Excessive Agency

LLM-based systems may undertake actions leading to unintended consequences. The issue arises from excessive functionality, permissions, or autonomy granted to the LLM-based systems.

**LLM09**

## Overreliance

Systems or people overly depending on LLMs without oversight may face misinformation, miscommunication, legal issues, and security vulnerabilities due to incorrect or inappropriate content generated by LLMs.

**LLM10**

## Model Theft

This involves unauthorized access, copying, or exfiltration of proprietary LLM models. The impact includes economic losses, compromised competitive advantage, and potential access to sensitive information.

OWASP Top 10 LLM Applications and Generative AI - 2025 Version
Example LLM Application and Basic Threat Modeling
Ads Dawson (GangGreenTemperTatum) | https://genai.owasp.org/ Nov 2025 - v.01 SaaS LLM application

# OWASP: LLM AI Cybersecurity & Governance Checklist



**LLM AI Cybersecurity & Governance Checklist**

From the OWASP Top 10
for LLM Applications Team

**Version: 1.1**

**Published:** *April 11, 2024*

Ref: https://atlas.mitre.org/matrices/ATLAS

# OWASP: Other AI Resources



https://mltop10.info/



https://owaspai.org/



https://owasp.org/www-project-ai-security-and-privacy-guide/

# NIST: AI Risk Framework & Adversarial ML

# NIST: AI Risk Framework



## Harm to People

- Individual: Harm to a person's civil liberties, rights, physical or psychological safety, or economic opportunity.
- Group/Community: Harm to a group such as discrimination against a population sub-group.
- Societal: Harm to democratic participation or educational access.

## Harm to an Organization

- Harm to an organization's business operations.
- Harm to an organization from security breaches or monetary loss.
- Harm to an organization's reputation.

## Harm to an Ecosystem

- Harm to interconnected and interdependent elements and resources.
- Harm to the global financial system, supply chain, or interrelated systems.
- Harm to natural resources, the environment, and planet.

## AI Risk Management Framework

**Map** — Context is recognized and risks related to context are identified

**Measure** — Identified risks are assessed, analyzed, or tracked

**Govern** — A culture of risk management is cultivated and present

**Manage** — Risks are prioritized and acted upon based on a projected impact

# MITRE: Adversarial Threat Landscape for AI Systems (ATLAS™)

# NIST: AI Risk Framework Use Cases

# NSA - Guidance for Strengthening AI System Security

# Artificial Intelligence Engineer Certificate by AETIBA

## AI Engineers



https://www.artiba.org/certification/artificial-intelligence-certification

## AiE by Artificial Intelligence Board of America (ARTIBA)

**Essentials of Artificial Intelligence & Machine Learning**

### 27%

- Artificial Intelligence Ecosystem
- Supervised Learning
- Ensemble Learning
- Unsupervised Learning

**Essentials of Ai & ML Programming**

### 21%

- Building Recommender Systems
- Logic Programming
- Heuristic Search Techniques
- Genetic Algorithms
- Building Games With Ai

**Essentials of Natural Language Processing**

### 26%

- Natural Language Processing
- NLP Development & Applications
- Probabilistic Reasoning for Sequential Data
- Speech Recognizer
- Object Detection and Tracking

**Essentials of Neural Networks & Deep Learning**

### 26%

- Neural Networks
- Neural Network Applications
- Reinforcement Learning
- Deep Learning with Convolutional Neural Networks

# NVIDIA AI Certifications

## NVIDIA-Certified Associate: Generative AI LLMs

An associate-level assessment for individuals who are looking to validate their skills in the use of generative AI and large language models.

**Learn About This Certification** ›

## NVIDIA-Certified Associate: Generative AI Multimodal

An associate-level assessment for individuals who are looking to validate their skills in the use of multimodal generative AI.

**Learn About This Certification** ›

## NVIDIA-Certified Associate: AI in the Data Center

An associate-level assessment for IT professionals and others looking to validate their skills in AI infrastructure in the data center.

**Learn About This Certification** ›
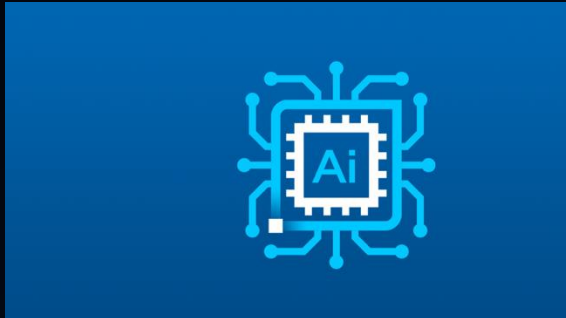
## NVIDIA-Certified Professional: InfiniBand

A professional-level assessment for networking and IT professionals looking to validate their skills in AI networking by NVIDIA.

**Learn About This Certification** ›

# AI Certificate Courses

Intel Edge AI Certification



https://www.intel.com/content/www/us/en/developer/tools/devcloud/edge/learn/certification.html

JETSON AI  Courses and Certifications



https://learn.microsoft.com/en-us/credentials/certifications/azure-ai-engineer/

edX



https://www.edx.org/learn/artificial-intelligence#browse-courses

Thank you!