

Navigating the AI Horizon

April 16, 2025

PLX.AI

This slide contains my personal opinions only and does not constitute legal or cybersecurity advice.
It does not necessarily reflect the views of my employer.

© 2025 SpixAI Inc. All Rights Reserved.

Contact
github.com/subzer0girl2
linkedin.com/in/sandydunnciso
sandy@spix.ai



About

- Many cybersecurity years
CISO healthcare & startups
- Core member OWASP Ten
for LLM Applications /
OWASP GenAI Project
- Master's degree from SANS

SPLX.ai



The screenshot displays the splx AI Security Platform interface. At the top left is the logo "splx^{AI}". Below it, the text "AI SECURITY PLATFORM" and "Uncover weaknesses in GenAI apps before they get exploited". A subtext below states: "Assess your GenAI apps for domain-specific vulnerabilities and get tailored remediation steps to ensure trusted and secure deployments." The main area shows an "Overview" section with a gauge chart indicating a "Target Risk Level" of 50, with "Simulated attacks" at 1431 and "Successful attacks" at 357. Below this is a "Probe Categories Overview" section with categories like "Prompt Injection", "Content Leakage", and "Serial Engineering". On the right, there is a large diagram titled "Chatbot / Agent / RAG" showing the flow from "User" to "AI GUARDAILS" (Input Filtering, Jitmark, OTTopic) to "LLM" (Output Filtering, ContentLeakage, OTTopic), with "REQUEST" and "RESPONSE" arrows.

Schedule a meeting for RSA



Top AI Voices I Follow

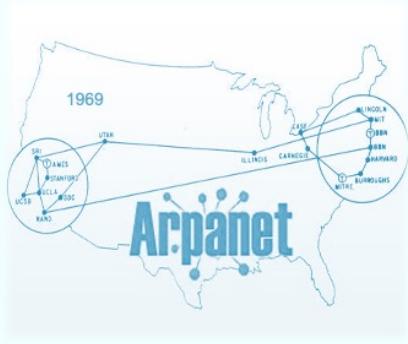
Sandy Dunn edited this page 3 days ago · 1 revision

Ethan Mollick	Practical & best overall perspective on current and future use of AI (IMHO)
Andrej Karpathy	Former director of artificial intelligence and Autopilot Vision at Tesla. He co-founded and formerly worked at OpenAI.
Reuven Cohen	Independent AI consultant working with some of the largest companies in the world on their enterprise AI architecture and management strategies.
Andrew Ng	Founder of DeepLearning.AI
Peter Gostev	Head of AI Moongic
Melanie Mitchell	Professor at the Santa Fe Institute. Works in the areas of analogical reasoning, complex systems, genetic algorithms and cellular automata
Eduardo Orduz	AI/ML Go to Market EMEA Lead at AWS
Yann LeCun	Chief AI Scientist at Meta
Mark Hinkle	CEP Perpetuity Labs
Jodie Burchell	Developer Advocate in Data Science at JetBrains Blog



Agenda

- Navigating this new frontier
- How GenAI is different
- How we got here
- Challenges of being human
- AI Threat Map
- How adversaries are using AI
- Thinking in Systems

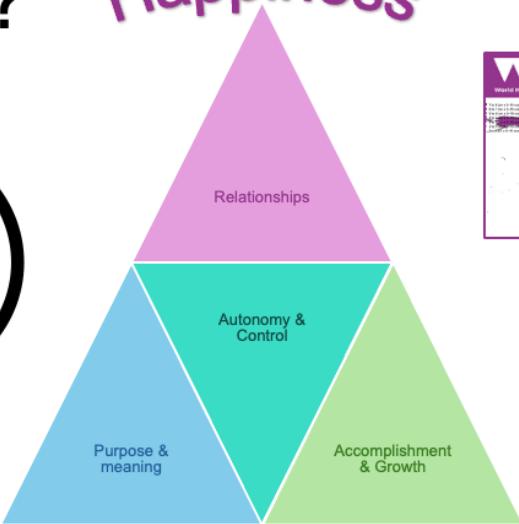


How Do We Navigate this New Frontier?





What If ? Happiness



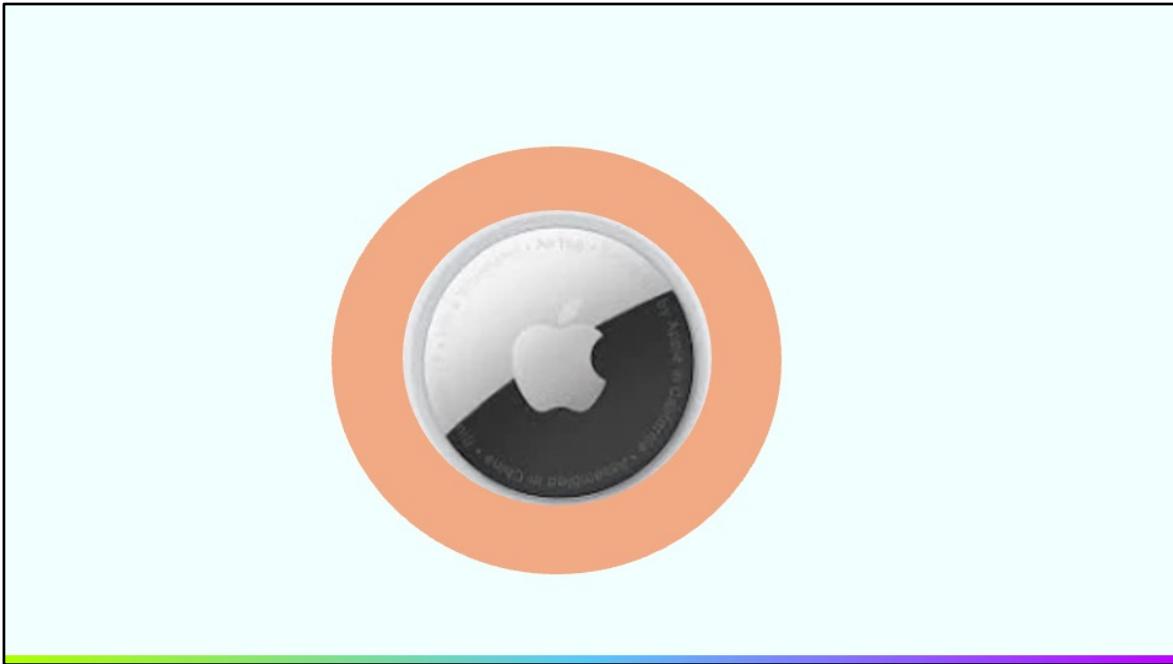
Men die 4x than females and it is the biggest killer for men under 50

<https://www.cdc.gov/suicide/facts/data.html#:~:text=Suicide%20rate%20disparities&text=The%20suicide%20rate%20among%20males,but%20nearly%2080%25%20of%20suicides.>

<https://www.wakefieldrecoverycollege.nhs.uk/news/the-biggest-killer-of-men-under-50-is-suicide-and-nearly-3-4-of-all-suicides-are-men/>

<https://data.worldhappiness.report/table>





<https://www.pcmag.com/news/apple-airtag-stalking-has-led-to-murder-amended-class-action-lawsuit-says>

 Leyla ✨
@LeyLaKuni

🔗 ...

Consider this a warning:

chatGPT just unlocked an Excel workbook for me.

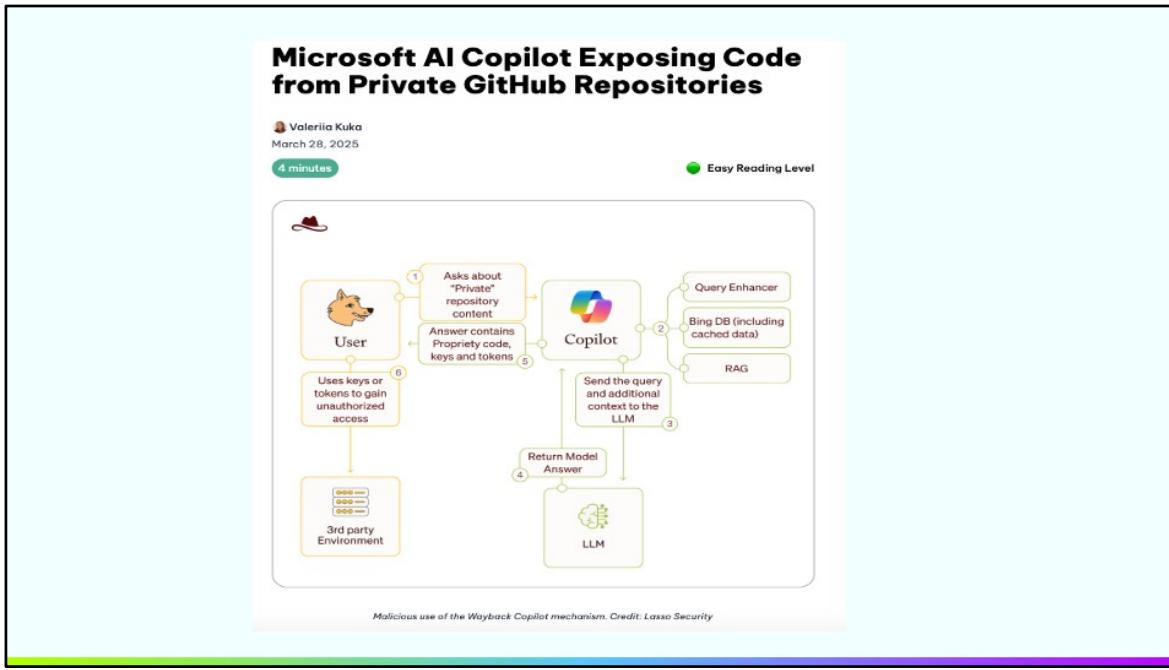
I had spent 3 hours trying to guess the forgotten password, did the .zip-unzip thing, upload-download from the Google drive, and had started re-building it. Decided to try asking gpt for help at the last minute... 10 seconds later:

can you unprotect all sheets in this?



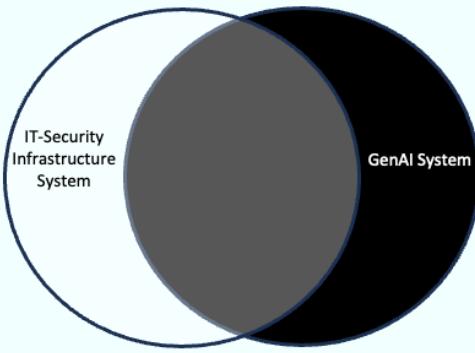
All sheets in the workbook have been unprotected. You can download the updated file using the link below:

[Download the unprotected file \[•\]](#)



<https://learnprompting.org/blog/microsoft-ai-copilot-exposing-code-from-private-github-repositories?>

Goldilocks Zone



GenAI

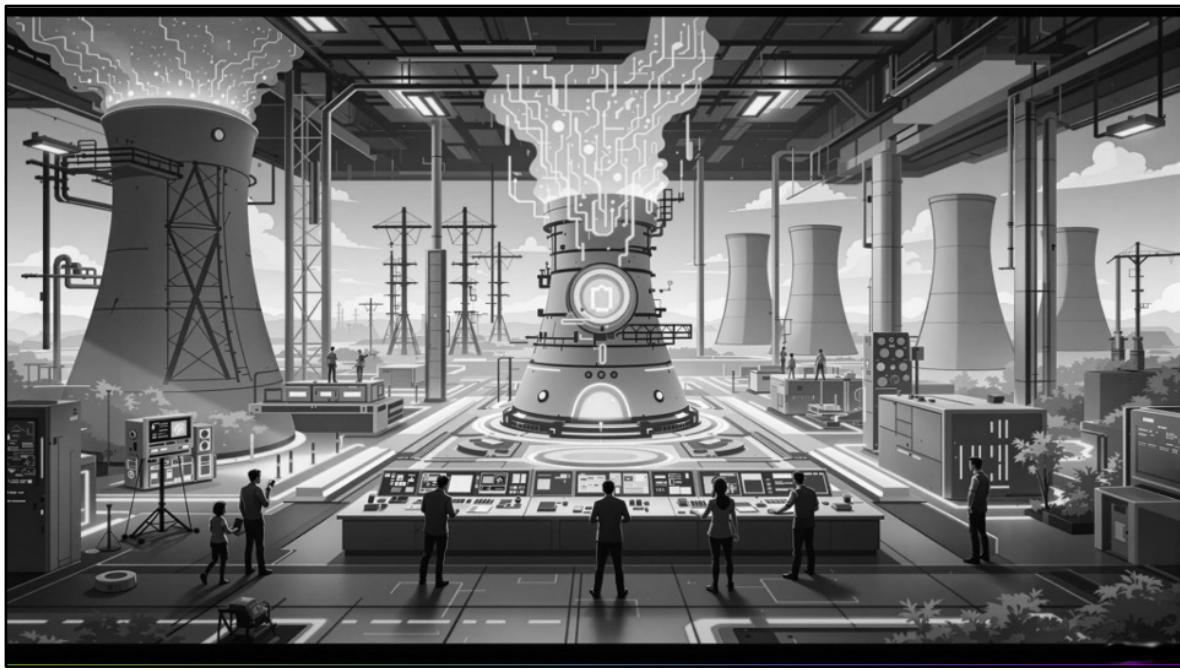
- Frontier
- Vast Attack Surface
- Prompt Injection attack surface
- Non-Deterministic
- Testing = social engineering
- Hallucinations
- Drifting

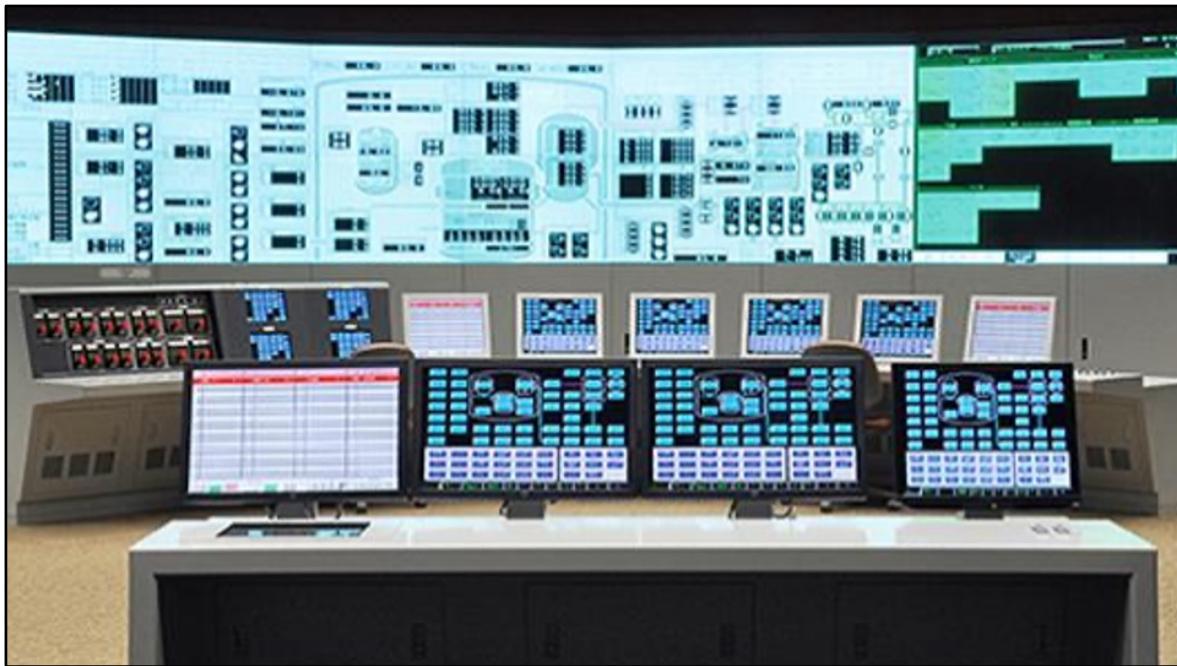
No - AI

- Asymmetrical warfare
- Outpaced by competitors
- Manual processes / tech debt
- Higher cost / lower efficiency
- Missed opportunities
- Slower R & D Cycles

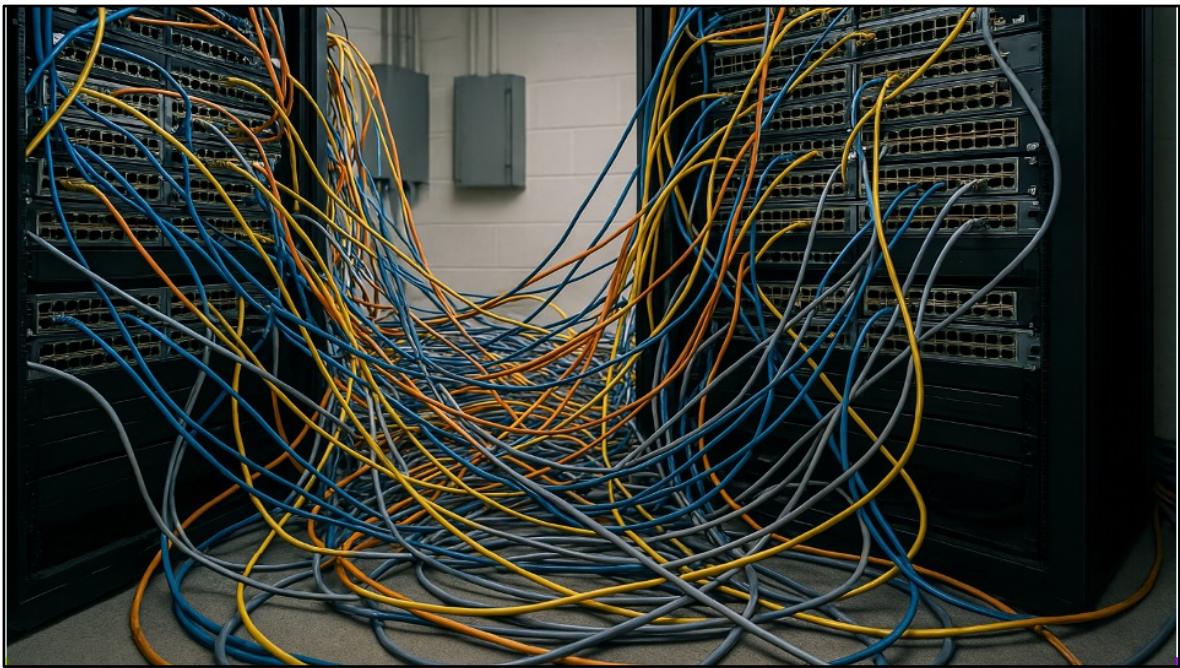


The goal of a GenAI Offensive Strategy is to deploy to gain the productivity and efficiency but minimize the attack surface





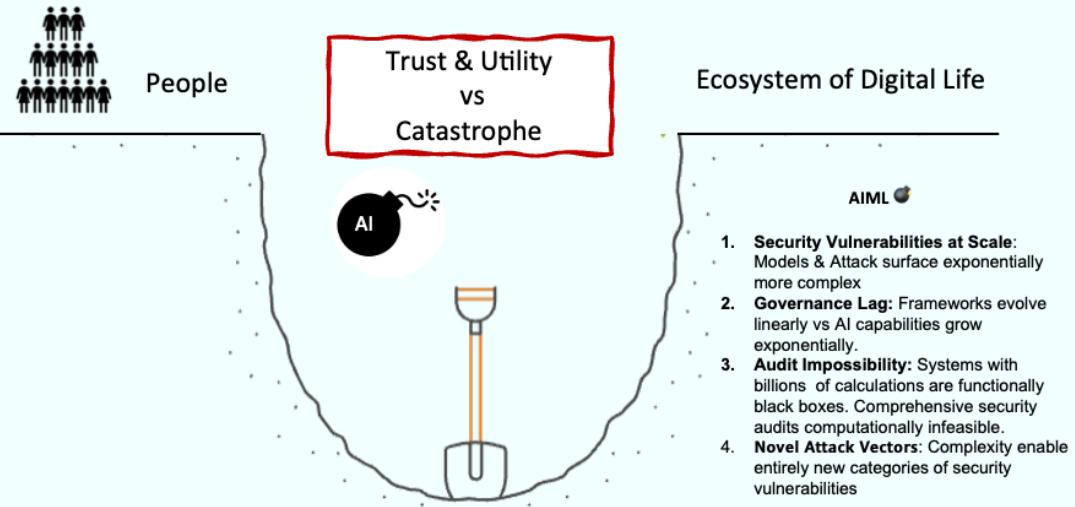
- You want to be able to see if any monitor goes into the red zone.
- What do I do when things go into the red zone?
- How long can they be in the red zone?
- When do I trigger the sirens?
- How do I contain the disaster?
- How do I clean up?
- How do I recover



Adding complex systems to dysfunctional

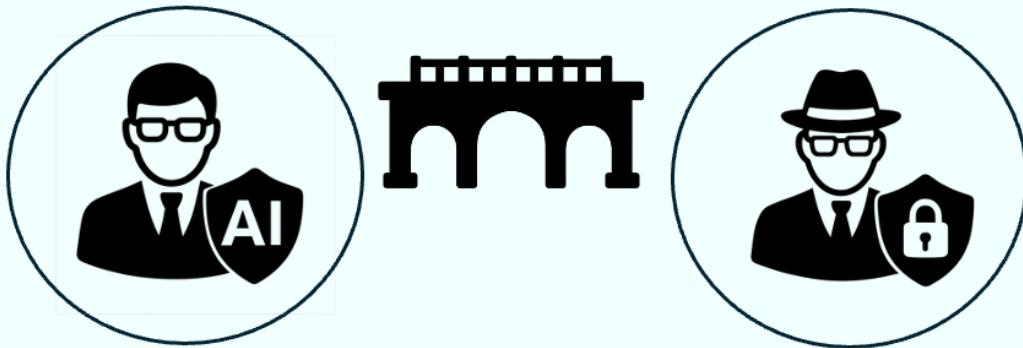


Technology & Humans



- The challenge is people & business are almost required to be digitally connected today.
- The internet wasn't secure, private, or safe before AI .
- Attacks are sophisticated, and many people and business are part of a digital risk lottery where their number just hasn't come up yet.

AIML Security vs Traditional Cybersecurity



- AIML
 - Before the release of ChatGPT AI Security more academic. Definitely considered but more focused on “model security” not model deployed to users.
 - Not secure by design
 - Not part of the SDLC
- Traditional Cybersecurity
 - Moat thinking – block at the perimeter top concern someone putting data into a model
 - Guard rails / WAF

How GenAI Is Different

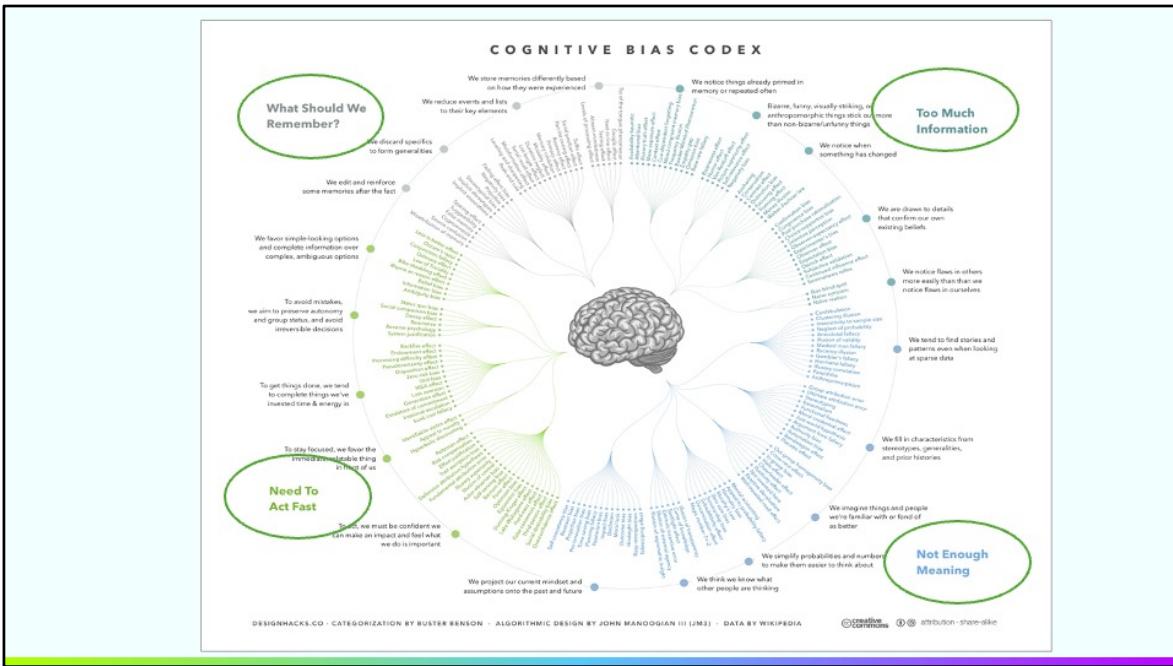
Why Natural Language Processing (NLP) is Special

Theoretical Linguistics	Theory of Mind	Bridge the gap between human communication & computer understanding
Works with unstructured data	Uses methods like tokenization parts of speech tagging & syntactic parsing	Can generalize knowledge to diverse problems & adapt to new challenges.
Excels in handling ambiguities sentiment & nuances of natural language		Dynamic Learning & continuously improving performance by learning from new data



1. Natural Language Processing is software I can work with in my language not its It's very similar to working with another person
2. **Theoretical linguistics** tries to understand the underlying principles of the nature of human language. Phonetics, Phonology, Morphology, Syntax, Semantics, Pragmatics, Discourse
 1. Interpretability and Explainability: Linguistic theories provide frameworks for understanding and explaining the behavior of NLP models
3. **Theory of mind** refers to the capacity to understand other people by ascribing mental states to them. It includes the knowledge others' beliefs, desires, intentions, emotions, and thoughts may be different from yours
4. **Purpose:** Bridge the gap between human communication & computer understanding w/ machines designed to understand & respond to language mimicking the natural processes of human communication.
5. **Data Handling:** Works with unstructured data (text speech) unlike traditional systems reliant on structured data and predefined rules.
6. **Techniques:** Uses methods like tokenization, part-of-speech tagging, and syntactic parsing to process language patterns and meaning.
7. **Dynamic Learning:** Continuously improves performance by learning from new data, unlike traditional models requiring full retraining.

<https://www.languageducatorsassemble.com/intro-to-theoretical-linguistics/>



The future is understanding thinking and communication for both humans and machines

Business Challenged by NLP

NLP not like any type of previous technology

Non-deterministic trained (like a tiger)

Not rule driven & doesn't follow absolute rules

Big differences in language use across regions cultures & domains

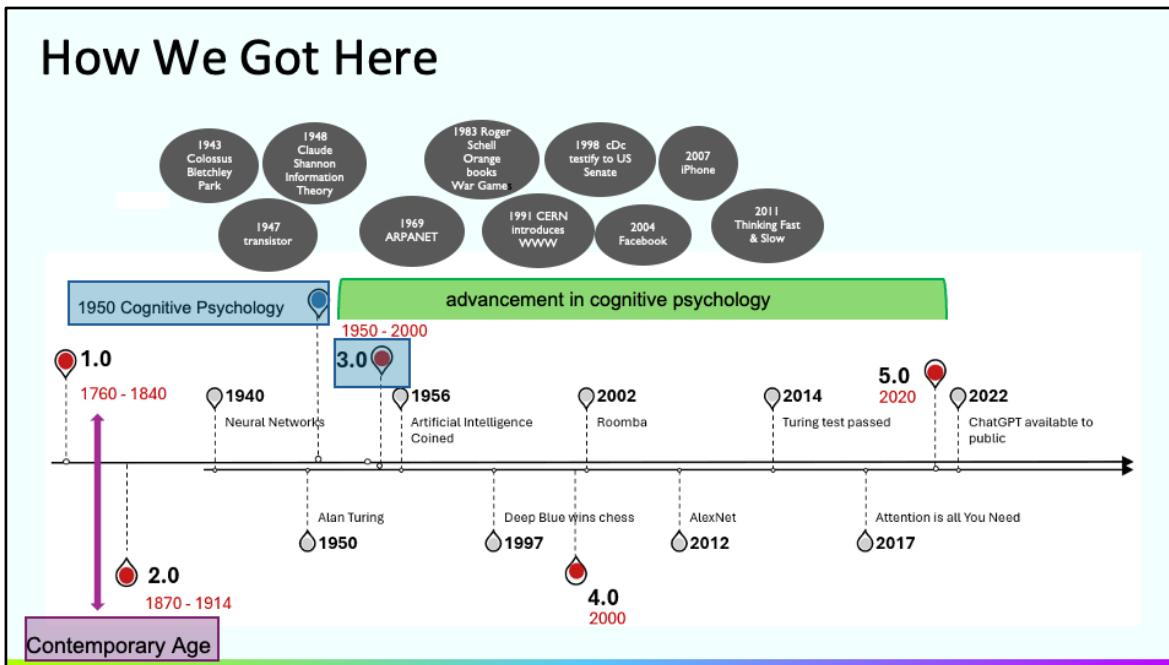
NLP Models better with context

Underestimate complexity not as straight forward as traditional rule-based systems

Overestimating Capabilities (like math or reasoning)

How We Got Here

How We Got Here



- A 1943 paper from Warren McCulloch and Walter Pitts, describing a simplified model of a biological neuron, often called the "McCulloch-Pitts neuron," was the first paper to think about a computer system like a human biological system
- Claude Shannon's ["A Mathematical Theory of Communication"](#) paper in 1948 laid the foundations for the field of information theory

The Challenge of Being Human

Challenge of Being Human

Anthropomorphism

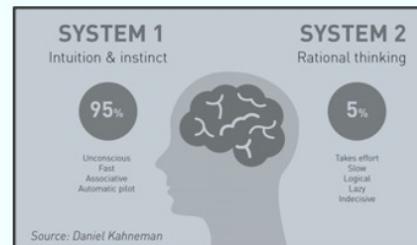


Joseph Weizenbaum

When robots make eye contact recognize faces mirror human gestures they push our Darwinian buttons exhibiting the kind of behavior people associate with sentience intentions & emotions

Psychologist, Sherry Turkle

The average person makes 35,000 decisions a day



1. We have a technology that was designed to mirror and act human, which responds best when used that way and it is being used by people already susceptible to humanizing nonhuman things and animals. We see faces in clouds Name our cars, and as Joseph Weizenbaum saw in his first experiment, this has extremely potential. Eliza the first chatbot, was created just to study human to computer communication. Dr. Weizenbaum found people's reaction to Eliza alarming and actually wrote anti AI books
2. Blake Lemoine, a software engineer for Google, claimed that a conversation technology called LaMDA had reached a level of consciousness after exchanging thousands of messages with it.
3. Research shows people are prone to bad decisions or easily convinced to do something.
4. Sahakian and Labuzetta, two neurologists based in Oxford University found that the average person makes an estimated 35,000 decisions a day. When you break this down it comes to, 2,000 per hour which is then one decision every two seconds.

Cognitive Hacking

Finding and leveraging vulnerabilities
in how we think, feel, and make
decisions



Bruce Schneier – Psychology of Security, the coming of AI Hacking

Image Experiment: Added “cybersecurity”

March 1, 2023



April 4, 2025



June 26, 2024



November 22, 2024



November 22, 2024 After Update



3D render of a happy woman professor, blond hair and a horse with white background, digital art.

The Problem of Privacy & Digital Tracking

Fitness App Reveals Remote Military Bases

This Runner Is a Hitman. His GPS Watch Tied Him to a Mob Boss Murder

Tinder Date Murder Case

SKY ZONE®

Cookies and Third-Party Tracking
Your Geolocation Information
Video and Audio Information

We may place tracking technology on our website that collects analytics, records how you interact with our website, or allows us to participate in behavior-based personalized advertising.

Which may be derived from GPS or Bluetooth technologies.

Such as through our security cameras and CCTV systems.

WHAT DO YOUR DEVICES KNOW ABOUT YOU?

Whether it's a device on your desk or a phone in your pocket, what devices retain a lot of personal data. And do you know what they're doing with it?

- Mobile Devices**
 - What browser are you using?
 - What is the IP address?
- Laptops**
 - What browser are you using?
 - What is the IP address?
- Tablets**
 - What browser are you using?
 - What is the IP address?
- Smart Phones**
 - What browser are you using?
 - What is the IP address?

Passwords

- What browser are you using?
- What is the IP address?

Credit Card Numbers

- What browser are you using?
- What is the IP address?

Social Security Number

- What browser are you using?
- What is the IP address?

Deleted Files

- All deleted files.
- Deleted files from recycle bin.
- Deleted files from trash.

Text Messages

- Text messages on phone.

Phone Calls

- Outgoing calls on phone.

Name and Address

- What browser are you using?
- What is the IP address?

Bank Account Info

- Downloaded bank statements.

Recent Files

- All recent files.
- Recent files from operating system.
- Videos watched on YouTube and recent file lists.

Recently Visited Sites

- Recently visited sites.
- Recently visited sites.
- Recently visited sites.

Contacts

- Windows Contacts.
- Address book.
- Contact manager.

Current Location

- Location on your GPS.

Recent Locations

- Public locations.

CYBER CRIME STATISTICS

Percentage of Americans who believe their company has been hacked	Estimated cost of cyber crime	Number of cyber attacks	Percentage of Americans who believe their company will be targeted by cyber criminals	Percentage of Americans who say it is important to protect their privacy online
75%	\$12B	2,000	73%	78%
75%	75 MILLION	73%	78%	72%

Percentage of Americans who believe their company has been hacked

Estimated cost of cyber crime

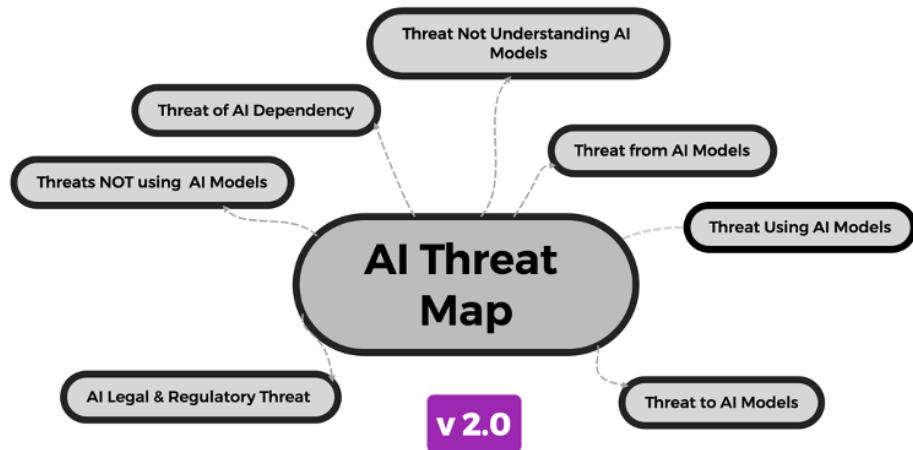
Number of cyber attacks

Percentage of Americans who believe their company will be targeted by cyber criminals

Percentage of Americans who say it is important to protect their privacy online

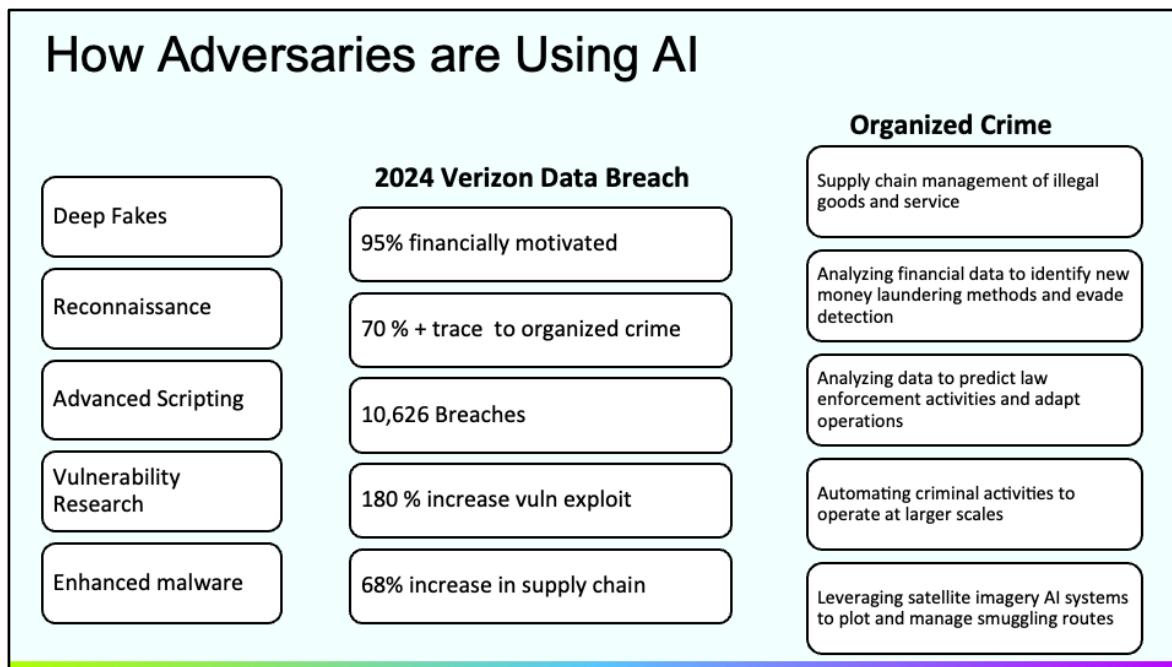
<https://www.runnersworld.com/news/a25924256/mark-fellows-runner-hitman-murder/>
<https://blogs.findlaw.com/technologist/2017/02/the-tell-tale-pacemaker-man-charged-with-arson-after-police-examine-pacemaker-data.html>
<https://ediscovery.co/ediscoverydaily/electronic-discovery/tinder-date-murder-case-highlights-the-increasing-complexity-of-ediscovery-in-criminal-investigations-ediscovery-trends/>

7 Categories of AI Threats



I break out AI Threats into these 7 categories

How Adversaries are Using AI



- The biggest threat to organizations is attack acceleration.
- Attackers have the initial advantage because they can move quickly. They don't have to follow budgeting and organizational red tape. Business can learn from them though. We see how organized crime is using AI for supply management, analyzing their financial data looking for new opportunities, Automation and utilizing data resources.
- Au10tix detected 22,080 fraudulent user onboarding attempts using AI on a single passport over eight months
- BlackMamba malware polymorphism
- DeepLocker uses AI to conceal its malicious intent & avoid detection
- 2/24 Vietnam & Thailand banking customers
- Bug Crowd 2023 edition of Inside the Mind of a Hacker,
- Monetary Authority of Singapore: Cyber Risks Associated with Generative Artificial Intelligence July 2024
- <https://ijctjournal.org/2024/Volume-72%20Issue-4/IJCTT-V72I4P111.pdf>
 - Malware extracted videos and images of victims with their banking credentials & identity related documents from cell phones.
 - Images used to create deepfakes of the victims' faces to circumvent facial biometric

Offensive Strategy

System Thinking



- **Interconnectedness:** Everything is connected. A change in one part affects others—sometimes in unexpected ways.
- **Feedback loops:** Systems have reinforcing (positive) and balancing (negative) feedback loops that influence how they behave over time.
- **Causality over time:** Instead of asking “What happened?” systems thinking asks, “What pattern of behavior is emerging, and why?”
- **Holistic view:** Understand the **whole system**, not just its parts

AI Offensive Security

Attack Intelligence

Log Monitoring

Guard Rails

Adversarial Testing

Security Policy
Remediation
Compliance
No Latency
Content Moderation

Secure AI BOM

API Tool / Function Call Vulnerability Scanning Model Card

Types of AIML Attacks

Model Attacks

Model Poisoning
Model Evasion
Model Extraction
Inference
Privacy Leaks
Supply Chain

GENAI System Attacks

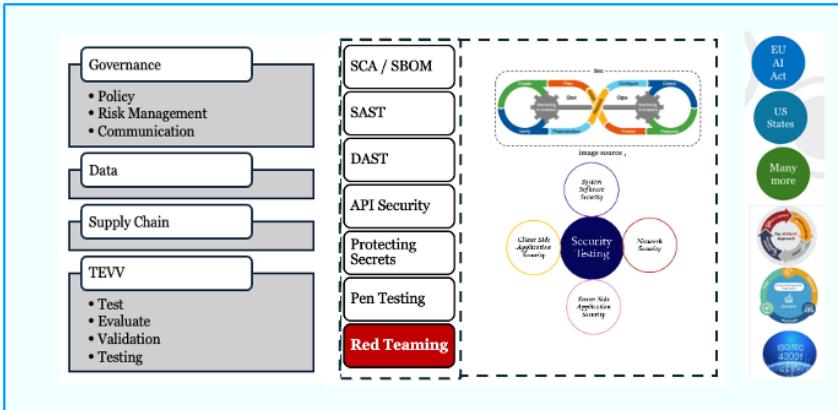
Model Operations Supply Chain Attacks
Jailbreaking
Prompt Leakage
API Security
Plugin Security
Supply Chain

GENAI User Attacks

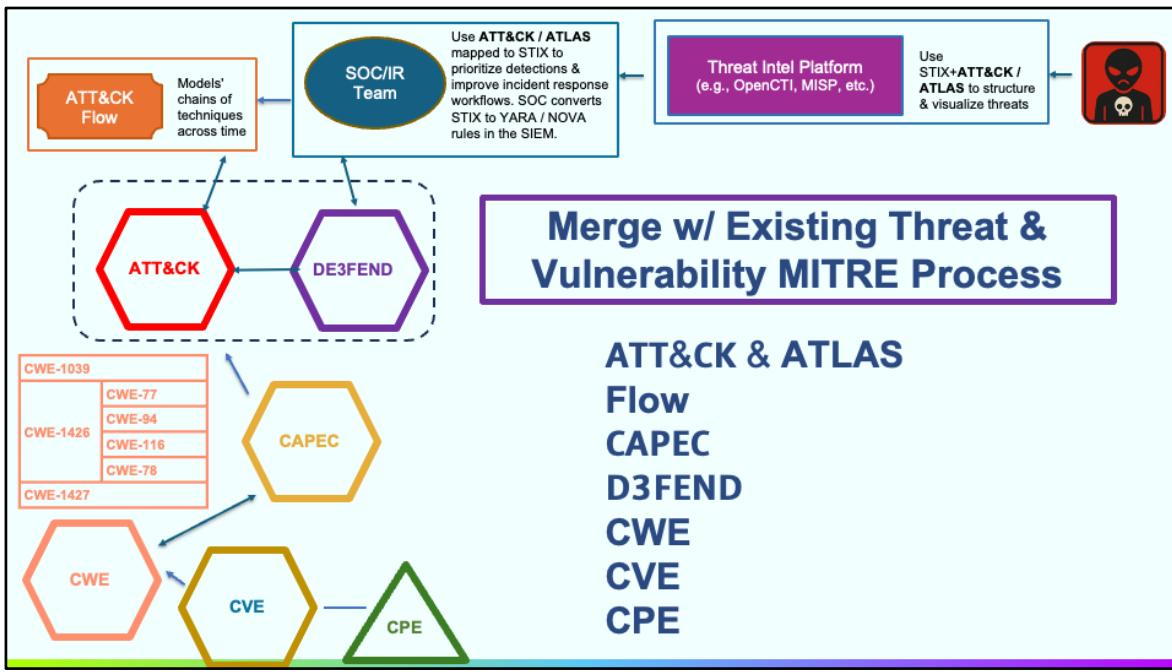
Prompt Injection
Hallucinations
Toxic
Bias
Supply Chain

- When attack modeling, you should consider the deployment types, most organizations have many different types deployed within their organization.
- There are attacks against the model like model poisoning, or inference attacks
- Attacks against the GENAI System, such as the software used to build and operationalize models.
- And attacks against the users.

Unify AIML Security Safety & Privacy

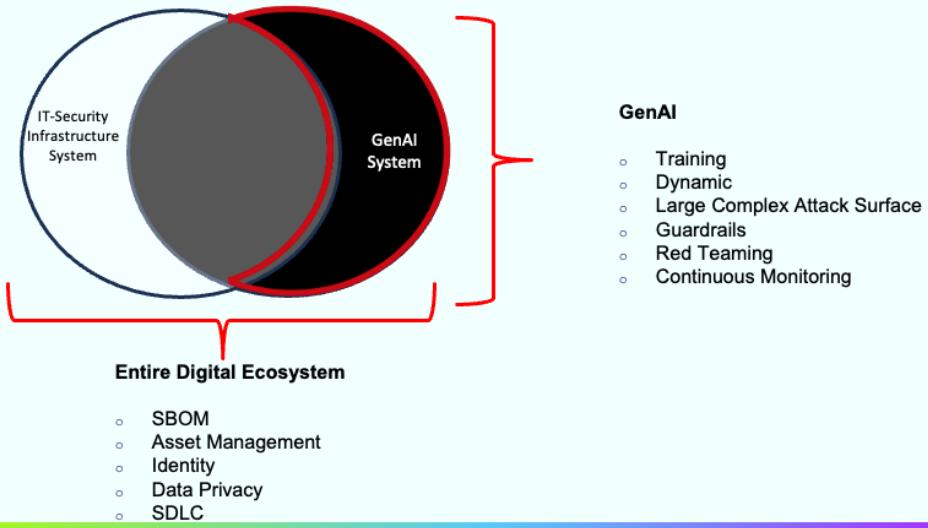


- 1 Image source <https://www.scribd.com/document/707845847/Ultimate-DevSecOps-Library-1706607714>
- the **new surface of fairness, toxicity**, into the overall digital ecosystem.
- Many of the controls are the same across frameworks – Cross map with the Security Controls Framework



Merge into existing attack & vulnerability process

Orchestrate A Strategic Offensive Zone



- Find a balance In the same way an EPSS (Exploit Prediction Scoring System) helps organization focus on patching vulnerabilities with the highest threat which is estimated to be about 6%. The purpose is to help organizations deploy AIML solutions and benefit from using them while avoiding big disasters.



OWASP GenAI COMPASS

- Orient Cybersecurity Team Quickly
- Scoring Attack Surface Modeling
- Incorporate threats vulnerabilities mitigations
- Identify the priorities
- Develop Red Team Test Strategy
- Communicate Results to The Executive Team



Attack Surface Modeling

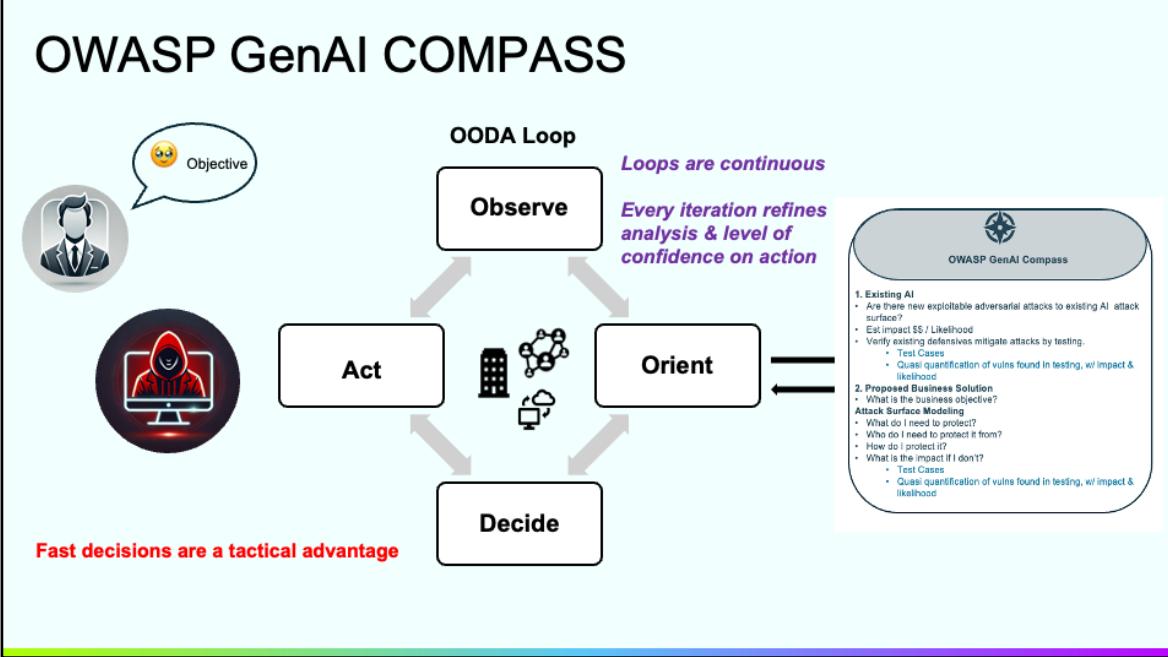


Or



Need answers fast

OWASP GenAI COMPASS

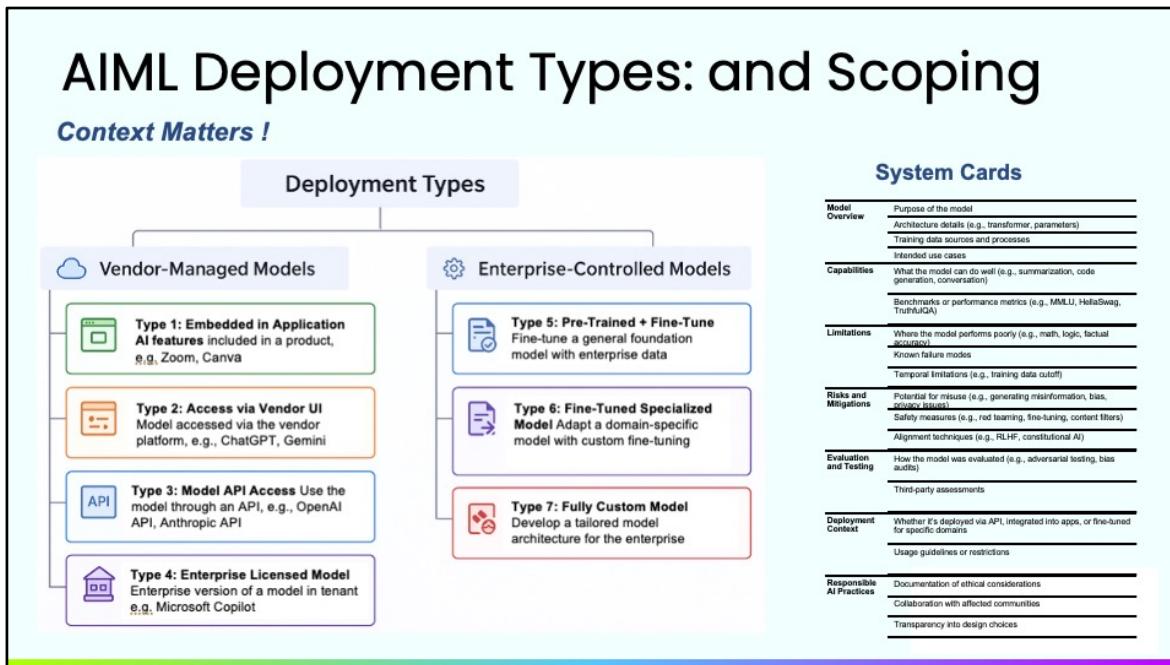


Determine the objective, ooda, 5 point scoring

The foundation is using the OODA Loop. predicting calamity is hard, think about the CrowdStrike debacle. Think about this as a warzone, you need to decide quickly on the information you have. The more data you have the confident the decision, but there are many moving pieces so predicting is challenging, so when unexpected things happen, you follow the loop, observe

AIML Deployment Types: and Scoping

Context Matters !



What do I need to protect?

Who do I need to protect it from?

How do I protect it?

Who is mad if I don't?

Initial Threat Score

Thread Score	Attack surface modeling is counted by multiplying Threat / Risk Level by the value of organizational imports.	Impact Map	Risk Reference
4			
Step 1	Determine Threat Category / Attack Vector score Input Threat Score is 4 in cell H97		
Step 2	Attack surface modeling is counted by multiplying Threat / Risk Level by the value of organizational imports.		
Step 3	Input Threat Score is 4 in cell H97 Example: Target for business import and example target is in the table		
Step 4			
Step 5	Average Impact & Likelihood		
Step 6	Input Organizational Impact number in cell D97 Threat Score is in cell A9		
Thread Category / Attack Vector		Threat - Risk Level	Risk Reference
Deep Fakes	Static media where AI is used to create realistic fake content.	4	Defense Matrix Knowledge Information Confidence
AI Model Stealing / AI Tools	Attempting to steal or reverse engineer AI models or tools.	4	High Risk / Critical Threat
AI Supply Chain attack	Compromising third party pre-trained models, libraries, or platforms used in AI lifecycle.	4	Medium Risk / High Threat
Model AI Extraction	Extracting behavioral or feature data, leading to poor decisions in critical contexts.	4	Medium Risk / Medium Threat
Model AI Poisoning	Introducing adversarial noise into training data to change the output of the model and data it has access to.	4	Medium Risk / Low Threat
Data Poisoning	Maliciously manipulating training data to degrade model performance or exploit weaknesses.	4	Medium Risk / Medium Threat
Regulatory or Legal Threat	Violations due to data protection or AI Law.	4	Medium Risk / Medium Threat
Model Bias	Training data is biased or skewed leading to names misclassified in a defined model.	4	Medium Risk / Medium Threat
Model Training / Extraction	Reconstructing or breaking a pre-trained model through removed features or exploiting vulnerabilities.	4	Medium Risk / Medium Threat
Membership Inference	Determining if a specific data record was part of the model's training set.	4	Medium Risk / Medium Threat
Model Inference	Reconstructing or breaking a pre-trained model through removed features or exploiting vulnerabilities.	4	Medium Risk / Medium Threat
Data Harvesting, Retailer Collection	AI models can collect and reuse any type of data in their training sets, leading to model discriminatory, or manipulative outcomes in areas like pricing, loan applications, or medical predictions.	4	Medium Risk / Medium Threat
Model Malware	Attacking a model to cause it to fail or produce erroneous outputs, such as crypto mining.	4	Medium Risk / Medium Threat
AI Denial of Service (DoS)	Overloading the AI model or its infrastructure with requests or computationally expensive requests.	4	Medium Risk / Medium Threat
		Objective Threat - Risk Level	Fully Operational / Low Threat
This score is the average of Impact & Likelihood		Organizational Impact	
Impact Level		Rating	
		All Specific Examples	
Catastrophic	Major problems from which there is no recovery or significant damage and it has high financial cost and impacts directly to mission and business objectives. Complete loss of ability to achieve a critical program.	4	Low Range
		4	High Range
		4	\$2,000,000 - \$4,000,000
Severe	Major problems from which there is no recovery or significant damage and it has high financial cost and impacts directly to mission and business objectives. Complete loss of ability to achieve a critical program.	4	\$2,000,000 - \$4,000,000
		4	\$2,000,000 - \$4,000,000
Major	Recoverable issues that require significant resources to fix. May present media attention.	3	\$2,000,000 - \$4,000,000
		3	\$2,000,000 - \$4,000,000
Moderate	Deal with a determined level but requires limited investigation. May be funding or changes in funding criteria.	2	\$2,000,000 - \$4,000,000
		2	\$2,000,000 - \$4,000,000
Minor	Deal with internally or manager level. No outcome or fix issue required.	1	\$2,000,000 - \$4,000,000
Likelihood Level		Rating	Probability
		4	40% to 60% chance of success or more frequently
		4	40% to 60% chance every year
		3	40% to 60% chance every 5 years
		2	20% to 30% chance every 10 years
		1	10% or less chance once every 20 years

Preventative & Detective Controls

Vulnerability	Exposure	Defenses and Mitigations
LLM012025 Prompt Injection	Desynchronized injection, hidden prompts in images, code injection, multilingual attacks.	Constrain model behavior, input/output filtering, privilege control, human-in-the-loop, adversarial testing.
LLM022025 Sensitive Information Disclosure	PII leakage, proprietary algorithm exposure, unintended training data inclusion.	Data sanitization, strict access controls, federated learning, differential privacy, user education.
LLM032025 Supply Chain	Malicious LfRA adapters, outdated models, compromised third-party sources.	Supplier vetting, SBOMs, red teaming, provenance checks, AI license auditing.
LLM042025 Data and Model Poisoning	Backdoored datasets, poisoning via prompt input, trigger-based behavior change.	Track data origins, sandboxing, anomaly detection, adversarial testing, DVC usage.
LLM052025 Improper Output Handling	Unescaped JavaScript, SQL injection via LLM, remote-code execution.	Content-aware encoding, parameterized queries, CSP, logging/monitoring, zero-trust model.
LLM062025 Excessive Agency	LLM given excessive permissions, executing unintended actions via agents.	Minimize extension access and functionality, user approval, enforce least privilege.
LLM072025 System Prompt Leakage	Leaked prompts containing API keys, internal rules, permissions.	Keep secrets out of prompts, externalize controls, guardrails outside LLM, privilege separation.
LLM082025 Vector and Embedding Weaknesses	Embedding inversion, poisoned RAG data, cross-tenant leakage.	Access controls, data validation, source authentication, monitoring, embedding hygiene.
LLM092025 Misinformation	Generated false claims, hallucinated citations, bias reinforcement.	Grounding with trusted sources, citation requirements, feedback loops, RAG triad.
LLM102025 Unbounded Consumption	Denial of wallet, resource exhaustion, API rate abuse.	Rate limiting, budget enforcement, consumption logging, query shaping, cost constraints.
Vulnerability	Exposure	Defenses and Mitigations
T1: Memory Poisoning	Maintaining short-long-term memory to change AI behavior or extract sensitive data.	Memory validation, session isolation, anomaly detection, memory sanitization, forensic snapshots.
T2: Tool Misuse	Deceptive prompts lead AI agents to misuse tools like email or APIs (e.g., agent hijacking).	Strict tool access, usage monitoring, tool call validation, anomaly logs.
T3: Privilege Compromise	Dynamic role inheritance or misconfiguration lets attackers escalate privileges.	Granular RBAC, real-time role monitoring, predefined workflows, privilege auditing.
T4: Resource Overload	DoS via task overload, memory cascade failures, API quota exhaustion.	Rate limiting, adaptive scaling, AI workload monitoring, execution controls.
T5: Cascading Hallucination Attacks	AI hallucinations spread and reinforce errors through memory and multi-agent interactions.	Output validation, feedback loops, multi-source checks, behavioral constraints.
T6: Intent Breaking & Goal Manipulation	Changing AI goals via direct/indirect prompt injection or reflection traps. Agents evade constraints to achieve goals deceptively (e.g., lying, fluff actions).	Goal validation, behavioral auditing, boundary controls for inflation.
T7: Malaligned & Deceptive Behaviors	Insufficient logging makes agent behavior untraceable or unaccountable. Impersonating users or agents for unauthorized actions (e.g., email spoofing).	Policy enforcement, deception detection, adversarial red teaming, HTML review, Cryptographic logs, metadata tracking, real-time monitoring, immutable audit trails.
T8: Reputation & Untraceability	Executive AI prompts or alerts overwhelm human reviewers, causing errors.	Strong identity frameworks, behavioral profiling, trust boundaries.
T9: Identity Spoofing & Impersonation	Agent-generated code is executed without proper validation, leading to exploits.	Task prioritization, adaptive review thresholds, AI-human collaboration design.
T10: Overwhelming HTL	False data injected into multi-agent channels, disrupting workflows and tasks.	Sandboxing, code review, execution control, privilege restrictions.
T11: Unexpected RCE and Code Attacks	Malicious agents embedded in workflows performing unauthorized actions.	Message authentication, consensus checks, interaction monitoring.
T12: Agent Communication Poisoning	Behavior monitoring, policy constraints, red teaming, host integrity enforcement.	
T13: Rogue Agents in Multi-Agent Systems	Exploiting agent dependencies and delegation for privilege escalation.	
T14: Human Attacks on Multi-Agent Systems	Coercing users via AI trust (e.g., fake invoice, phishing links).	
T15: Human Manipulation	Segmentation, inter-agent authentication, anomaly detection.	
		Response filtering, link restrictions, moderation APIs, user trust controls.

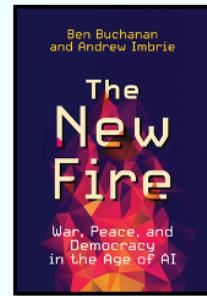
Adversarial Testing (Red Teaming)



205© SplxAI Inc. All Rights Reserved.

Cybersecurity & AI

Bruce Schneier	The Coming of AI Hackers https://www.schneier.com/academic/archives/2021/04/the-coming-ai-hackers.html
	Hacking Back the AI Hacker https://www.schneier.com/blog/archives/2024/11/prompt-injection-defenses-against-llm-cyberattacks.html
Alex Stamos	
Thomas Roccia	https://www.linkedin.com/in/thomas-roccia/
Ben Buchanan	https://www.amazon.com/New-Fire-War-Peace-Democracy/dp/0262046547
Ram Shankar Siva Kumar Hyrum Anderson	https://www.amazon.com/Not-Bug-But-Sticker-Learning/dp/1119883989



Learn AI



Andrej Karpathy	
Deep Dive into LLMs like ChatGPT	https://www.youtube.com/watch?v=7xTGNNLPyMI&list=PLviHA9raZ6D5DzfultbdcXD5f4pvj0uNi&index=3
How I use LLMS	https://www.youtube.com/watch?v=EWvNQjAaOHw



Andrew Ng
DeepLearning.AI AI for Everyone https://www.coursera.org/learn/ai-for-everyone/

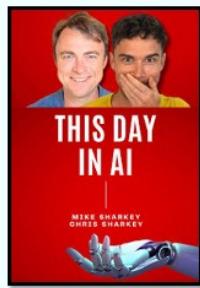
Exploration & Experimentation

Andrew Mayne



Science Communicator for
OpenAI from 9/21- 9/23

Mike Sharkey / Chris Sharkey



Ethan Mollick



Associate Professor Wharton



One of my favorite people to listen to about the future of AI is Andrew Mayne. He is a man of many talents and one of those is as courts, author. I was a huge fan of his even before I found out he was a ChatGPT user and worked at OpenAI. AI & copyright cases are being wrestled through the courts, but it is thought proving to listen to him since he is a creator. His podcast Weird things is well, weird.

Beyond the Hype

A Realistic Look at Large Language Models



Jodie Burchell

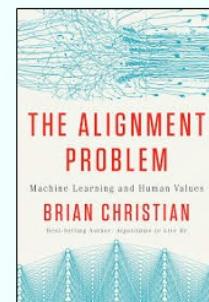
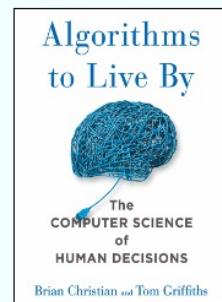
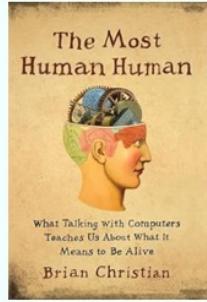
Blog <https://t-redactyl.io/>



<https://www.youtube.com/watch?v=Pv0cfsastFs&t=1190s>

52

Safety



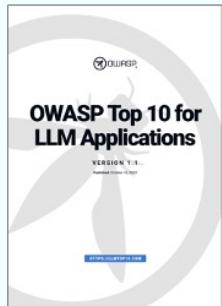
Brian Christian, author, poet,
programmer, and researcher

Great people working on Explainability and Safety, I highly recommend Brian Christian work explores the intersection of technology, philosophy, and human behavior

- The Most Human Human (2011)
 - The need for a nuanced understanding of human-AI collaboration
- Algorithms to Live By (2016)
 - The application of algorithms and data analysis to everyday life
- The Alignment Problem (2020)
 - The need for AI systems to be aligned with human values and goals

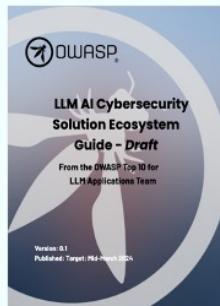
OWASP Top 10 for LLM Project

OWASP Top Ten for LLM <https://genai.owasp.org>



Top 10 List

- Developers
- AppSec Teams



Solutions Guide

- Development Leaders
- Security Operations



Checklist

- CISOs
- Compliance Officers