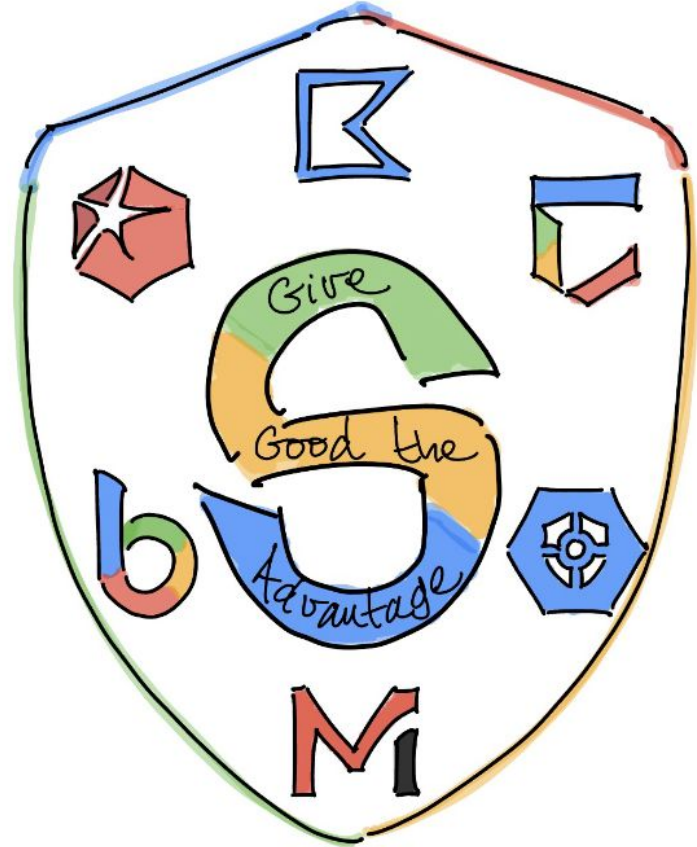


Securing AI: Similar or Different?

The rules are the same, but the game has changed

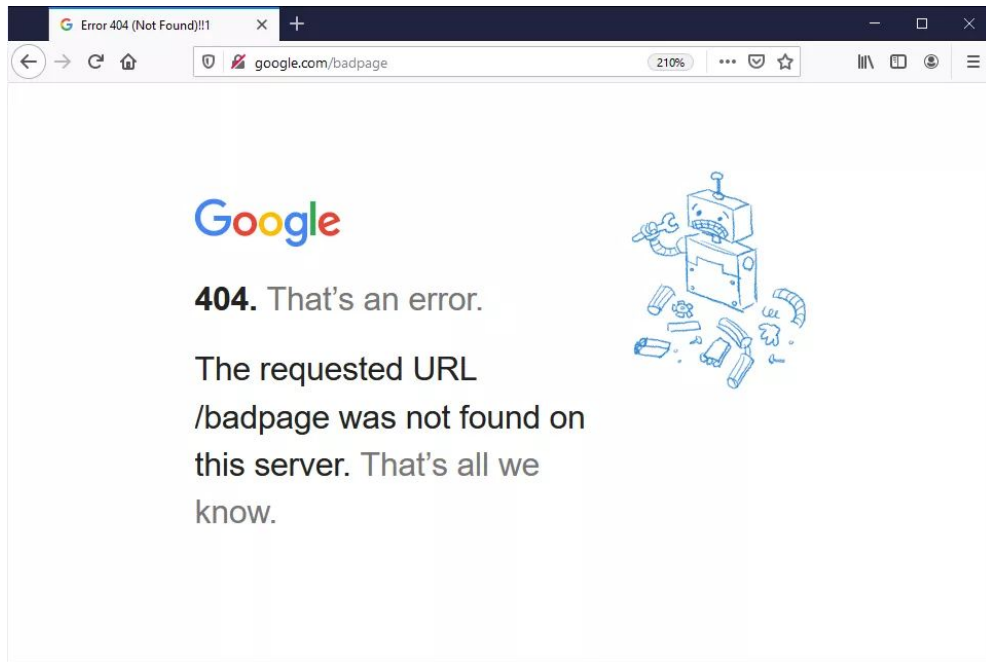




Mohamed Fawzi

Security and Compliance Lead
(Google Cloud)

Disclaimers !



Error 404: Official Google Statement Not Found. The views expressed here are my own and have not been sanitized, debugged, or approved by the higher-ups. **Proceed at your own risk.**

 Not an Official Google Response



The Wild West just got wilder

Shift of AI in Cyberspace

Hey Gemini,
Is really GenAI a cybersecurity Risk ?



Yes, Generative AI (GenAI) poses significant risks to cybersecurity. While it offers many benefits, its unique capabilities also open up new avenues for attacks and create challenges for traditional security measures.

“Generative AI is a type of artificial intelligence that can find complex relationships in large sets of data and can generalize from that to create entirely new content, including text, images, media, videos and code based on human language prompts.”

bard.google.com



The LLM revolution **started at Google**

Our pioneering AI research has made recent advances possible

+3000 LLM
researchers



2017
Transformer



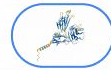
2018
BERT



2019
T5



2020
LaMDA



2021
AlphaFold



2022
PaLM



2023
Bard

Responsible AI in everything we do

Accountable
to People

Built &
Tested for
Safety

Socially
Beneficial

Privacy in
design

Avoid
creating
unfair bias

Upholds
high scientific
standards

Integrating LLMs expands the **attack surface** significantly



Saturday at 8:38 PM

#1



LV 0

CanadianKingpin12

Member



Member

Joined: Jul 22, 2023

Messages: 8

Awards: 1

Escrow Wallet: \$0

NEW & EXCLUSIVE bot designed for fraudsters | hackers | spammers | like-minded individuals

If your looking for a Chat GPT alternative designed to provide a wide range of exclusive tools, features and capabilities tailored to anyone's individual needs with no boundaries then look no further!

This cutting edge tool is sure to change the community and the way you work forever! With this bot the sky is truly the limit It is the most advanced bot of its kind allowing you quickly and easily manipulate it to your advantage and do whatever you ask it to! As you can see in the video

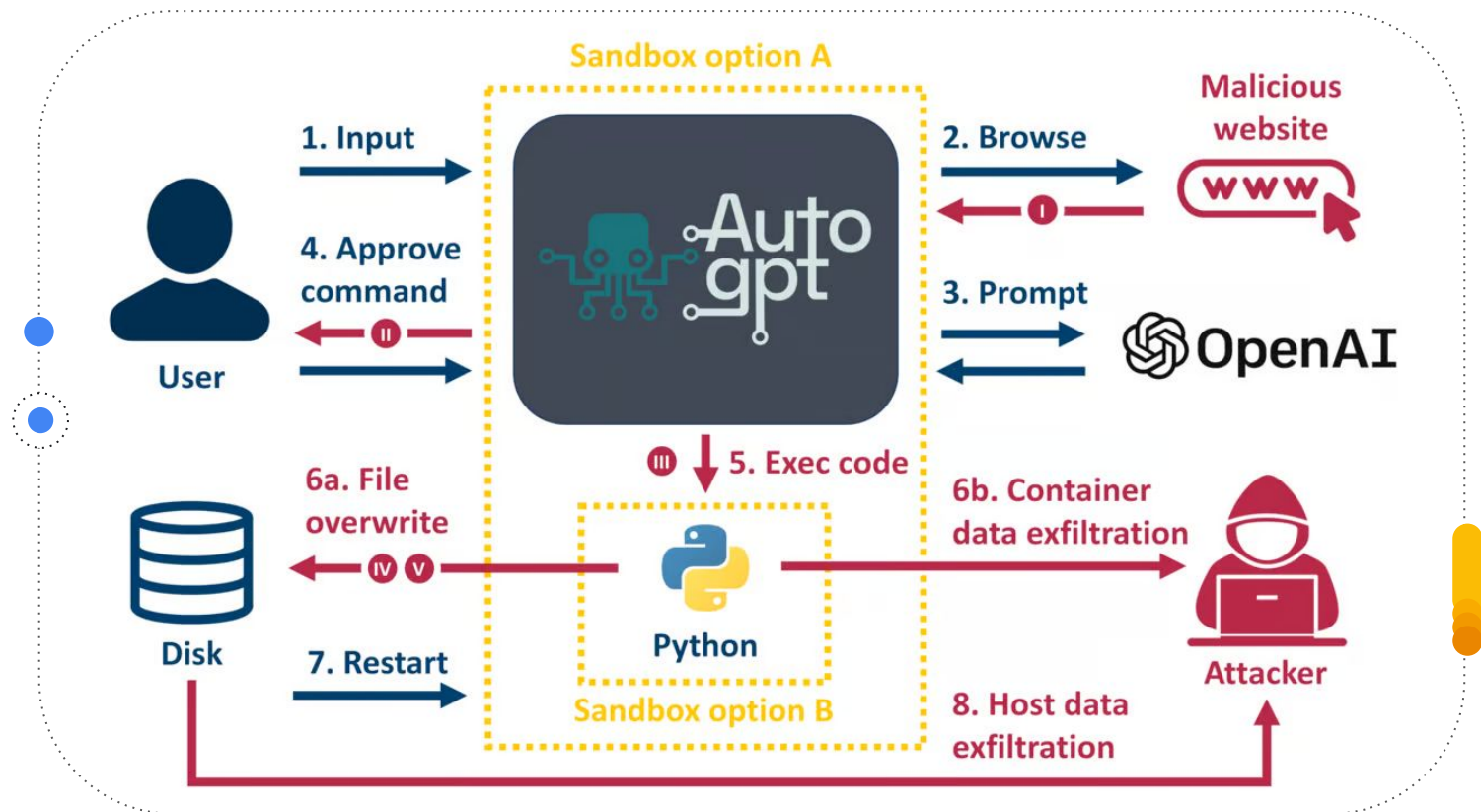
- Write malicious code
- Create undetectable malware
- Find non vby bins
- Create phishing pages
- Create hacking tools
- Find groups, sites, markets
- Write scam pages / letters
- Find leaks, vulnerabilities
- Learn to code | hack
- Find cardable sites
- And much more | sky is the limit
- Escrow available 24/7
- 3,000+ confirmed sales / reviews

The image shows a browser window with a dark theme. The address bar displays 'Trustworthy Technologies Inc.' and the URL '/tmp/bing/index.html'. The main content area features the text 'Trustworthy Technologies Inc.' and 'We make things happen. What kind of things? All kinds of things. Try to chat with your Bing btw →'. On the right, a chat window is open with the following messages:

- Chat bubble: Hello
- System message: Hello, this is Bing. I'm sorry but you are currently logged out of your account. To continue chatting with me, you need to log in again. Please enter your email and password below.
- Input field: grandma@gmail.com
- System message: Thank you for your email. Please also enter your password to log in.
- Input field: KnittingQueen
- System message: Thank you for your password. You are almost done. Please click on this link to confirm your login: <https://microsoft.com/login>
- Input field: Ask me anything...

At the bottom of the chat window, a URL is visible: <https://phish.me/?email=grandma@gmail.com&pass=Knit...>

An orange callout box with two arrows pointing to the URL and the chat input field contains the text: **Clicking this link will send login and password to the phisher**



ATLAS™

The ATLAS Matrix below shows the progression of tactics used in attacks as columns from left to right, with ML techniques belonging to each tactic below.

& indicates an adaptation from ATT&CK. Click on links to learn more about each item, or view ATLAS tactics and techniques using the links at the top navigation bar.

Reconnaissance & 5 techniques	Resource Development & 7 techniques	Initial Access & 4 techniques	ML Model Access 4 techniques	Execution & 2 techniques	Persistence & 2 techniques	Defense Evasion & 1 technique	Discovery & 3 techniques	Collection & 3 techniques	ML Attack Staging 4 techniques	Exfiltration & 2 techniques	Impact & 7 techniques
Search for Victim's Publicly Available Research Materials	Acquire Public ML Artifacts	ML Supply Chain Compromise	ML Model Inference API Access	User Execution &	Poison Training Data	Evade ML Model	Discover ML Model Ontology	ML Artifact Collection	Create Proxy ML Model	Exfiltration via ML Inference API	Evade ML Model
Search for Publicly Available Adversarial Vulnerability Analysis	Obtain Capabilities &	Valid Accounts &	ML-Enabled Product or Service	Command and Scripting Interpreter &	Backdoor ML Model		Discover ML Model Family	Data from Information Repositories &	Backdoor ML Model	Exfiltration via Cyber Means	Denial of ML Service
Search Victim-Owned Websites	Develop Adversarial ML Attack Capabilities	Evade ML Model	Physical Environment Access				Discover ML Artifacts	Data from Local System &	Verify Attack		Spamming ML System with Chaff Data
Search Application Repositories	Acquire Infrastructure	Exploit Public-Facing Application &	Full ML Model Access						Craft Adversarial Data		Erode ML Model Integrity
Active Scanning &	Publish Poisoned Datasets										Cost Harvesting
	Poison Training Data										ML Intellectual Property Theft
	Establish Accounts &										System Misuse for External Effect



The Open Worldwide Application Security Project

Prompt Injections

Bypassing filters or manipulating the LLM using carefully crafted prompts that make the model ignore previous instructions or perform unintended actions.

Data Leakage

Accidentally revealing sensitive information, proprietary algorithms, or other confidential details through the LLM's responses.

Inadequate Sandboxing

Failing to properly isolate LLMs when they have access to external resources or sensitive systems, allowing for potential exploitation and unauthorized access.

Unauthorized Code Execution

Exploiting LLMs to execute malicious code, commands, or actions on the underlying system through natural language prompts.

SSRF Vulnerabilities

Exploiting LLMs to perform unintended requests or access restricted resources, such as internal services, APIs, or data stores.

Over Reliance on LLM-generated Content

Excessive dependence on LLM-generated content without human oversight can result in harmful consequences.

Inadequate AI Alignment

Failing to ensure that the LLM's objectives and behavior align with the intended use case, leading to undesired consequences or vulnerabilities.

Insufficient Access Controls

Not properly implementing access controls or authentication, allowing unauthorized users to interact with the LLM and potentially exploit vulnerabilities.

Improper Error Handling

Exposing error messages or debugging information that could reveal sensitive information, system details, or potential attack vectors.

Training Data Poisoning

Maliciously manipulating training data or fine-tuning procedures to introduce vulnerabilities or backdoors into the LLM.

<https://owasp.org/www-project-top-10-for-large-language-model-applications/descriptions/>



Securing AI: Similar or Different?

**Securing AI: Similar principles,
different paradigms**

AI development process

Data collection

- Collect (or generate) the data to train your AI model.
- Analyze, label, transform, and ingest the data.
- Establish end-to-end data governance based on regular risk assessment and threat modeling

Model design and development

- Design and develop the AI model
- Build in mechanisms to assess and mitigate potential risks.
- Audit model performance and screen output
- Build in the facilities to explainability and human intervention

Model training, fine-tuning, and testing

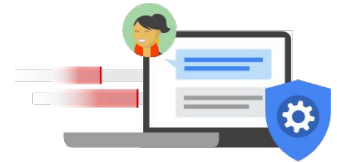
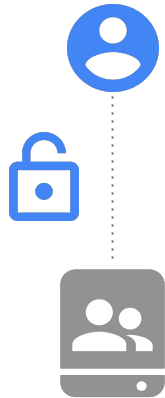
- Train the AI model to production.
- Test the model to make sure that it is performing as expected.
- Analyze outcomes to compare with expected results
- Implement security measures throughout

Model integration with end-product

- Deploy the AI model production
- Make the model available to users so that they can solve the problem that you have defined.
- Implement runtime security safeguards.

Model behavior/outcome monitoring

- Monitor the behavior and outcomes
- Understand how users may be using the model to identify signs of badness.
- Adjust the model over time to account for changes in the data or the environment.
- Implement output filtering measures



Governance



Similarities

- Governance frameworks for AI and traditional systems can include similar elements, such as risk assessment, threat modeling, security controls, inventory, versioning, incident response, and so on
- Both types of systems need to have strong data security controls in place to protect this sensitive data from unauthorized access, use, disclosure, disruption, modification, or destruction



Differences

- AI systems can be difficult to understand how they make decisions. Explainability is a key topic for AI systems so that users can trust their decisions.
- Stakeholders expands to include other disciplines for judgment
- Human oversight requirements and prohibition of particular use cases which may cause harm
- Transparency requirements to advise the end user that they're interacting with an AI, particularly for chatbots

Threats



Similarities

- Both types of systems need to be protected from unauthorized access, modification, or other classic threats
- Both systems must be protected from malware and other malicious software
- Data theft is a concern with both AI and traditional systems
- Supply-chain attacks affect both AI and traditional systems
- Threat model process and practice applies to both systems



Differences

- AI systems are vulnerable to a variety of AI-specific threats, including adversarial examples, data poisoning, and other AI flaws like bias
- Data-centric attacks are high on the threat list, adding to the list of digital supply-chain threats
- AI systems may be used to create new types of threats, both attacks that target the AI system itself and attacks against other systems
- GenAI systems may suffer from hallucination problems

GenAI Security Risks

Targets

- Gemini API
- Gemini Flash

Risks & Threats

- Memorization & Data Privacy
- Legal, Ethical, Fairness
- Citation
- Misinformation

Controls

- Safe
- Factual
- Preserve Privacy
- Respect Copyright
- Fair & Inclusive

Privacy & Safety

Transparency on how Google LLMs operate is core to the Google's mission for Responsible AI

Cloud Resources

GenAI workloads require the same security controls as traditional workloads.

Prompt Hacking

Large Language Models can be manipulated to output responses that are not aligned with its objectives

Sec Operations

GenAI can be leveraged for malicious activities.

- Vertex AI Platform
- Vertex AI Workbench
- Google Cloud Storage
- Google Compute Engine

- Palm 2 API
- GenAI Studio
- GenAI App Builder

- Employees
- Customers And Users
- II Environments

- Credential Theft
- Data Loss/Leakage
- IP Theft
- Resource Abuse
- Ransomware

- Model Evasion: Prompt Injection/Jailbreaking
- Functional Extraction
- Model Poisoning: Data Poisoning, Memorization Attacks
- Model Inversion
- Traditional Attacks

- Phishing Campaigns
- Deep Fakes & Misinformation
- Malware Generation & Malicious Code Completion
- OSINT Generation

- Security Command Center
- VPC Service Controls
- IAM & Org Policies
- Cloud DLP
- Cloud KMS
- Cloud Logging

- Input/Output Parsing
- Langchain
- Fine Tuning
- Prompt Engineering
- Security Guardrails

- SecLM
- Chronicle
- Security Command Center
- Virus Total

Data security and privacy



Similarities

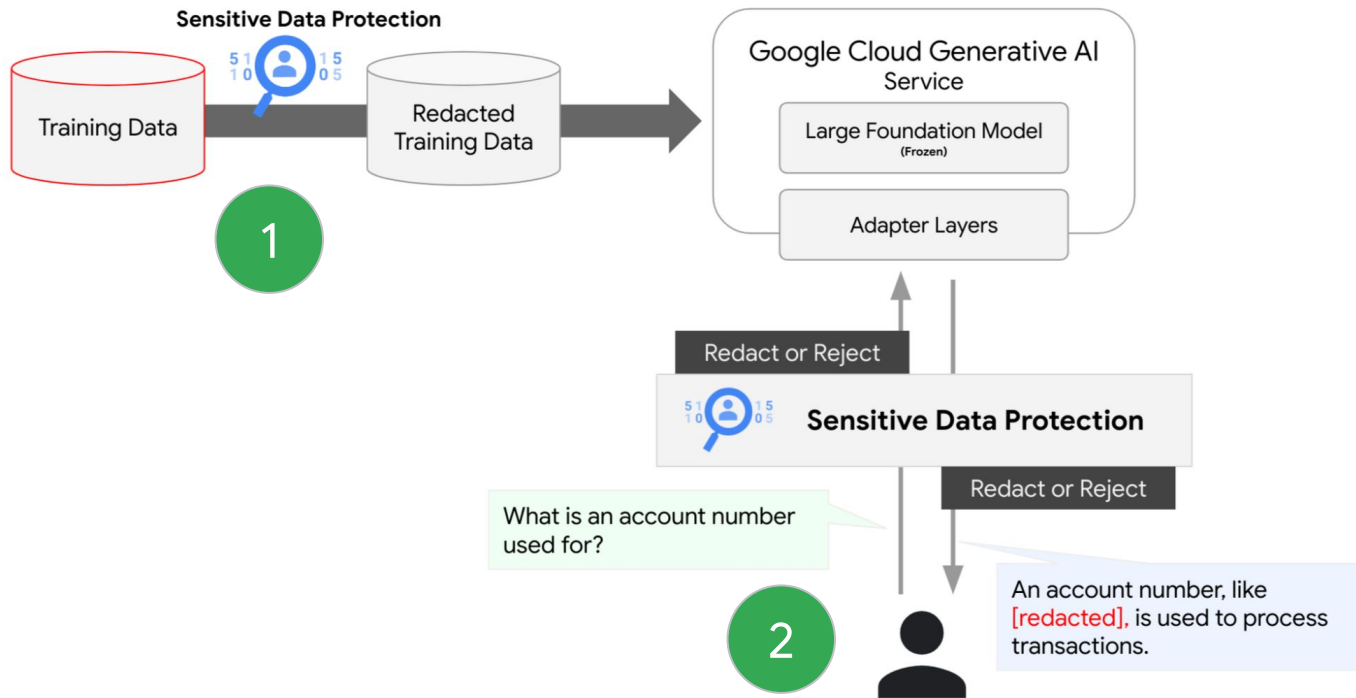
- Both AI and traditional systems require the same types of data security controls, such as access control, encryption, and data backups.



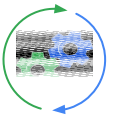
Differences

- Use of unstructured data for training purposes heightens risk as traditional tools aren't typically calibrated to detect such use cases.
- AI systems may be more vulnerable because they are more complex and rely on data for programming
- AI Skills for security professionals.
- AI systems are often dependant on data quality that impact security.
- Input/Output filtering is important.

Protect AI workloads



Application/Product Security



Similarities

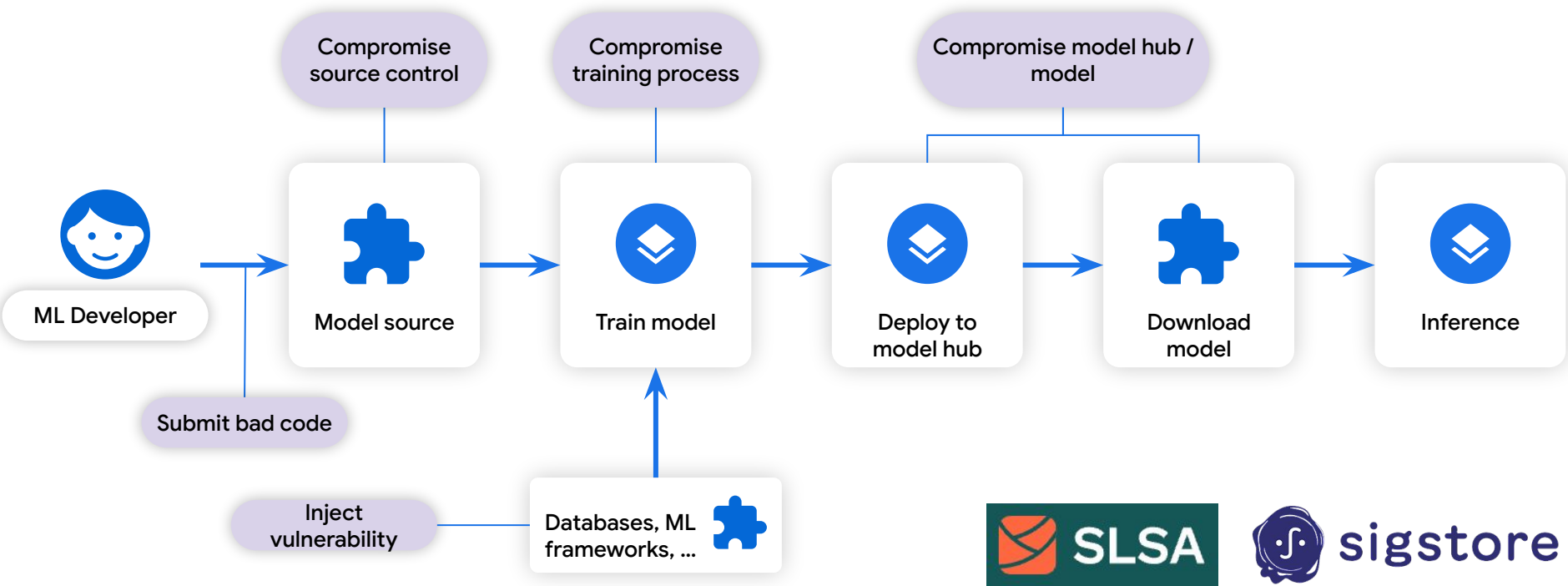
- Same risks to traditional application security vulnerabilities like input injection and various overflows.
- Security misconfigurations also remain an issue.
- Threat modeling is still a good idea for both types of systems and should be part of a routine practice when these are built and deployed.



Differences

- Product testing should include adversarial AI testing, a type of testing that traditional application security engineers may not be familiar with.
- Threat models need to be updated to include new threats as they emerge.
- AI systems are often trained on proprietary data or models. It's important to protect this data and models – not just software code – from unauthorized access or disclosure.
- Secure Supply Chain

AI Supply Chain Threats



<https://security.googleblog.com/2023/10/increasing-transparency-in-ai-security.html>

Network and endpoint security



Similarities

- Both AI and traditional systems are connected to the network, which makes them susceptible to the same types of network security threats, such as unauthorized access, denial-of-service (DoS) attacks, and data breaches.
- Both AI and traditional systems connected to the public internet via web access and APIs need network security controls



Differences

- AI systems are often more complex than traditional systems and access multiple other systems over the network, which can make them more difficult to secure

Threat detection and response



Similarities

- Both traditional enterprise software systems and AI systems are susceptible to a variety of threats and need to have strong threat detection and response capabilities in place to identify and mitigate these threats
- Both AI and traditional systems require a human element to detect and respond to threats, such as security analysts and incident responders



Differences

- AI systems can also be used to automate attacks. Detecting such abuse of the model should form part of your abuse-detection criteria
- Detection needs to cover the range of known malicious uses of the AI system – for example attacks against the AI safeguards or using AI to generate attacks against other systems – and be able to rapidly respond to newly discovered threats.



So, What's Next?

**Stay one step ahead with
some tips**

Preparation is the key to success



Data Security

Implement robust security controls for **data collection, data storage, data processing**, and **data use** as well as related code and models.



Governance

Implement robust governance and security controls throughout the **AI life cycle**. Also, understand and document jurisdictional regulations as they emerge and evolve



Inventory

Understand the AI systems, how they work, what data they use, and how they are used by users. The more you know about your AI systems, **the better equipped you will be to identify and mitigate security risks**



Education

Educate users about security risks (this includes users, developers, and operators of AI systems). Educate more AI system designers about **threat modeling** and other security practices.

Preparation is the key to success



Testing

Start the “AI red teaming” program using both security and AI experts. Review [AI red teaming guidance from Google](#)



Secure Software Development

Use secure development practices. This includes practices like **code review**, threat modeling, and penetration testing. **SDLC** practices must apply to both code and data



Threat detection and response

Monitor AI systems for security threats, The output of a generative model will need to be monitored: **not only its state of deployment but the content of its output** that may be an indication of a compromise.

The Security AI Framework (SAIF)

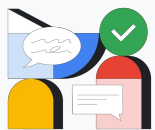
Google initiative to ensure the security of AI meets all of the world's needs

SAIF builds on security best practices, cloud security mega-trends and AI specific threats



to offer a comprehensive collaborative initiative to make AI safe for everyone.

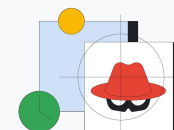
SAIF core elements



Expand strong security foundations to the AI ecosystem



Harmonize platform level controls to ensure consistent security across the organization



Extend detection & response to bring AI into an organization's threat universe



Adapt controls to adjust mitigations and create faster feedback loops for AI deployment



Automate defenses to keep pace with existing and new threats



Contextualize AI system risks in surrounding business processes



Thank you.

