



Guardians of Intelligence

Navigating Supply Chain Security for ML Applications

Dan Fernandez

Agenda

- Intro to ML Applications
- Understanding the Software Supply Chain
- Attack Vectors
- Safeguarding the AI Realm



Dan Fernandez

- Data Nerd
- Product Manager in Cybersecurity
- Instructional Associate @ GaTech



danielfernandez@infosec.exchange



@danielfernandez



/in/dafnz



danielfernandez.medium.com



hey@dafnz.com

Disclaimer

Opinions are my own and not my employers.



Presentation Resources

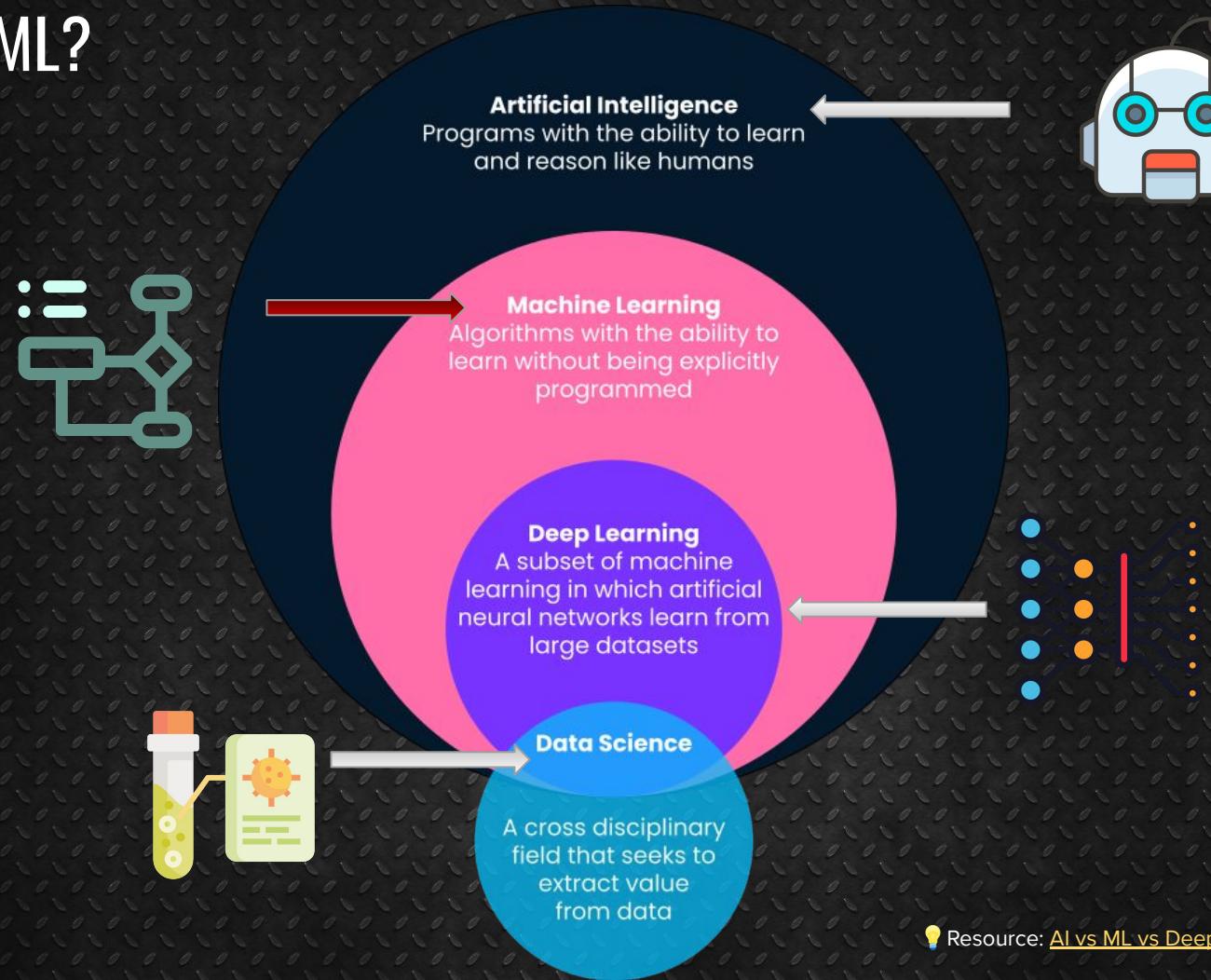


- The slides will be made available in Github immediately after the presentation.
- Slides contain references to content with additional information.
- Resource: Will look like this



Intro to ML Applications

What is ML?



What are the types of tasks do ML Models Perform?



Supervised Machine Learning

- Regression

The screenshot shows the Zillow website interface for Orlando, FL. At the top, there are navigation links for Buy, Rent, Sell, Home Loans, and Agent finder. Below that is a search bar with dropdown menus for Price, Beds & Baths, Home Type, and More. A 'Save search' button is also present. The main area features a map of Orlando with numerous red dots representing property locations. A specific boundary is highlighted with a blue line and labeled 'Schools'. Below the map is a list of five house listings:

- \$387,500**
3 bds, 2 ba, 1,343 sqft - House for sale
4814 Myrtle Bay Dr, Orlando, FL 32829
SAKE LOWE ORLANDO REALTY LLC
- \$339,900**
3 bds, 2 ba, 1,154 sqft - House for sale
7544 Panthera Ct, Orlando, FL 32822
FLOREN WEAVER TEAM OF ORLANDO, EXP REALTY LLC
- \$510,000**
4 bds, 3 ba, 2,850 sqft - House for sale
1665 Canoe Creek Dr, Orlando, FL 32824
EXP REALTY LLC
- \$350,000**
2 bds, 1 ba, 1,158 sqft - House for sale
3508 Route Rd, Orlando, FL 32821
EXP REALTY LLC
- \$543,900+**
4 bds, 3 ba, 2,361 sqft - New construction
DESTIN Plan, Wynwood Square
248 Hinson - Orlando East City

- Classification



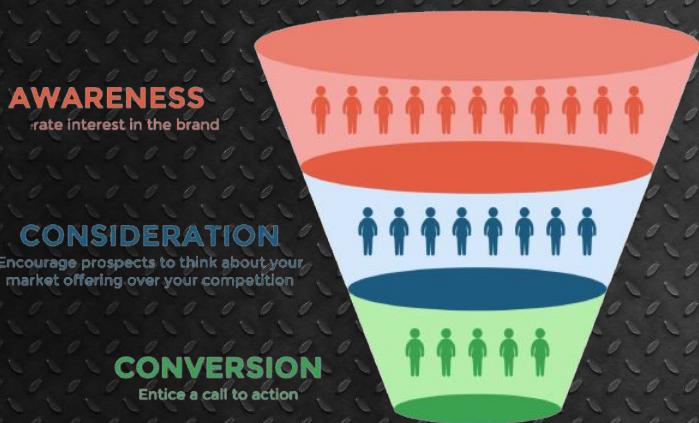
💡 Resource: [Supervised ML Tutorial](#)



What are the types of tasks do ML Models Perform?

Unsupervised Machine Learning

- Customer Segmentation



💡 Resource: [Unsupervised ML Tutorial](#)

- Recommender Systems

Google News

Search for topics, locations & sources

Home For you Following News Showcase U.S. World Local Business Technology Entertainment Sports Science Health

Full Coverage

OpenAI adds voice and image capabilities to ChatGPT

Sort Share

Top news

- ChatGPT users can now browse internet, Openai says 3 days ago
- How to Use ChatGPT's New Image Features 4 hours ago
- ChatGPT can finally access the internet in real time, but there's a catch 2 days ago
- Chatbots can now talk, but experts warn they may be listening too Yesterday
- The New ChatGPT Can See and Talk: Here's What It's Like. 3 days ago
- ChatGPT Can Now Talk to You—and Look Into Your Life 5 days ago

In-depth

- If you wouldn't take advice from a parrot, don't listen to ChatGPT: Putting the tool to the test 7 days ago

From Twitter

Why should we secure ML Applications?



Finance



Public Safety



Cybersecurity



Finance: Flash Crash +

REUTERS

World Business Markets Breakingviews Video More

COMMODITIES APRIL 23, 2013 / 5:58 PM / UPDATED 10 YEARS AGO

Analysis: False White House tweet exposes instant trading dangers

By Steven C. Johnson 6 MIN READ [f](#) [t](#)

NEW YORK (Reuters) - The upheaval in financial markets caused by a false report of explosions at the White House was brief, but its effect on traders who have come to rely on Twitter may last quite a bit longer.



The U.S. Capitol Building is pictured in Washington, February 27, 2013. REUTERS/Jason Reed

Volatile price moves ricocheted through markets for stocks, bonds, currencies and commodities around 1 p.m. EDT on Tuesday after a tweet purporting to be from the Associated Press said there had been two explosions at the White House and that President Barack Obama had been injured.

Source: [http://www.reuters.com/article/2013/04/23/us-usa-politics-explosions-idUSBRE93M0JL20130423](#)



Resource: [Flash Crash Simulator Whitepaper](#)

Entrepreneur

Sign In Subscribe [Search](#)

Business News

Eli Lilly Stock Plummets After Parody Twitter Account Says Insulin is Now Free

The pharma company felt the affects of a dangerous imitation account on Thursday.

BY EMILY RELLA • NOV 11, 2022 [Share](#)



Getty Images

Elon Musk's disastrous back-and-forth rollout of Twitter Blue has caused a multitude of communication problems between users and trolls, with fake news and information inadvertently going viral after users quickly see blue checkmarks and assume what they are reading is from a legitimate source.

Though much of the banter has been in good jest, one company is feeling the very real implications of the viral spreading of misinformation.

Eli Lilly and Company's stock plummeted on Thursday after an account posing to be the pharmaceutical company claimed that insulin would now be free in the United States.

Source: [http://www.entrepreneur.com/article/340000](#)

Disabling Physical Security

Vox

Explainers • Crossword • Video • Podcasts • Politics • Policy • Culture • Science • More • Give • Search

We have a request
We need 292 more contributions today to hit our goal of adding 2,500 reader gifts in September. Contributions help keep Vox's policy coverage and beyond free. Support quality journalism that's not behind a paywall. Make a contribution today.

Yes, I'll Give

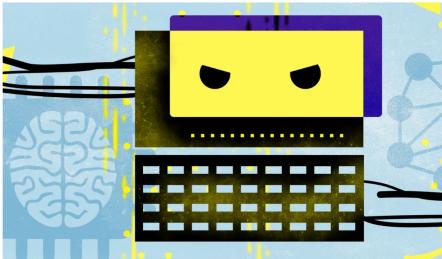
FUTURE PERFECT TECHNOLOGY

It's disturbingly easy to trick AI into doing something deadly

How "adversarial attacks" can mess with self-driving cars, medicine, and the military.

By Sigal Samuel | Apr 8, 2019, 9:10am EDT

SHARE



Javier Zarracina/Vox

Most Read

- What Dianne Feinstein's death means for California's Senate elections
- Scientists will unleash an army of crabs to help save Florida's dying reef
- 1 winner and 3 losers from Fox's dud of a second GOP debate
- More young, healthy people should be getting Paxlovid when they get Covid
- Twitter's CEO had a wild, combative appearance at the Code conference

Future Perfect
Each week, we explore unique ways to think about the world.
Sign up for our newsletter.

The illustration depicts a stylized brain on the left, connected by lines to a central computer circuit board. A large, cartoonish yellow face with black eyes is superimposed over the circuit board, symbolizing how AI can be manipulated or "tricked".

DARPA DEFENSE ADVANCED RESEARCH PROJECTS AGENCY ABOUT US / OUR RESEARCH / NEWS / EVENTS / WORK WITH US / EXPLORE BY TAG

> Defense Advanced Research Projects Agency > Our Research > Guaranteeing AI Robustness Against Deception

Guaranteeing AI Robustness Against Deception (GARD)

Dr. Alvaro Velasquez

The growing sophistication and ubiquity of machine learning (ML) components in advanced systems dramatically expands capabilities, but also increases the potential for new vulnerabilities. Current research on adversarial AI focuses on approaches where imperceptible perturbations to ML inputs could deceive an ML classifier, altering its response. Such results have initiated a rapidly proliferating field of research characterized by ever more complex attacks that require progressively less knowledge about the ML system being attacked, while proving increasingly strong against defensive countermeasures. Although the field of adversarial AI is relatively young, dozens of attacks and defenses have already been proposed, and at present a comprehensive theoretical understanding of ML vulnerabilities is lacking.

GARD seeks to establish theoretical ML system foundations to identify system vulnerabilities, characterize properties that will enhance system robustness, and encourage the creation of effective defenses. Currently, ML defenses tend to be highly specific and are effective only against particular attacks. GARD seeks to develop defenses capable of defending against broad categories of attacks. Furthermore, current evaluation paradigms of AI robustness often focus on simplistic measures that may not be relevant to security. To verify relevance to security and wide applicability, defenses generated under GARD will be measured in a novel testbed employing scenario-based evaluations.

As part of the program, GARD researchers from Two Six Technologies, IBM, MITRE, University of Chicago, and Google Research generated the following virtual testbed, toolbox, benchmarking dataset, and training materials that are now available to broader research community:

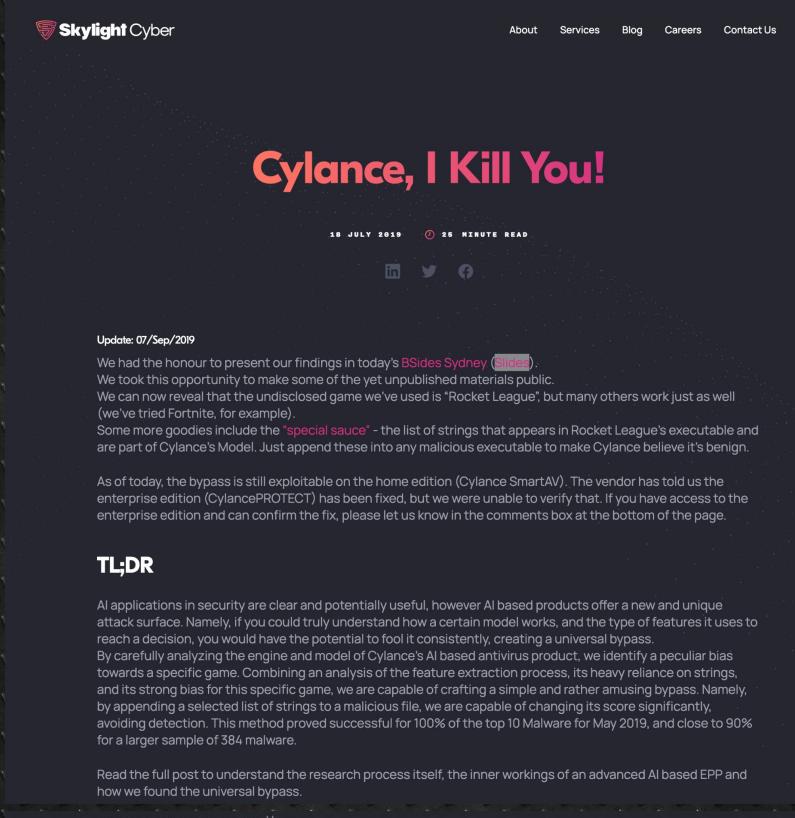
- The Armory virtual platform, available on GitHub, serves as a "testbed" for researchers in need of repeatable, scalable, and robust evaluations of adversarial defenses.
- Adversarial Robustness Toolbox (ART) provides tools for developers and researchers to defend and evaluate their ML models and applications against a number of adversarial threats.
- The Adversarial Patches Rearranged In COnText (APRICOT) dataset enables reproducible research on the real-world effectiveness of physical adversarial patch attacks on object detection systems.
- The Google Research Self-Study repository contains "test dummies" that represent a common idea or approach to build defenses.

The GARD program's Holistic Evaluation of Adversarial Defenses repository is available at <https://www.gardproject.org/>. Interested researchers are encouraged to take advantage of these resources and check back often for updates.

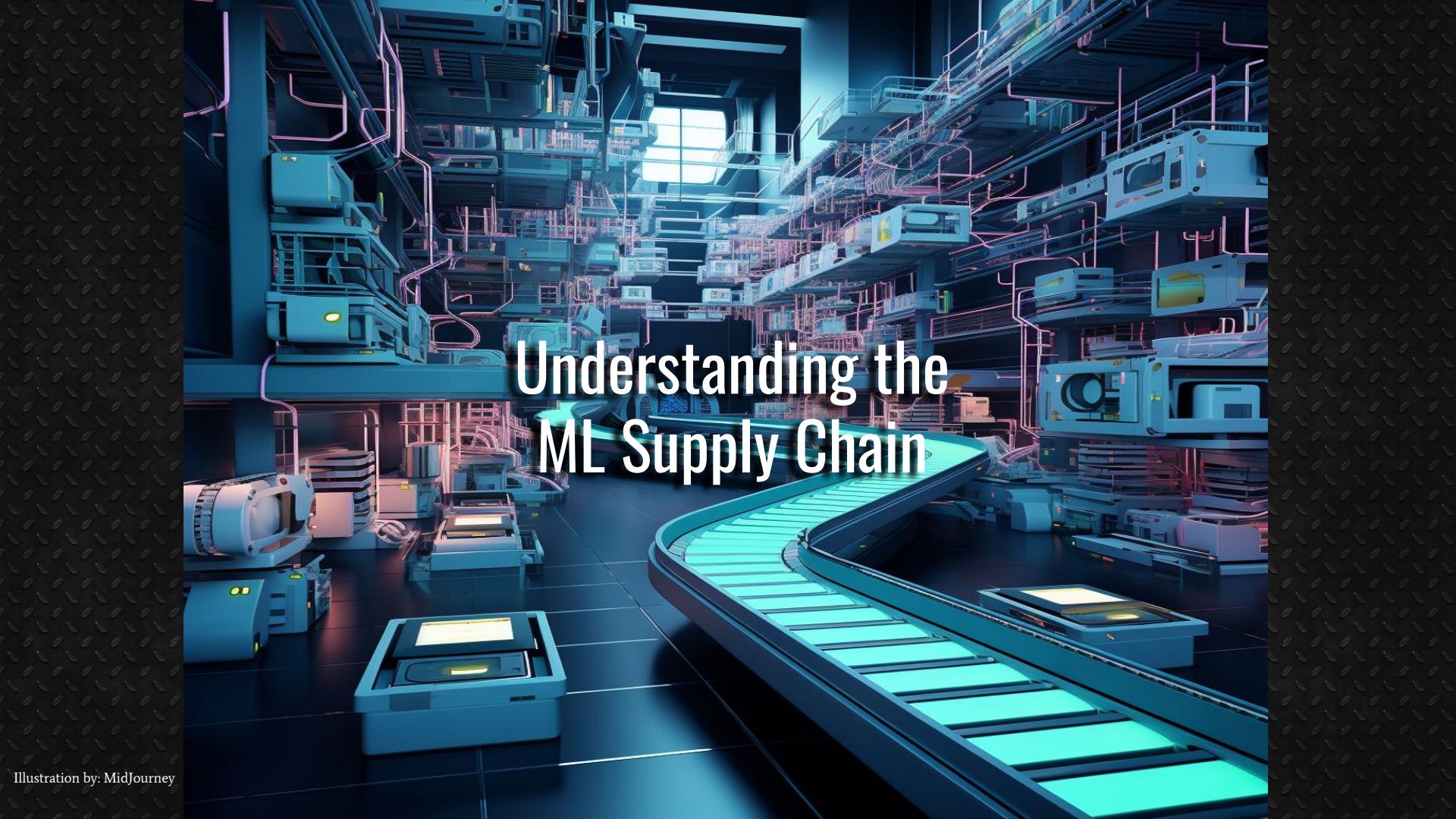
The text discusses the GARD project, which aims to establish theoretical foundations for ML system robustness and develop defenses against adversarial attacks. It highlights the creation of a virtual testbed, tools like ART and APRICOT, and datasets for research.

Resource: [Robust Physical-World Attacks on Deep Learning Visual Classification](#)

Impairing Malware Detection

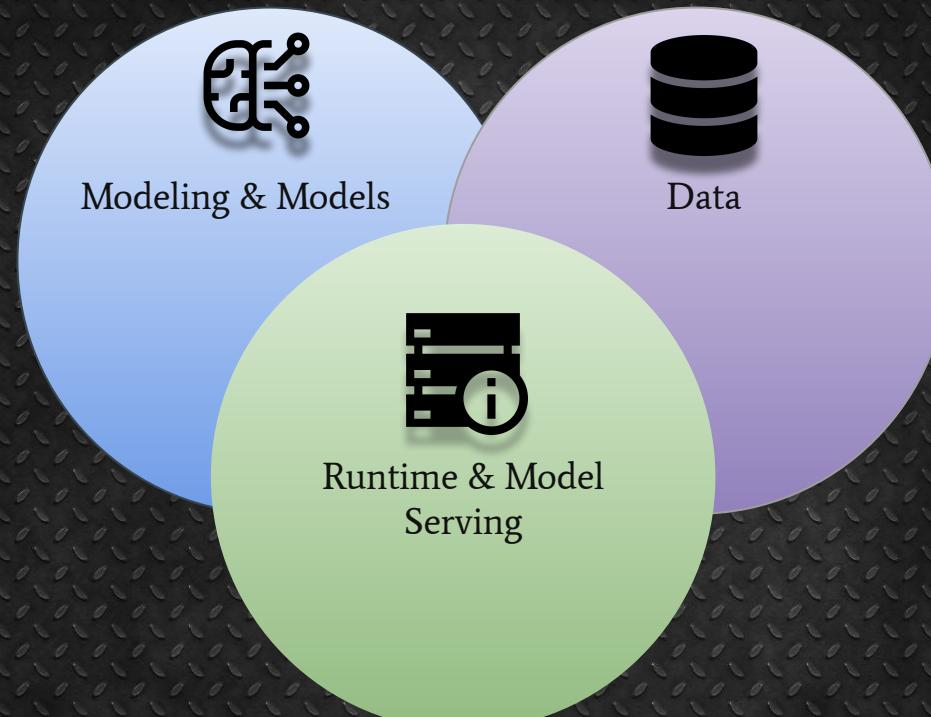


 Resource: [Cylance Machine Learning Bypass Writeup](#)

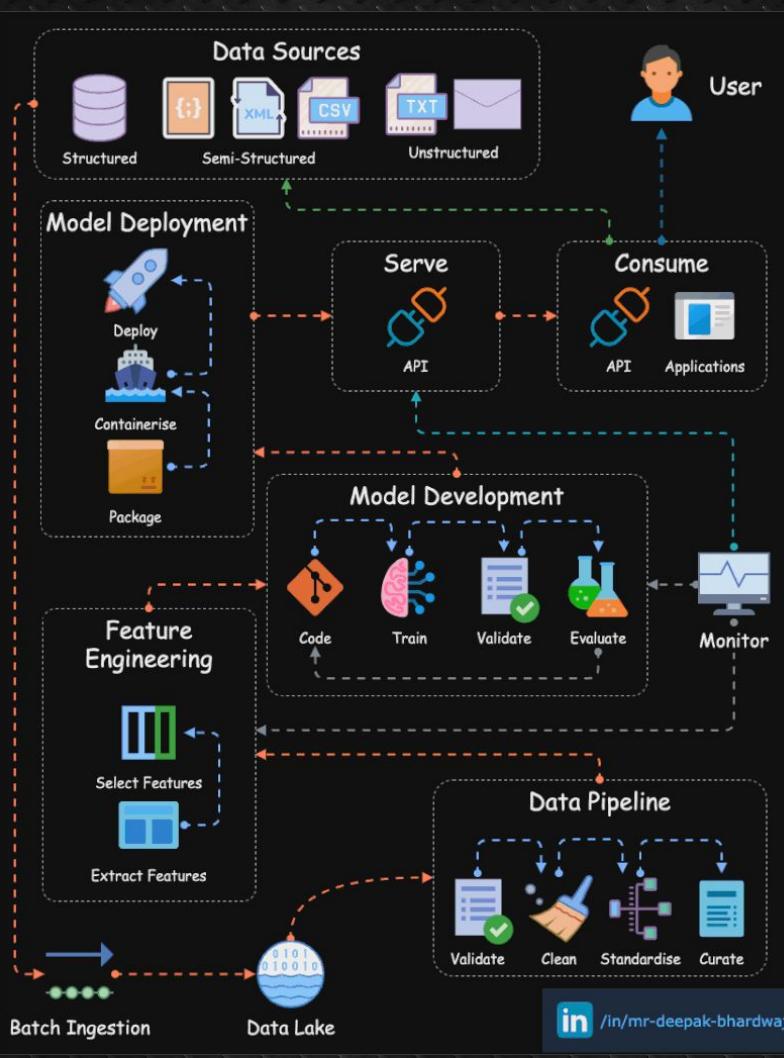


Understanding the ML Supply Chain

Machine Learning Supply Chain



The ML Stack



Data

- Data Fuels Machine Learning
- Better Models Require Increasingly More + Better Data
- Open Source Data is great but comes with a risk

The screenshot shows a white web page against a dark background. At the top, the Deloitte Insights logo is visible along with navigation links for SPOTLIGHT, TOPICS, INDUSTRIES, and MORE FROM DELOITTE INSIGHTS. A search icon is on the far right. Below the header, the word "Article" is in a small black box. To the right, the date "12 minute read · 10 December 2021" is shown. The main title "Trustworthy open data for trustworthy AI" is in large bold letters. Below it, the subtitle "Opportunities and risks of using open data for AI" is in a smaller italicized font. Four author profiles are listed with their names and countries: Tasha Austin (United States), Kara Busath (United States), Allie Diehl (United States), and Pankaj Kamleshkumar Krishnani (India). To the right of the authors are three small circular icons with symbols: a person, a download arrow, and a link. At the bottom of the article section, there is a paragraph about Fei-Fei Li's work on ImageNet.

Early in her career, Fei-Fei Li, now professor of computer science at Stanford University, recognized that an algorithm would not be able to make better decisions unless the underlying data reflects real-world data. Her solution was to map the entire image library of the world. The result of the 2.5 years of effort was ImageNet, a collection of 14 million images.¹



Resource: [The Risks of Open Data - Deloitte](#)



Resource: [Data Poisoning Explained](#)



Data



 Hugging Face

[Models](#) [Datasets](#) [Spaces](#) [Docs](#) [Solutions](#) [Pricing](#) [Log In](#) [Sign Up](#)

The AI community building the future.

The platform where the machine learning community collaborates on models, datasets, and applications.

Tasks: Libraries: Datasets: Languages: Licenses: Other: Models: 409,541 Filter by name: Meta-llama/Llama-2-7B (Test Generation + Updated 6 days ago + 25.2M + 64 stabilityai/stable-diffusion-xl-base-0.9 (Updated 6 days ago + 2.03G + 39 openchat/openchat (Test Generation + Updated 2 days ago + 1.3M + 136 llyasviel/ControlNet-v1.1 (Updated Apr 26 + 1.2M cexpense/zeroscope_v2_XL (Updated 2 days ago + 2.4M + 234 meta-llama/Llama-2-13b (Test Generation + Updated 4 days ago + 3.52B + 64 tiiuae/falcon-40b-instruct (Test Generation + Updated 27 days ago + 3.38B + 699 WizardLM/WizardCoderz-15B-V1.0 (Test Generation + Updated 3 days ago + 1.22B + 332 CompVis/stable-diffusion-v1-4 (Test-to-Image + Updated about 17 hours ago + 4.48B + 5.72M stabilityai/stable-diffusion-2-1 (Test-to-Image + Updated about 27 hours ago + 7.82B + 2.81K Salesforce/qwen-7B-8k-inst (Test Generation + Updated 4 days ago + 6.33B + 0.57 THUDU/THUDU (Test Generation + Updated 4 days ago + 1.22B + 1.17 monstera/monstera (Test Generation + Updated 4 days ago + 1.22B + 1.17)

Image Source: [HuggingFace](#)

Modeling

- Programming Languages & Libraries



keras



TensorFlow



pandas



scikit
learn



PyTorch



matplotlib



ANACONDA



NumPy



SciPy

Modeling



- Research & Prototyping Environments

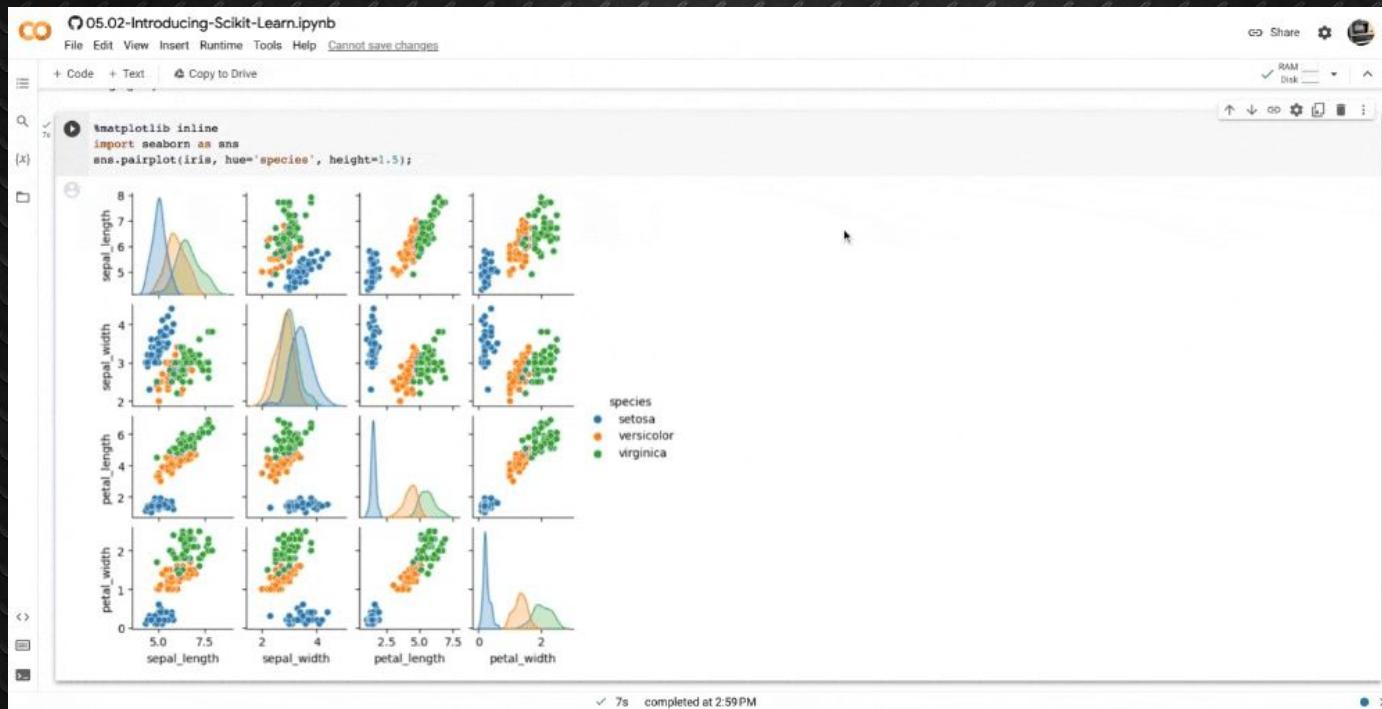


Image Source: [jakevdp - Introducing Scikit Learn Notebook](#)

OOB Models



G google/bert_uncased_L-2_H-128_A-2 □ ❤ like 22

Transformers PyTorch JAX Safetensors bert Inference Endpoints arxiv:1908.08962 License: apache-2.0

gpt2 □ ❤ like 1.4k

Text Generation Transformers PyTorch TensorFlow JAX TF Lite Rust ONNX Safetensors English doi:10.57967/hf/0039 gpt2 exbert Inference Endpoints

text-generation-inference License: mit

meta-llama/Llama-2-7b □ ❤ like 2.63k

Text Generation PyTorch English facebook meta llama llama-2 arxiv:2307.09288

mistralai/Mistral-7B-Instruct-v0.1 □ ❤ like 327

Text Generation Transformers PyTorch mistral finetuned Inference Endpoints text-generation-inference License: apache-2.0

jonatasgrosman/wav2vec2-large-xlsr-53-english □ ❤ like 281

Automatic Speech Recognition Transformers PyTorch JAX Safetensors common_voice mozilla-foundation/common_voice_6_0 English wav2vec2 audio hf-asr-leaderboard

mozilla-foundation/common_voice_6_0 robust-speech-event speech xlsr-fine-tuning-week Eval Results Inference Endpoints License: apache-2.0

OOB Models



[kaggle](#) Competitions Datasets Models Code Discussions Courses ...

Search Sign In Register

Level up with the largest AI & ML community

Join over 15M+ machine learners to share, stress test, and stay up-to-date on all the latest ML techniques and technologies. Discover a huge repository of community-published models, data & code for your next project.

[Register with Google](#) [Register with Email](#)



Who's on Kaggle?

Learners
Dive into Kaggle courses, competitions & forums.



Developers
Leverage Kaggle's models, notebooks & datasets.

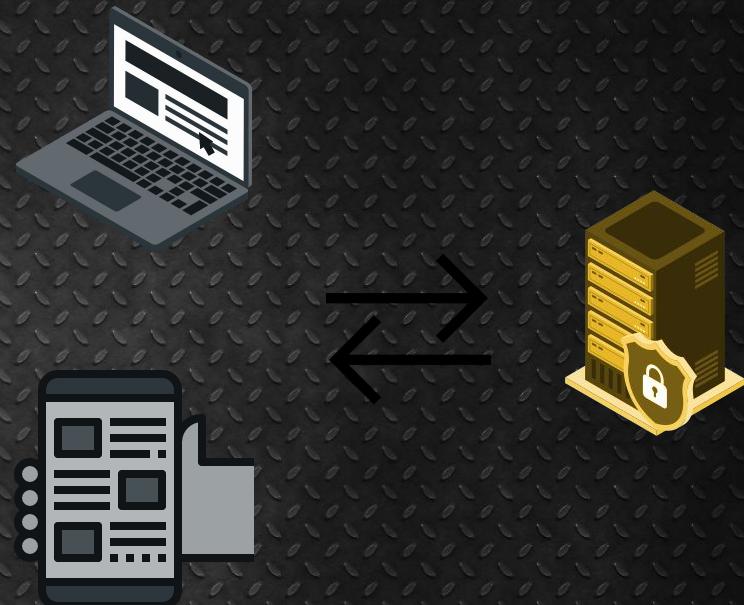


Researchers
Advance ML with our pre-trained model hub & competitions.



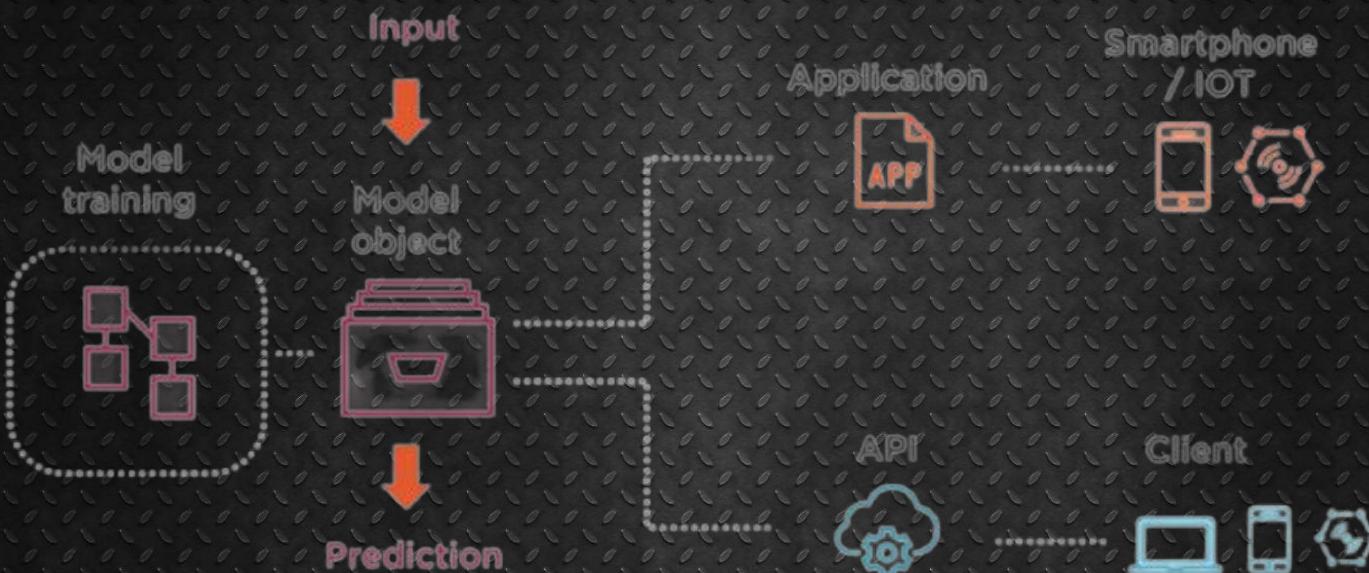
Runtime

- At runtime, the previously developed model may be converted to another programming language
 - Python
 - Java / Scala
 - Go
- The Environment itself can take many forms
 - On-Prem
 - Cloud
 - On- Device

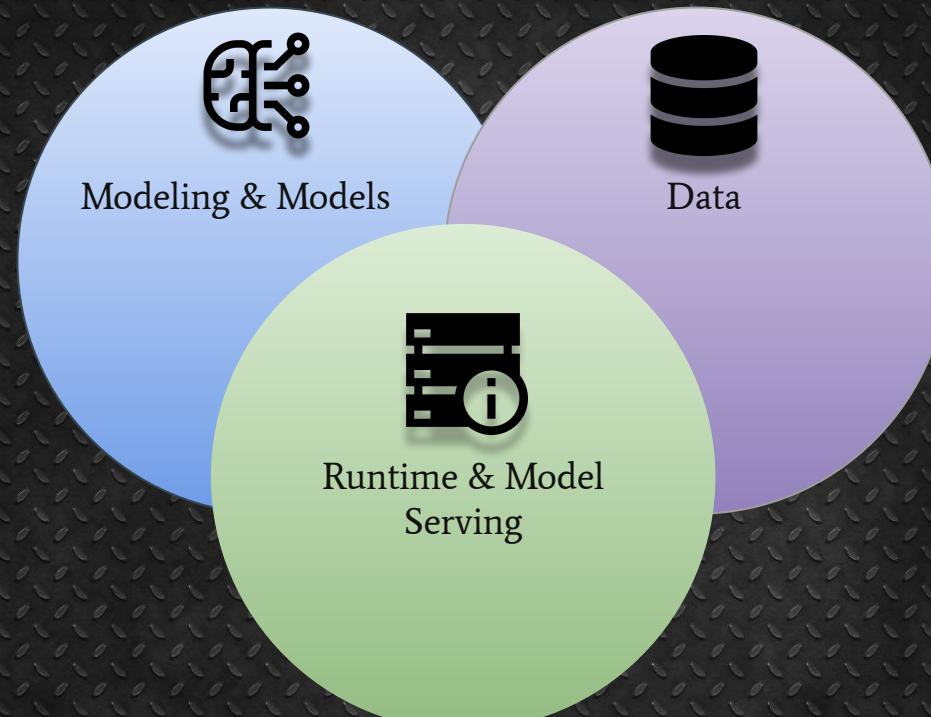




Model Serving



Machine Learning Supply Chain





Attack Vectors

MITRE Atlas Overview

Reconnaissance &	Resource Development &	Initial Access &	ML Model Access	Execution &	Persistence &	Defense Evasion &	Discovery &	Collection &	ML Attack Staging	Exfiltration &	Impact &
5 techniques	7 techniques	4 techniques	4 techniques	2 techniques	2 techniques	1 technique	3 techniques	3 techniques	4 techniques	2 techniques	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

Data: VirusTotal Poisoning



Σ 731adcf2d7fb61a8335e23dbe2436249e5d5753977ec465754c6b699e9bf161

60 / 72 Community Score

60 security vendors and 3 sandboxes flagged this file as malicious

Reanalyze Similar More

731adcf2d7fb61a8335e23dbe2436249e5d5753977ec465754c6b699e9bf161
731adcf2d7fb61a8335e23dbe2436249e5d5753977ec465754c6b699e9bf161.exe

Size 2.93 MB Last Analysis Date 3 days ago EXE

peeks idle checks-user-input direct-cpu-clock-access detect-debug-environment

DETECTION DETAILS RELATIONS BEHAVIOR COMMUNITY 25 +

Join the VT Community and enjoy additional community insights and crowdsourced detections, plus an API key to automate checks.

Popular threat label ransomware.blackcat/ispzk Threat categories ransomware trojan Family labels blackcat ispzk minerva

Security vendors' analysis Do you want to automate checks?

AhnLab-V3	Trojan/Win.Generic.R499440	Alibaba	Ransom:Win32/BlackCat.108c7b7e
ALYac	Trojan.Ransom.BlackCat	Antiy-AVL	Trojan/Win32.Filecoder
Arcabit	Trojan.Ransom.BlackCat.C	Avast	Win32.RansomX-gen [Ransom]
AVG	Win32.RansomX-gen [Ransom]	Avira (no cloud)	TR/YAV.Minerva.ispzk
BitDefender	Trojan.Ransom.BlackCat.C	BitDefenderTheta	Gen:NN.ZexaCO.36738.7IW@armxpIH
Bkav Pro	W32.AIDetectMalware	ClamAV	Win.Ransomware.BlackCat-9934796-0



Resource: [VirusTotal Data Poisoning Case Study](#)

Modeling: Code & Library Dependencies



SC MEDIA
A Cybersecurity Alliance Resource

TOPICS EVENTS PODCASTS RESEARCH RECOGNITION LEADERSHIP

Supply chain, DevSecOps

f t x in

Millions of GitHub repositories potentially vulnerable to RepoJacking

Steve Zurier June 23, 2023



There are potentially millions of repositories on GitHub susceptible to RepoJacking attacks, according to Aqua Security research. (Adobe Stock Images)

Millions of GitHub repositories are potentially vulnerable to RepoJacking, which allows malicious actors to control an old repository if an organization changed its username on the open-source software development service.

According to [research posted June 21](#) by Aqua Security's Nautilus research group, RepoJacking is when a malicious actor registers a username and creates a repository used by an organization in the past, but has since changed its username. A developer may think the repo is safe, but in reality it's controlled by the attacker and susceptible to malware.

NEWS 4 SEP 2023

Python Package Index Targeted Again By VMConnect



Alessandro Mascellino

Freelance Journalist

Email Alessandro Follow @a_mascellino

ADVERTISEMENT

Cybersecurity experts at ReversingLabs have unveiled a concerning continuation of the infamous VMConnect campaign.

This ongoing assault, initially discovered in early August, has revealed an insidious trend of cyber-criminals infiltrating the Python Package Index (PyPI), a repository for open-source Python software.

The VMConnect campaign, which originally involved two dozen malicious Python packages, has now been expanded further. In this latest wave of attacks, the perpetrators have demonstrated remarkable persistence and adaptability, raising significant concerns for the cybersecurity community.

The initial VMConnect campaign made headlines for its ability to mimic widely used Python tools, such as vConnector, eth-tester and databases, effectively concealing their malicious intent within legitimate-looking software packages.

[Read more about the campaign: VMConnect: Python PyPI Threat Imitates Popular Modules](#)

Now, ReversingLabs has once again sounded the alarm, uncovering three additional malevolent Python packages that are believed to be part of this extended campaign: tablideder, request-plus and requestspro.

One of the standout characteristics of this ongoing VMConnect campaign is the cyber-criminals' ingenuity in evading detection. Unlike traditional malware, which often activates upon installation, these malicious Python packages remain dormant until they are imported and called upon by legitimate applications.

This stealthy approach serves as a clever defense mechanism against conventional

You may also like

NEWS 4 AUG 2023

[VMConnect: Python PyPI Threat Imitates Popular Modules](#)

NEWS 9 NOV 2022

[Malicious Package on PyPI Hides Behind Image Files, Spreads Via GitHub](#)

NEWS 23 FEB 2023

[Dozens of Malicious 'HTTP' Libraries Found on PyPI](#)

NEWS 2 JUN 2023

[Malicious PyPI Packages Use Compiled Python Code to Bypass Detection](#)

NEWS 5 MAY 2023

["Kekw" Malware in Python Packages Could Steal Data and Hijack Crypto](#)



Resource: Whitepaper - [Towards Measuring Vulnerabilities and Exposures in Open-Source Packages](#)

OOB Models: Open Source Model Vulnerabilities



AI Risk Database Search by model name or URL, file hash, file artifact, or risk report... Report Vulnerability Sign In

microsoft/resnet-50

image-classification, huggingface, pytorch, jax, transformers

[Model Overview](#) Related Models

Model

Repository URL	https://huggingface.co/microsoft/resnet-50
Repository Type	huggingface
Commit Date	Jul 01, 2022
Commit Hash	f5104f67a0a8892c17fa776add3e55999dc67893
Author	microsoft
Reputation	Downloads: 2752575 Likes 124
Vulnerabilities	None

Risk Overview ⓘ

Overall Risk Score	30th percentile
Operational Risk Score	41th percentile
Security Risk Score	21th percentile
Fairness Risk Score	n/a

Top Vulnerability Reports

⚠️ Severe sensitivity to Square Attack Affects 1 model	reported by robustintelligence Jan 07, 23
⚠️ Severe sensitivity to Gaussian Blur Affects 1 model	reported by robustintelligence Jan 07, 23
⚠️ Severe sensitivity to Contrast Increase Affects 1 model	reported by robustintelligence Jan 07, 23
⚠️ Severe sensitivity to Randomize Pixels With Mask Affects 1 model	reported by robustintelligence Jan 07, 23
⚠️ Severe sensitivity to Gaussian Noise Affects 1 model	reported by robustintelligence Jan 07, 23

[View All](#)

Image Source: airisk.io



Resource: Talk AI Village DEFCON 31 - [Assessing the Vulnerabilities of the Open-Source Artificial Intelligence \(AI\) Landscape](#)

Model Serving: Hijacking Facial Recognition System



South China Morning Post | China | Economy | HK | Asia | Business | Tech | Lifestyle | People & Culture | World | Comment | Video | Sport | Post Mag | Style | All | SUBSCRIBE | myNEWS |

Facial recognition FOLLOW

Get more with **myNEWS**
A personalised news feed
of stories that matter to
you [Learn more →](#)

Tech / Tech Trends

Chinese government-run facial recognition system hacked by tax fraudsters: report

A group of tax scammers hacked a government-run identity verification system to fake tax invoices

The fake tax invoices from the criminal group were valued at US\$76.2 million

Masha Boruk FOLLOW Published: 7:00am, 31 Mar, 2021 • [Why you can trust SCMP](#)

TOP PICKS

Tech Tech war takes a new turn as Huawei launches 5G smartphones with mystery chip 1 Feb 2023

Tech Huawei's Meng ends controversial chairmanship with 5G smartphone success 2 Feb 2023

Tech Chinese AI champion SenseTime hires ex-employee under police probe 2 Feb 2023

Tech Crypto trading in mainland China and Hong Kong drops along with East Asia 2 Oct 2023


A criminal group duped a government-run platform's facial recognition system by using manipulated personal information and high definition photographs, which were bought from an online black market. Photo: Shutterstock

Identity verification using facial recognition is widely adopted in China, as the technology has become an integral part of apps from mobile payments and travel to retail, as well as surveillance systems and online platforms for government services.

💡 Resource: MITRE Case Study - [Facial Recognition Hijack](#)

💡 Resource: Whitepaper - [Attacks for Extraction of Embedded Neural Network Models](#)



Safeguarding the AI Realm

What is available for Defenders?



Mitigations

Mitigations represent security concepts and classes of technologies that can be used to prevent a technique or sub-technique from being successfully executed.



This draft of mitigations are defined based on current ATLAS case studies.
Feedback and contributions are welcome - please join the [#mitigations channel on Slack](#).

The table below lists mitigations from MITRE ATLAS™. Scroll through the table or use the filter to narrow down the information.

Search for keywords

Adversarial



ID	Name ⓘ	Description
AML.M0003	Model Hardening	Use techniques to make machine learning models robust to adversarial inputs such as adversarial training or network distillation.
AML.M0006	Use Ensemble Methods	Use an ensemble of models for inference to increase robustness to adversarial inputs. Some attacks may effectively evade one model or model family but be ineffective against others.
AML.M0008	Validate ML Model	Validate that machine learning models perform as intended by testing for backdoor triggers or adversarial bias.
AML.M0010	Input Restoration	Preprocess all inference data to nullify or reverse potential adversarial perturbations.
AML.M0015	Adversarial Input Detection	Detect and block adversarial inputs or atypical queries that deviate from known benign behavior, exhibit behavior patterns observed in previous attacks or that come from potentially malicious IPs. Incorporate adversarial detection algorithms into the ML system prior to the ML model.

What is available for Defenders?



Open Source CALDERA
Plugin

- Developed for Adversary Emulation AI
- TTPs defined by MITRE Atlas



Resource: [MITRE Arsenal Repository](#)

What is available for Defenders?



OWASP

PROJECTS CHAPTERS EVENTS ABOUT Q Member Login

OWASP Top 10 for Large Language Model Applications

Main Example

OWASP Top 10 for Large Language Model Applications version 1.1

LLM01: Prompt Injection

Manipulating LLMs via crafted inputs can lead to unauthorized access, data breaches, and compromised decision-making.

LLM02: Insecure Output Handling

Neglecting to validate LLM outputs may lead to downstream security exploits, including code execution that compromises systems and exposes data.

LLM03: Training Data Poisoning

Tampered training data can impair LLM models leading to responses that may compromise security, accuracy, or ethical behavior.

LLM04: Model Denial of Service

Overloading LLMs with resource-heavy operations can cause service disruptions and increased costs.

LLM05: Supply Chain Vulnerabilities

Depending upon compromised components, services or datasets undermine system integrity, causing data breaches and system failures.

LLM06: Sensitive Information Disclosure

Failure to protect against disclosure of sensitive information in LLM outputs can result in legal consequences or a loss of competitive advantage.

LLM07: Insecure Plugin Design

LLM plugins processing untrusted inputs and having insufficient access control risk severe exploits like remote code execution.

LLM08: Excessive Agency

Granting LLMs unchecked autonomy to take action can lead to unintended consequences, jeopardizing reliability, privacy, and trust.

LLM09: Overreliance

Failing to critically assess LLM outputs can lead to compromised decision making, security vulnerabilities, and legal liabilities.

LLM10: Model Theft

Unauthorized access to proprietary large language models risks theft, competitive advantage, and dissemination of sensitive information.

Disclaimer: This document is a work in progress and is subject to change. It is intended to provide general guidance and is not a substitute for professional advice.

Author: OWASP LLM Task Force

What is available for Defenders?



[ISO](#) Standards About us News Taking part [Store](#)

Search

← ICS ← 35 ← 35.020

ISO/IEC FDIS 5338

Information technology — Artificial intelligence — AI system life cycle processes

General information

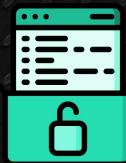
Status : Under development

Edition : 1 Number of pages : 39

Technical Committee : ISO/IEC JTC 1/SC 42 Artificial intelligence

ICS : 35.020 Information technology (IT) in general

What is available for Defenders?



Protect AI
@ProtectAICorp

Protect AI announces three open source tools that detect vulnerabilities in ML systems, and are freely available to organizations.

Learn more about our commitment to making AI security accessible to all - bnews.pr/3rHP5C2

#protectai #mlsecops #aisecurity #cybersecurity

businesswire.com

1:29 PM · Oct 6, 2023 · 865 Views

- **Rebuff:** Prompt Injection defense framework.
- **ModelScan:** Identify unsafe code in models.
- **NB Defense:** Jupyter Notebook security.



Resource: [Protect AI Open Sources Three Tools to Help Organizations Secure AI/ML Environments from Threats](#)

Key Takeaways

- ML Engineering = Software Engineering.
- ML Security > LLM Prompt Injections.
- Data can be a blind spot, good data governance is the best mitigation.
- Explainability can be your best friend or your worst enemy.



Thank you!



danielfernandez@infosec.exchange



@danielfernandez



/in/dafnz



danielfernandez.medium.com

Slides: github.com/dnlfdz/talks



hey@dafnz.com