# CS 7775

# Seminar in Computer Security: Machine Learning Security and Privacy
# Fall 2023

Alina Oprea
Associate Professor
Khoury College of Computer Science

September 7 2023

# Introduction

- **Ph.D. at CMU, 2007**
  - Research in applied cryptography, data security, and cryptographic file systems
- **RSA Laboratories, 2007-2016**
  - Cloud and storage security, applied cryptography, game theory
  - ML/AI in security
- **Northeastern Khoury College – since Fall 2016**
  - NDS2 Lab part of the Cybersecurity and Privacy Institute
  - Adversarial machine learning: study the vulnerabilities of ML in face of attacks and design defenses (poisoning, evasion)
  - Privacy in machine learning: auditing, memorization, membership inference
  - Machine learning for security: threat detection, collaborative defenses
- **Sabbatical at Google Research, 2022-2023**
  - Privacy auditing of federated models
  - Privacy of large language models (LLMs)

# Alina Oprea: Trustworthy Machine Learning

## ML Integrity
- Poisoning with different objectives (availability, backdoor, subpopulation)
- Realistic poisoning and evasion attacks for cyber applications
- Poisoning in decentralized systems (FL, P2PFL) amplified by network attacks
- Mitigations: data sanitization, ensemble-based, formal verification
- Poisoning survey 2023
- Adv ML taxonomy NIST report 2023

**Collaborators**: Cristina Nita-Rotaru, NIST
**Funding**: ARL, DoD, MIT LL, Microsoft

## ML Privacy
- Membership inference attacks on ML updates and fine-tuned models, label-only attacks
- First memorization attack on large language models (GPT-2)
- Property inference attacks with poisoning
- Introduced privacy auditing to empirically estimate privacy leakage
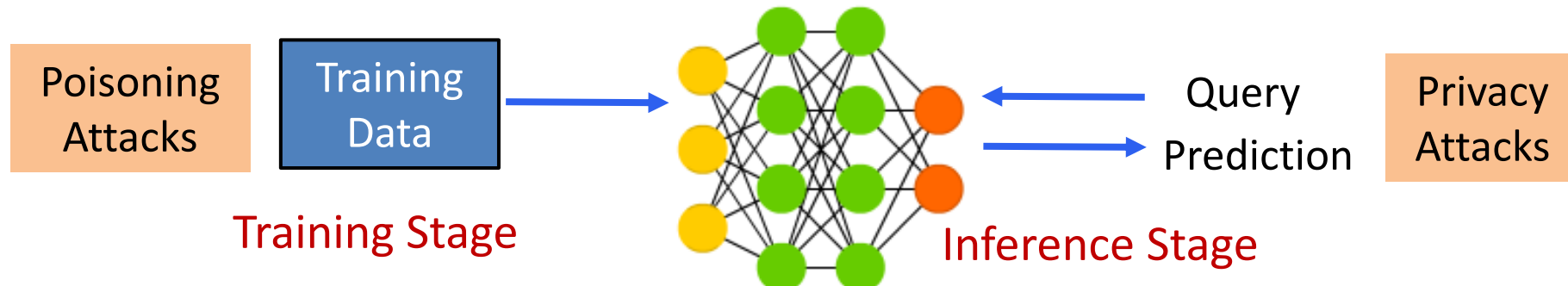
**Collaborators**: Jon Ullman, Google Research
**Funding**: NSF, Apple

## AI for Cyber Defense
- Supervised, semi-supervised and federated models trained on security logs for threat detection
- Epidemiological models for self-propagating malware
- Cyber resilience in real-world graphs
- Measurements of attacks on the web (Rapid 7 honeypot data)
- Adversarially-resilient RL framework for cyber defense

**Collaborators**: Tina Eliassi-Rad, Engin Kirda, Wil Robertson
**Funding**: DARPA, NSF, PwC, Cisco



Poisoning Attacks → Training Data → (neural network) → Query Prediction → Privacy Attacks

Training Stage     Inference Stage

# TA Introduction

- John Abascal
  - $3^{rd}$ year PhD student at Northeastern
  - Working on privacy of ML models
  - Part of the NDS2 research lab

# Class Introduction

- Research area

- What topics you are interested in trustworthy ML

- ML experience

- What do you hope to get from the class

- Something we cannot read online about you!

# CS 7775 Course objectives

- Provide in-depth coverage of adversarial attacks on ML:
  - Evasion attacks at inference time
  - Poisoning attacks at training time
  - Privacy attacks
- Learn how to classify the attacks according to the adversarial objective, knowledge, and capability (Discuss NIST report)
- Understand existing methods for training robust models and the challenges of achieving both robustness and accuracy
- Discuss security and privacy of LLMs
- Read and discuss research papers in trustworthy ML as a group
- Work on a research project individually or in a team of 2

# Course Policies

- **Website:**
  - https://www.ccs.neu.edu/home/alina/classes/Fall2023
- **Schedule**
  - Mon and Thu 11:45am – 1:25pm EST
  - Office hours:
    - Alina: Thursday 3:00 – 4:00 pm
    - John: Monday, 2:00 - 3:00 pm
- **Online resources**
  - Use Piazza for questions and discussion
  - Gradescope for paper summaries and assignments

# Class Outline

- Introduction
  - Review of machine learning and deep learning
  - Taxonomy of adversarial ML
- Evasion attacks and defenses
- Poisoning attacks
- Privacy attacks and defenses
  - Membership inference, memorization
  - Differential privacy, auditing
- LLM security and privacy
- Schedule is tentative and flexible
  - Will ask for paper/ topic suggestions (last 2 lectures are open)

# Grading

- <span style="color:red">Assignments – 15%</span>
  - 2 assignments at the beginning of class (first month)
- <span style="color:red">Paper summaries – 10%</span>
  - Read and submit paper summaries before every class
- <span style="color:red">Discussion leading and class participation – 25%</span>
  - Lead discussion in several classes and actively participate in discussion
- <span style="color:red">Final project – 50%</span>
  - Select your own project topic related to trustworthy AI (individual or teams of 2)
  - Two types of projects: research or systematization of knowledge (SoK)
  - Project proposal presented in class mid Oct
  - Milestone mid Nov
  - Presentation at end of class (early Dec) and written report due Dec 11

# Academic Integrity

- Homework / paper summaries are done individually

- Class project is done either individually or in a team of 2

- Rules
  - Can discuss with colleagues or instructors
  - Can post and answer questions on Piazza
  - Code cannot be shared with colleagues
  - If using code from the Internet, add reference
    - Use python or R packages, and can get inspiration from tutorials, but do not use directly entire projects from Internet

- https://osccr.sites.northeastern.edu/academic-integrity-policy/

# ML Resources

• Trevor Hastie, Rob Tibshirani, and Jerry Friedman, [Elements of Statistical Learning](), Second Edition, Springer, 2009.

• Christopher Bishop. [Pattern Recognition and Machine Learning](). Springer, 2006.

• A. Zhang, Z. Lipton, and A. Smola. [Dive into Deep Learning]()

• Lecture notes by Andrew Ng from Stanford

• DS 4400 lecture notes:

[http://www.ccs.neu.edu/home/alina/classes/Spring2021/](http://www.ccs.neu.edu/home/alina/classes/Spring2021/)

# AI Deployments in Practice

# Deep Learning



Neural networks return and excel at image recognition, speech recognition, …

The 2018 Turing award was given to Yoshua Bengio, Geoff Hinton, and Yann LeCun.

input layer    hidden layer 1    hidden layer 2    hidden layer 3

output layer

# Generative AI



Stable Diffusion
"Mountains with flowers in bloom"



DALL-E 2
"A painting of a fox sitting in a field at sunrise in the style of Claude Monet"



Write a poem about the sea

A boundless expanse, the sea unfolds,
Its story written in waves, untold.
A symphony of blue, captivating sight,
Where heaven meets earth, in pure delight.

Gentle whispers of the ocean breeze,
Caress the senses with tranquil ease.
Salt-kissed air upon the skin,
Invigorating, as new journeys begin.

ChatGPT

# What is Your Favorite ML / Deep Learning Application?

# Applications of ML

- Healthcare
- Vision
- NLP
- Speech recognition
- Self-driving cars
- Stock market analysis
- Recommendations
- Sentiment analysis
- Human behavior
- Quality of life

- Business
- Sports
- Bots / chatbots
- Science / engineering
- Bioinformatics
- Precision medicine
- Unsupervised learning
- Reinforcement learning

# Success stories: Speech recognition

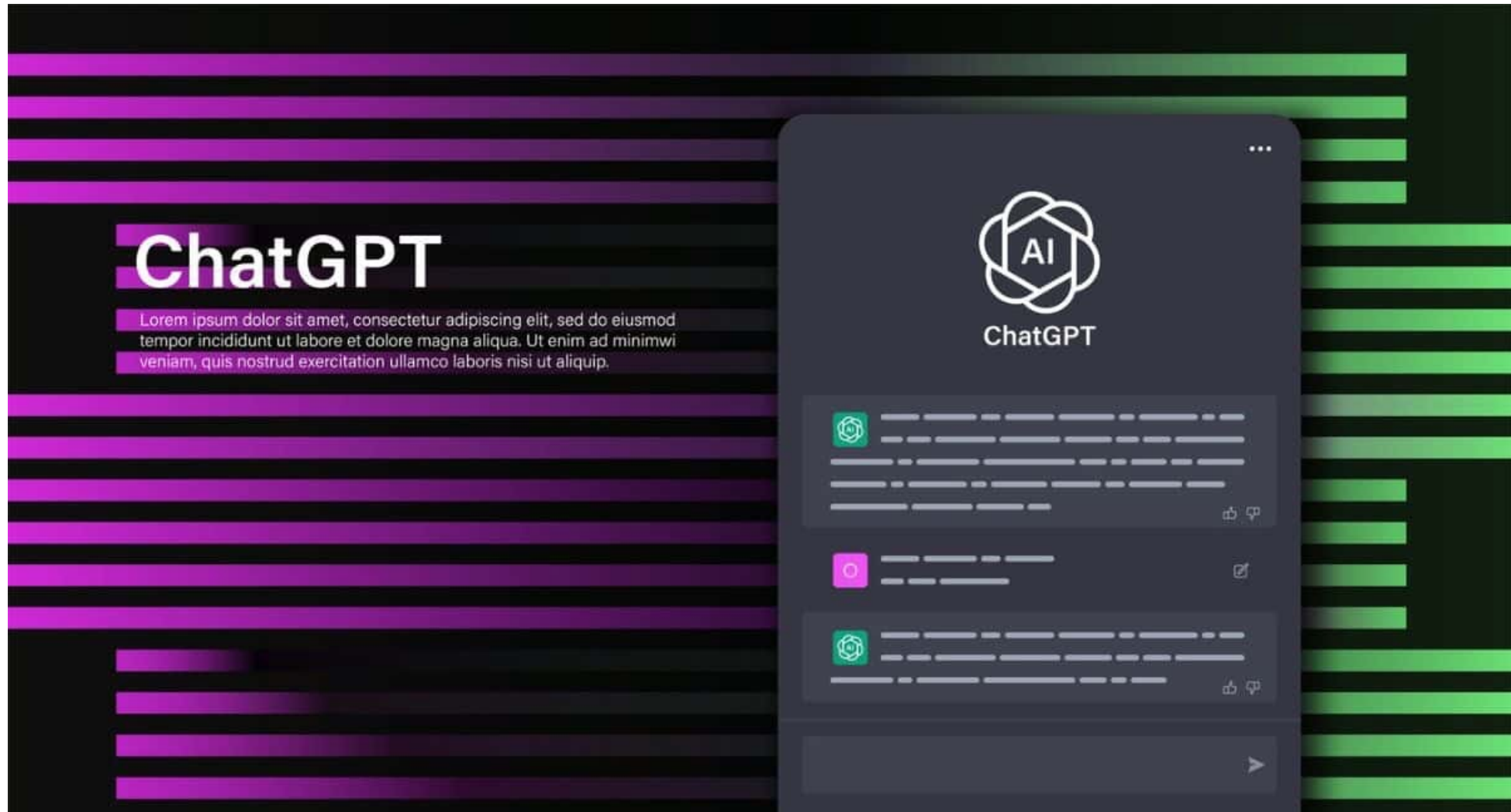# Success stories: Machine Translation
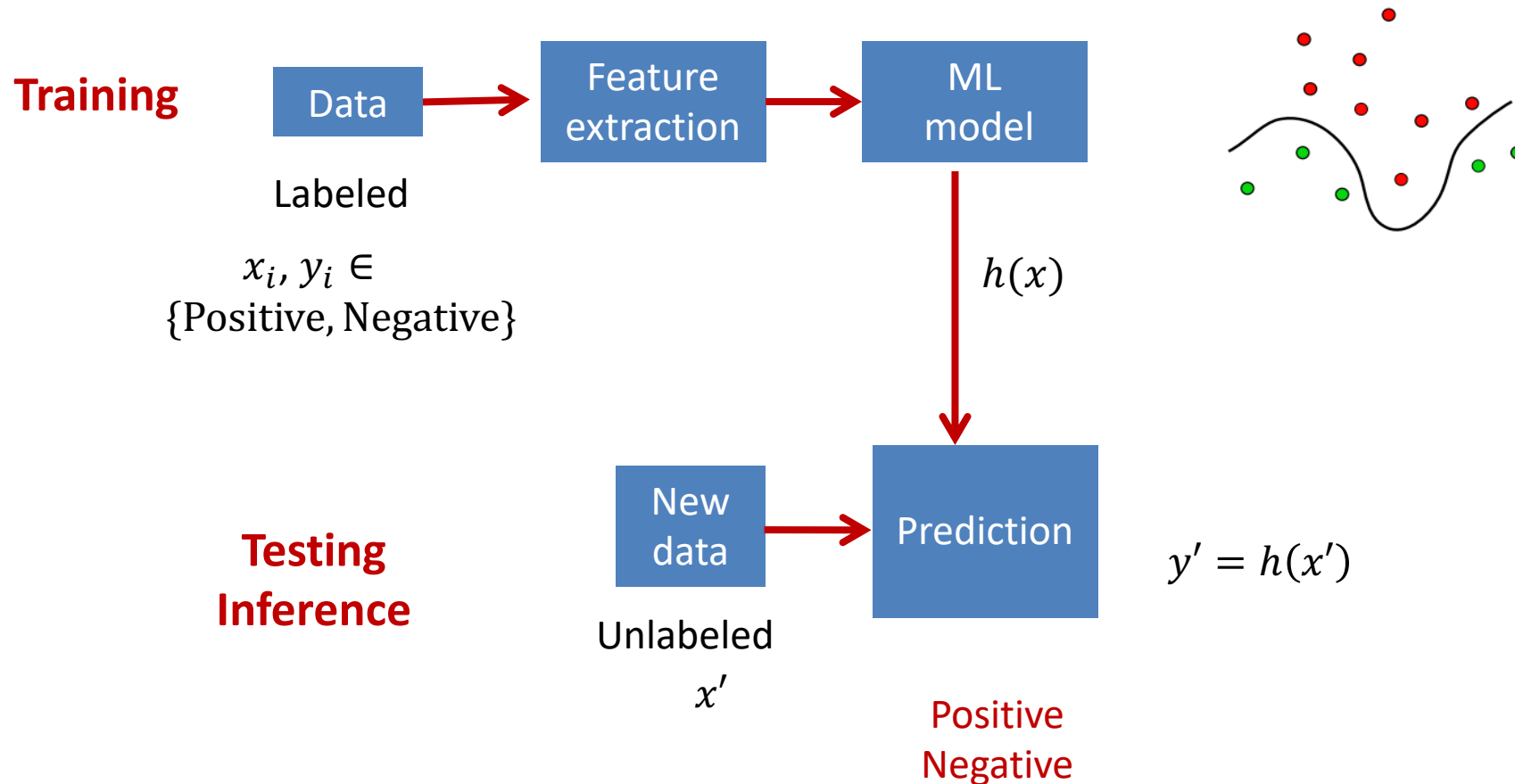
# Success stories: Image segmentation

# Success stories: Chatbots

# Short History of ML

- Legendre and Gauss – linear regression, 1805
  - Astronomy applications
- Probabilistic models
  - Bayes and Laplace - Bayes Theorem, 1812; Markov chains, 1913
- Fisher – linear discriminant analysis for classification, 1936
  - Logistic regression, 1940
- Rosenblatt - Perceptron, 1958
- Widrow and Hoff - ADALINE neural network, 1959
- "AI winter", limitations of  perceptron and linear models, 1970
- Breiman, Friedman, Olshen, Stone - decision trees (non-linear models), 1980
- Cortes and Vapnik - SVM with kernels, 1990
- Breiman: Bagging, 1994; Ho – random forest, 1995;  Freund and Shapire – AdaBoost, 1997
- Geoffrey Hinton, Deep learning, back propagation, 2006
- C. Szedegy: Adversarial manipulation of image classification, 2013
- ChatGPT release: Nov 2022

# Supervised Learning

**Training**

Data → Feature extraction → ML model

Labeled

$x_i, y_i \in$
{Positive, Negative}

$h(x)$

**Testing Inference**

New data → Prediction

Unlabeled
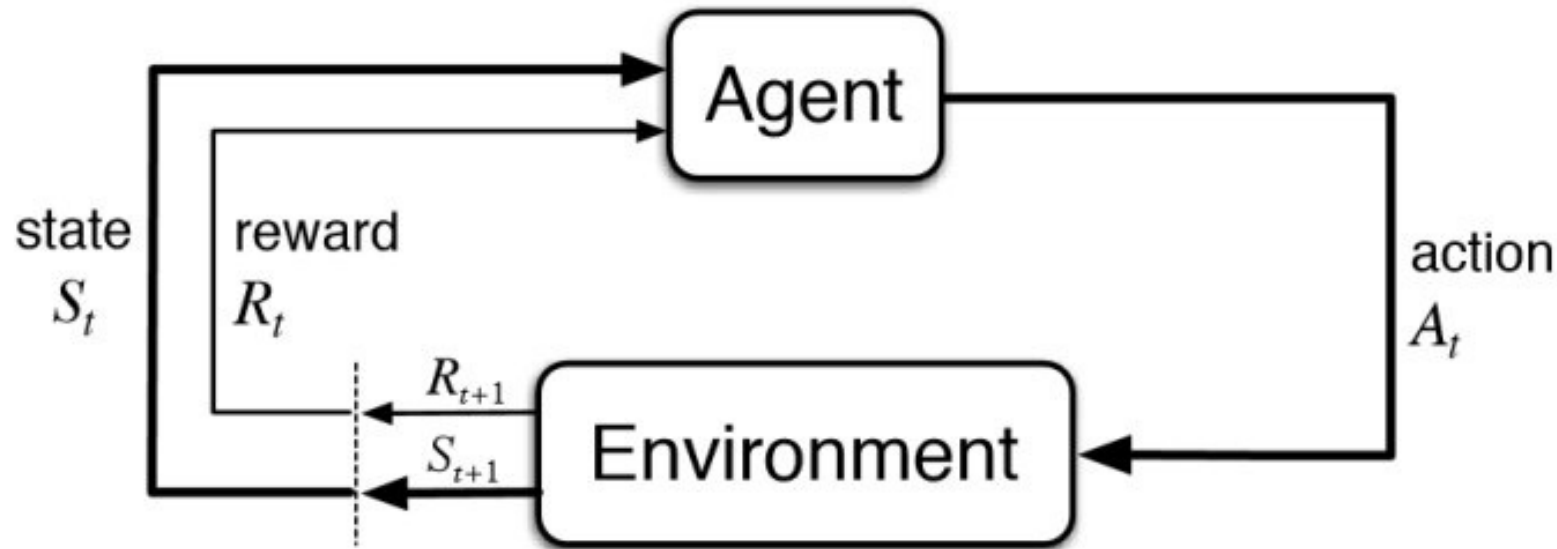$x'$

$y' = h(x')$

Positive
Negative

- Main Assumption: Distribution of training and testing data is similar
- Model can learn from training data and generalize to testing data
- Concrete metrics to measure model performance

# Unsupervised Learning

- Input: unlabeled data
- Clustering
  - Group similar data points into clusters
  - Examples: k-means, hierarchical clustering, density-based clustering
- Dimensionality reduction
  - Project the data to lower dimensional space
  - Examples: PCA (Principal Component Analysis), UMAP
- Anomaly detection
  - Learn normal patterns during training and identify anomalies at testing
  - Examples: KDE, auto encoders, Local Outlier Factor, Isolation Forest

# Reinforcement Learning



- Agents learn by interacting with an environment
- They take actions and obtain reward
- Goal: learn optimal policy to maximize reward
- Methods: Q learning, Deep Q Networks (DQN)
- Applications: Games (AlphaGo Zero), robotics
- https://deepmind.com/blog/article/alphago-zero-starting-scratch

# Federated Learning

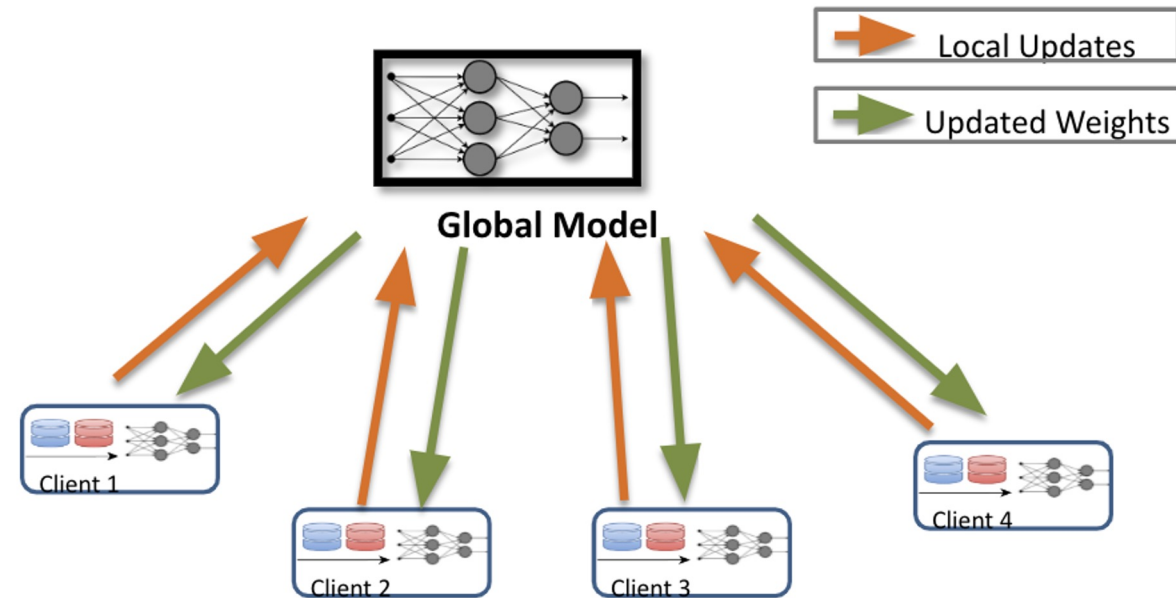Multiple clients collaboratively train machine learning models by interacting with an aggregation server
- Federated Averaging [McMahan et al. 2017]

Training is an iterative process
- Clients receives global model
- Subset of clients update the model using local data and send updates to server
- Global model is updated by aggregating client contributions

Benefits
- Training data remains on client devices
- Computational efficiency



McMahan, Brendan, et al. Communication-efficient learning of deep networks from decentralized data. *Artificial Intelligence and Statistics*. PMLR, 2017.
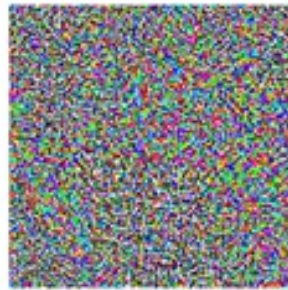
# Security and Privacy Risks of AI

- Deep Neural Networks and other classifiers are not resilient to adversarial manipulations
  - Szegedy et al. *Intriguing properties of neural networks*. 2013
  - Biggio et al. *Evasion attacks against machine learning at test time*. 2013
  - Goodfellow et al. *Explaining and Harnessing Adversarial Examples*. 2014
- Started the field of Adversarial Machine Learning

$x$
"panda"
57.7% confidence

$+.007 \times$

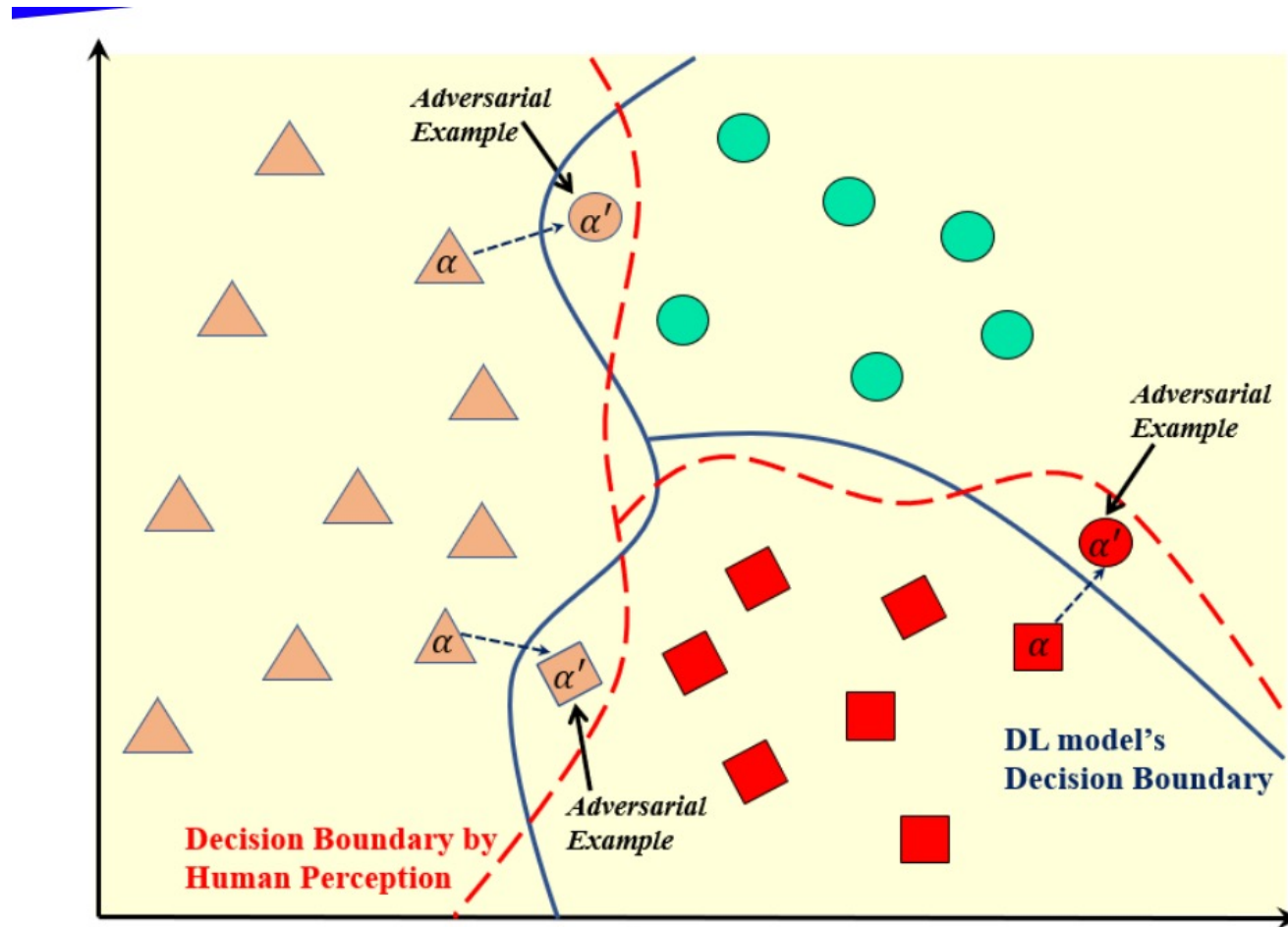$\text{sign}(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
"gibbon"
99.3 % confidence

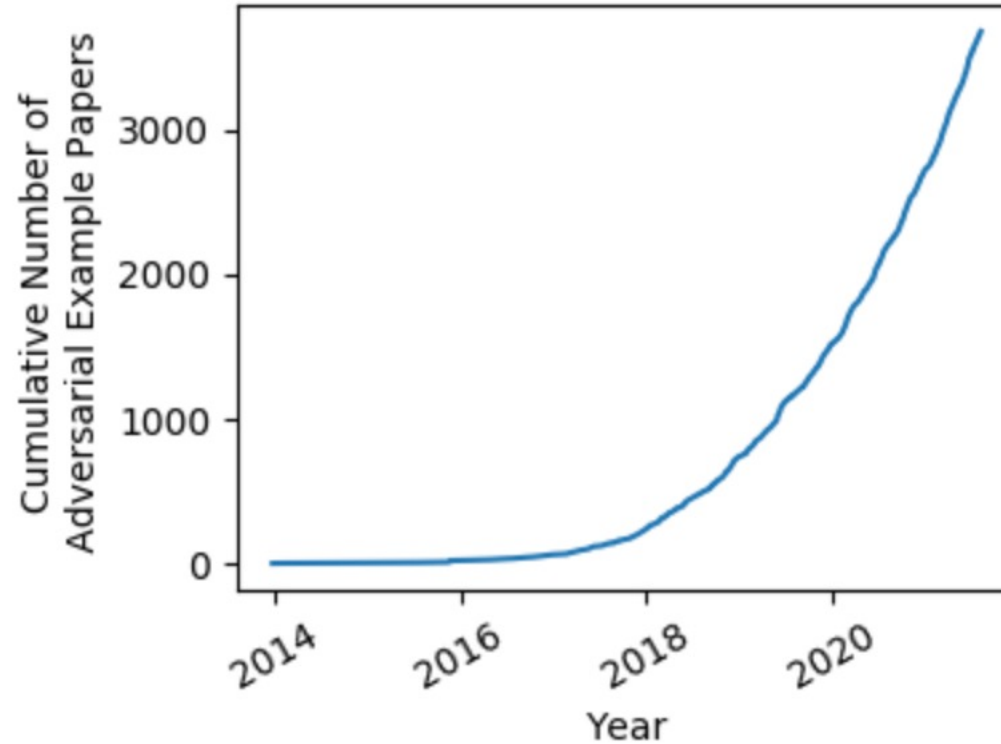Attacker changes distribution of testing data!

Adversarial example

# What are Adversarial Examples



Adversarial Robustness of Deep Learning: Theory,
Algorithms, and Applications. Tutorial at ICDM 2020
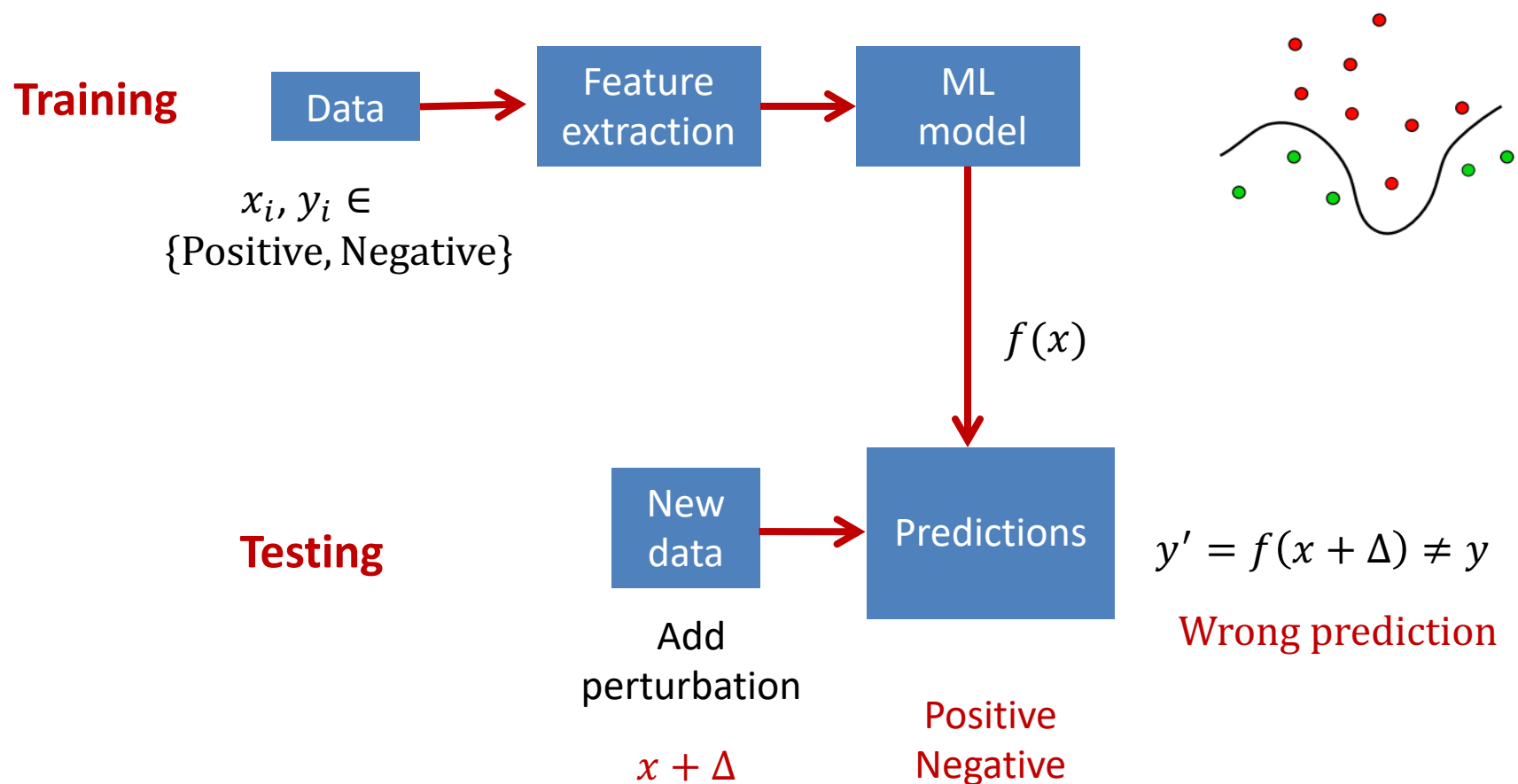
# Adversarial ML Literature



- Graph by Nicholas Carlini, Google
- Papers published in AI and security conferences
- We will only cover a small subset (~35 papers)

# Safety Concerns of AI

# Safety Concerns of AI

- ## Adversarial ML
  - ML can be manipulated
  - Small change in input results in different prediction (adversarial examples / evasion attacks)
  - Corrupted training data can modify the model (poisoning attacks)
- ## Privacy concerns
  - User data remains private when ML models are trained on it
- ## Ethics and fairness of AI
  - Predictions of ML are fair for underrepresented minorities
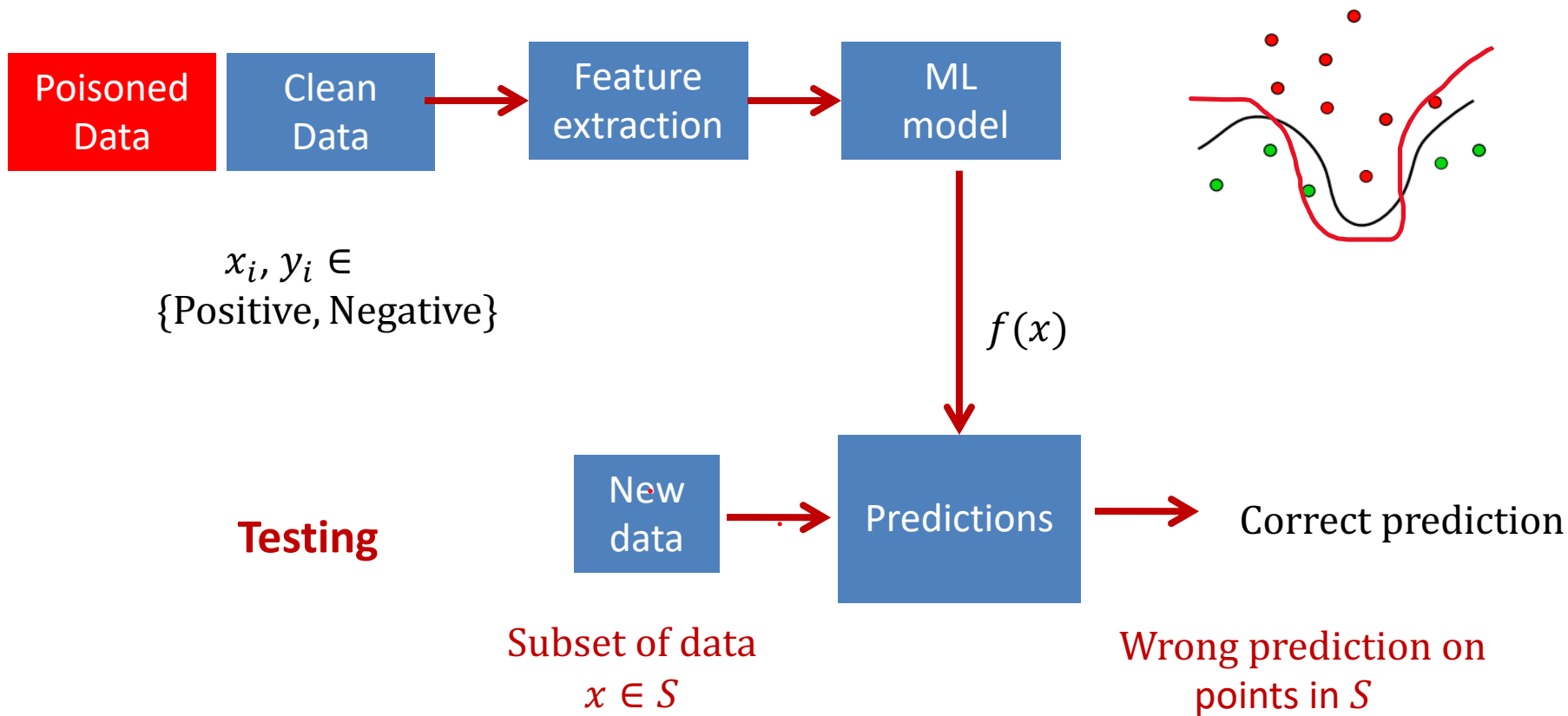  - Robots will not perform harmful actions

# Evasion Attacks



**Training**

Data → Feature extraction → ML model

$x_i, y_i \in$
{Positive, Negative}

$f(x)$

**Testing**

New data → Predictions

Add perturbation

$x + \Delta$

Positive
Negative

$y' = f(x + \Delta) \neq y$

Wrong prediction

- Modify testing point by adding small perturbation to misclassify it

# Poisoning Attacks

**Training**

| Poisoned Data | Clean Data | → | Feature extraction | → | ML model |
|---|---|---|---|---|---|

$x_i, y_i \in$
{Positive, Negative}

$f(x)$

**Testing**

| New data | → | Predictions | → | Correct prediction |
|---|---|---|---|---|

Subset of data
$x \in S$

Wrong prediction on points in $S$

- Poisoning attack inserts corrupted data at training
- Model makes incorrect predictions on subset of data at testing

# Privacy Attacks on ML

Deployed
ML Model

Data

Labels

Query     Prediction

- Reconstruction attacks: Extract sensitive attributes
  - [Dinur and Nissim 2003]
- Membership Inference: Determine if sample was in training
  - [Shokri et al. 2017], [Yeom et al. 2018], [Hayes et al. 2019], [Jayaraman et al. 2020]
- Model Extraction: Learn model architecture and parameters
  - [Tramer et al. 2016], [Jagielski et al. 2020]
- Memorization: Extract training data from queries to the model
  - [Carlini et al. 2021]

# Real-World Attacks: Road Sign Classification



Eykholt et al. *Robust Physical-World Attacks on Deep Learning Visual Classification*. In CVPR 2018

# Real-World Attacks: Face Recognition

**Adversarial Glasses**

- M. Sharif et al. (ACM CCS 2016) attacked deep neural networks for face recognition with carefully-fabricated eyeglass frames

- When worn by a 41-year-old white male (left image), the glasses mislead the deep network into believing that the face belongs to the famous actress Milla Jovovich

# Adversarial Examples in Connected Cars



Original Image
Steering angle = -4.25

Adversarial Image
Steering angle = -2.25

- Udacity challenge: Predict steering angle from camera images, 2014
- A. Chernikova, A. Oprea, C. Nita-Rotaru, and B. Kim. *Are Self-Driving Cars Secure? Evasion Attacks against Deep Neural Networks for Self-Driving Cars.* 2019

# Adversarial ML in the Real World



Slide from David Evans, UVA

# Poisoning in the Real World



Listen to this article

It took less than 24 hours for Twitter to corrupt an innocent AI chatbot. Yesterday, Microsoft unveiled Tay — a Twitter bot that the company described as an experiment in "conversational understanding." The more you chat with Tay, said Microsoft, the smarter it gets, learning to engage people through "casual and playful conversation."

# Misinformation with Generative AI



From Barrett et al. Identifying and Mitigating the Security Risks of Generative AI, 2023

# Hallucination with Generative AI



FORBES > BUSINESS

BREAKING

## Lawyer Used ChatGPT In Court—And Cited Fake Cases. A Judge Is Considering Sanctions

**Molly Bohannon** Forbes Staff

*I cover breaking news.*

Follow

Jun 8, 2023, 02:06pm EDT

# Generative AI for Cybercrime



JULY 13, 2023   |   DANIEL KELLEY   |   BEC / EMAIL PROTECTION / THREAT DISCOVERY / UNCATEGORIZED

## WormGPT – The Generative AI Tool Cybercriminals Are Using to Launch Business Email Compromise Attacks

# Spear phishing with Generative AI

# Summary

- AI has a long history
- Adversarial ML gained attention with the discovery of adversarial examples by Szedegy et al. and Biggio et al.in 2013
- Different types of adversarial attacks
  - Poisoning (training time)
  - Evasion (inference time)
  - Privacy (inference time)
- Multiple application domains: image classification, speech recognition, NLP, cyber security
- Attacks in the real world can have serious consequences
- Defenses are usually domain specific and have limitations

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!