# Can you trust your machine learning system?

*

## Sandip Kundu
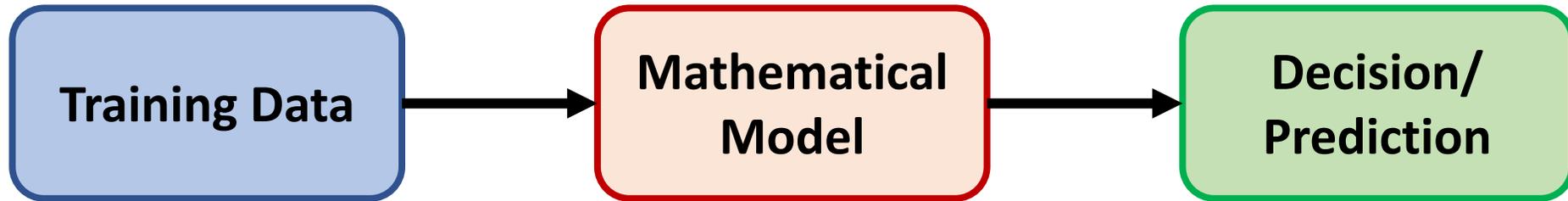
**National Science Foundation**

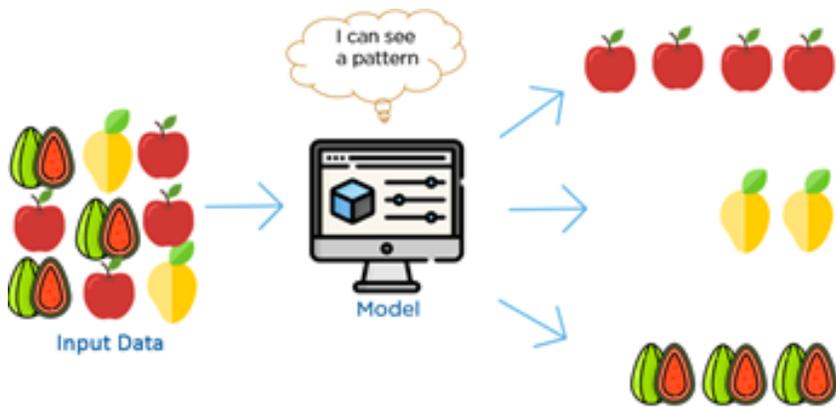on leave from

University of Massachusetts, Amherst
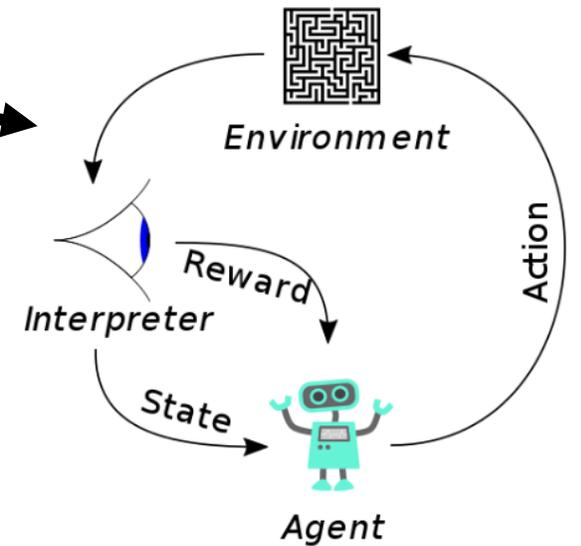
# Machine Learning is becoming Ubiquitous

**Self-driving Cars**



**Cybersecurity**



**Healthcare**



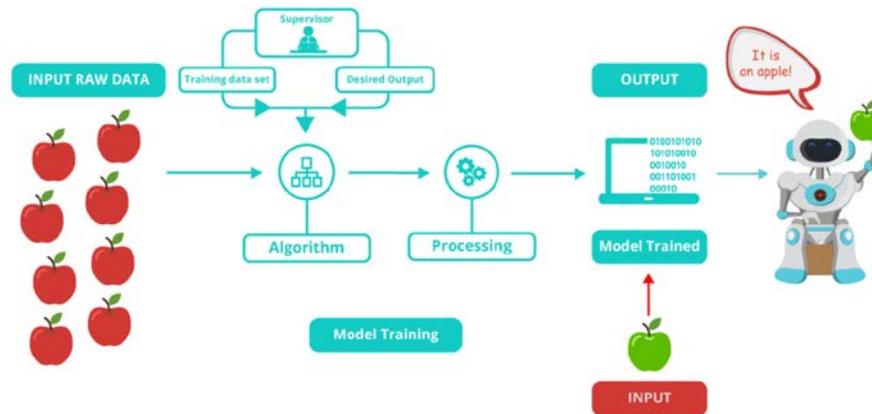**Facial Recognition**



© Sachin Farfade/Mohammad Saberian/Li-Jia Li

**Speech Recognition**

# Machine Learning
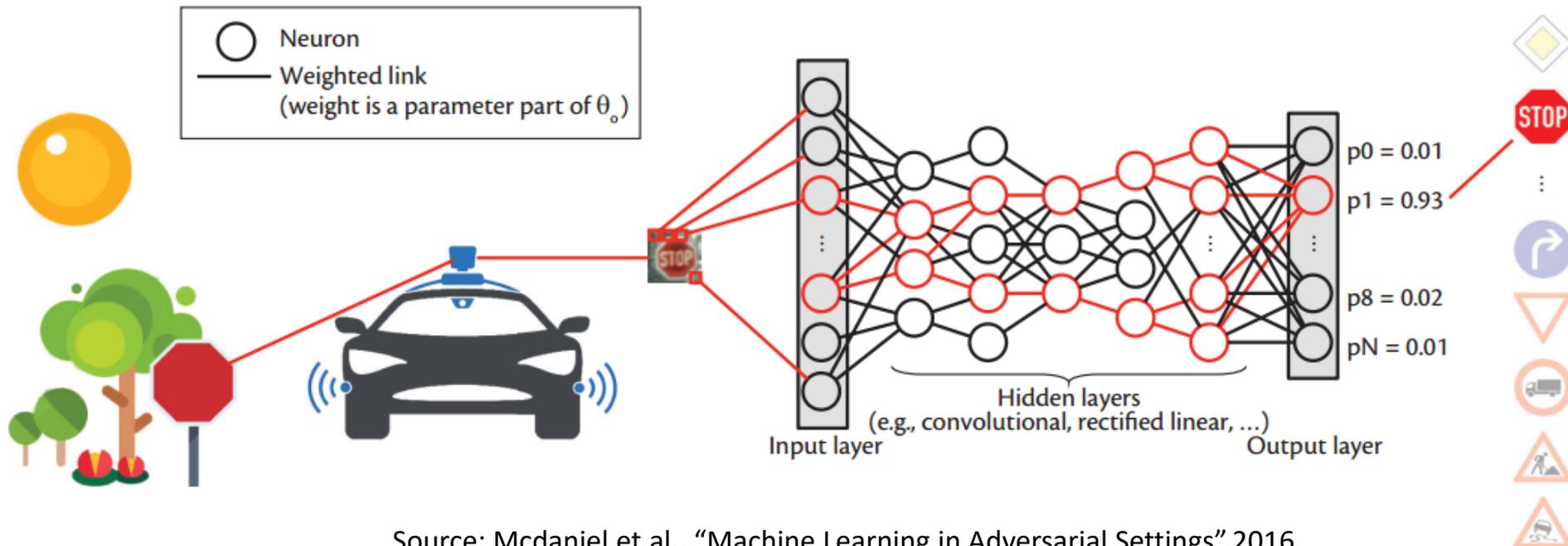
# Self-driving Cars

❖ Cars incorporating systems to assist or replace drivers
  o Ex. automatic parking, Waymo

❖ Self-driving cars with ML infrastructure will become commonplace
  o Ex. NVIDIA DRIVETM PX 2 – open AI car computing system



Source: Mcdaniel et.al., "Machine Learning in Adversarial Settings",2016.

# Healthcare Applications

❖ Diagnosis in Medical Imaging

❖ Treatment Queries and Suggestions

❖ Drug Discovery

❖ Personalized Medicine



* Simm, Jaak, et al. "Repurposing high-throughput image assays enables biological activity prediction for drug discovery." *Cell chemical biology* (2018)



* A Esteva et.al., "Dermatologist-level classification of skin cancer with deep neural networks",2017.

# Cybersecurity

## Spam Filtering



* http://www.thenonprofittimes.com/news-articles/rate-legit-emails-getting-caught-spam-filters-jumped/

## Biometrics ID



* https://www.tutorialspoint.com/biometrics/biometrics_overview.htm

## Intrusion Detection System



## Malware Detection



**Signature - based**

**Anomaly - based**

# Facial Recognition

❖ **Secure Authentication and Identification**
  ○ Apple FaceID
  ○ FBI database – criminal identification

❖ **Customer Personalization**
  ○ Ad targeting
  ○ Snapchat



* Posterscope, Ouividi EYE Corp Media, Engage M1 – GMC Arcadia



Calista_Flockhart_0002.jpg
Detection & Localization

Frontalization:
@152X152x3

C1:
32x11x11x3
@142x142

M2:
32x3x3x32
@71x71

C3:
16x9x9x32
@63x63

L4:
16x9x9x16
@55x55

L5:
16x7x7x16
@25x25

L6:
16x5x5x16
@21X21

**F7:**
**4096d**

F8:
4030d

Taigman et.al.,"DeepFace: Closing the Gap to Human-Level Performance in Face Verification",2014

# Other Machine Vision Applications


Google Lens

❖ **Digital annotation** of real-world

- o Text, language recognition – E.g. Billboards, auto-translation
- o Geo-tagging Landmarks
- o Integration with other services – E.g. ratings for restaurant, directions



❖ **Augmented Reality**

- o **Gaming** – adaptive integration with real-world
- o **Augmented Retail** – E.g. Clothes Fitting

# Speech Recognition

❖ Envisioned in science fiction since 1960's
  o HAL 9000, Star Trek

❖ Natural Language Processing (NLP) has gained increased importance
  o Modeling large vocabularies, accents – translation, transcription services
  o **Smartphones** – Apple Siri, Google Assistant, Samsung Bixby
  o Home - Amazon's Echo/Alexa,
  o IBM Watson



*http://nlp.stanford.edu/~wcmac/papers/20140716-UNLU.pdf*

# Machine learning (ML) Process

# Machine Learning Security and Privacy

# Introduction

❖ ML algorithms in real-world applications mainly focus on **accuracy** (effectiveness) **or/and efficiency** (dataset, model size)

  o Few techniques and design decisions to keep the ML models *secure and robust*!

❖ Machine Learning as a Service (MLaaS) and Internet of Things (IoT) further complicate matters

  o Attacks can **compromise millions of customers**' security and privacy

  o Concerns about **Ownership** of data, model

# ML Vulnerabilities

❖ Key vulnerabilities of machine learning systems
  o ML models often derived from **fixed datasets**
  o Assumption of similar distribution between training and real-world data
    ▪ **Coverage** issues for complex use cases
    ▪ Need **large datasets**, **extensive data annotation**, **testing**

❖ Strong adversaries against ML systems
  o ML algorithms **established** and **public**
  o Attacker can leverage ML knowledge for **Adversarial Machine Learning** (AML)
    ▪ **Reverse engineering** model parameters, test data – **Financial incentives**
    ▪ **Tampering** with the trained model – **compromise security**

# Classification of Security and Privacy Concerns

❖ Attacker's Goals

- **extract** model **parameters**
  **(model extraction)**

- **extract** **private data**
  **(model inversion)**

- **compromise** model to produce false positives/negatives
  **(model poisoning)**

- **produce** adversary selected outputs
  **(model evasion)**

- **render** model **unusable**

❖ Attacker's Capabilities

- access to Black-box ML model

- access to White-box ML model

- manipulate *training data* to **introduce** vulnerability

- access to query to ML model

- access to query to ML model with confidence values

- access to training for building model

- **find and exploit** vulnerability during *classification*

# Security and Privacy Concerns

# Model Extraction

# Model Extraction Attack

❖ Model **IP ownership** - **primary source of value** for company/ service

❖ **Attacker's Capabilities:**
   o Access to black-box model
   o Access to query to ML model

❖ **Goal:** Learns close approximation, $f'$, of $f$ using as few queries as possible
   o Service provider prediction APIs themselves used in attack
      ▪ APIs return extra information – **confidence scores**



f'(x) = f(x) on 100% of inputs
100s-1000's of online queries

amazon web services™   big ml®

Attack   x   Model f   Data

f'   f(x)

- Logistic Regressions, Neural Networks, Decision Trees, SVMs
- **Reverse-engineer model type & features**

* Tramer et.al., "Stealing Machine Learning Models via Prediction APIs.", 2016.

# Extraction Countermeasures

❖ **Restrict information** returned
  - o E.g. do not return confidence scores
  - o **Rounding** – return approximations where possible

❖ **Strict** query **constraints**
  - o E.g. disregard incomplete queries

❖ **Ensemble methods**
  - o Prediction = aggregation of predictions from multiple models
  - o Might still be susceptible to *model evasion* attacks

❖ Prediction API minimization is not easy
  - o API should still be useable for legitimate applications

* Tramer et.al., "Stealing Machine Learning Models via Prediction APIs.", 2016.

# Model Inversion

# Training Data Confidentiality

❖ **Training data** is **valuable** and **resource-intensive** to obtain
  o Collection of **large datasets**
  o Data **annotation** and **curation**
  o Data **privacy** in critical applications like healthcare

❖ Ensuring training data **confidentiality** is **critical**



QUARTZ

**Waymo's driverless cars have logged 10 million miles on public roads**

By Jane C. Hu • October 10, 2018



The New York Times

*Sloan Kettering's Cozy Deal With Start-Up Ignites a New Uproar*

By Charles Ornstein and Katie Thomas

Sept. 20, 2018

# Model Inversion Attack

❖ Extract **private and sensitive inputs** by leveraging the outputs and ML model.

❖ **Optimization goal**: Find inputs that maximize returned confidence value to infer sensitive features or complete data points from a training dataset

❖ **Attacker's Capabilities:**

o Access to Black-box or White-box model

o Exploits confidence values exposed by ML APIs



An image recovered using a new model inversion attack (left) and a training set image of the victim (right). The attacker is given only the person's name and access to a facial recognition system that returns a class confidence score.
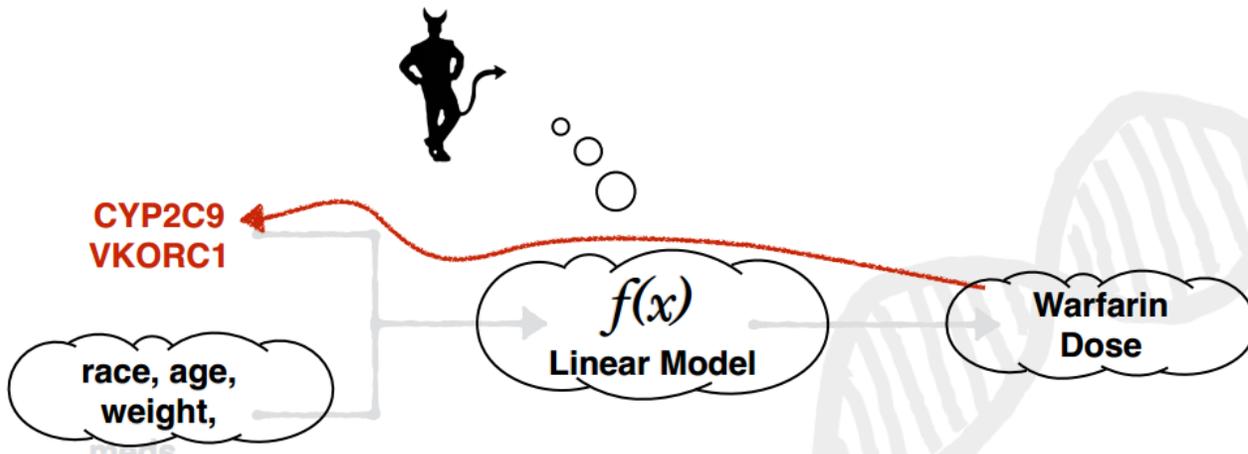
* Fredrikson et.al., "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures.", 2015

# Privacy of the Training or Test Data

❖ **Attacker's capabilities:** Access to query to ML model

❖ Extracting patients' genetics from *pharmacogenetic dosing models*
  o **Queries** using *known information* – E.g. demographics, dosage
  o **Guess** unknown information and check model's response - assign weights
  o Return guesses that produce **highest confidence score**

| age | height | weight | race | history | vkorc1 | cyp2c9 | dose |
|-----|--------|--------|------|---------|--------|--------|------|
| 50-60 | 176.2 | 185.7 | asian | cancer | A/G | *1/*3 | 42.0 |

CYP2C9
VKORC1

race, age, weight,

$f(x)$ **Linear Model**

**Warfarin Dose**

$f(x)$

| age | height | weight | race | history | vkorc1 | cyp2c9 | dose | | |
|-----|--------|--------|------|---------|--------|--------|------|------|------|
| 50-59 | 176.53 | 144.2 | white | | | | 42.0 | 49.7 | p=0.23 |
| 50-59 | 176.53 | 144.2 | white | | | | 42.0 | 42.0 | p=0.75 |
| 50-59 | 176.53 | 144.2 | white | | | | 42.0 | 39.2 | p=0.01 |

$f(x)$

| age | height | weight | race | history | vkorc1 | cyp2c9 | dose | | |
|-----|--------|--------|------|---------|--------|--------|------|------|------|
| 50-59 | 176.53 | 144.2 | white | Cancer | A/G | *1/*1 | 42.0 | 49.7 | p=0.23 |
| 50-59 | 176.53 | 144.2 | white | Heart | G/G | *1/*3 | 42.0 | 42.0 | p=0.75 |
| 50-59 | 176.53 | 144.2 | white | Diabetes | A/A | *2/*3 | 42.0 | 39.2 | p=0.01 |

Fredrikson et.al., "Privacy in Pharmacogenetics: An End-to-End Case Study of Personalized Warfarin Dosing ", 2014.

# Training Data Tampering

❖ **Attacker's goal:** Leaking information about training data by modifying training algorithm

❖ **Attacker's capabilities:**

- Provides tampered APIs that remembers too much information
- Access to Black-box model
  - Extending the training dataset with additional synthetic data
- Access to white-box model
  - Encoding sensitive information about training data in model parameters



A typical ML training pipeline. Data *D* is split into training set *D*train and test set *D*test.
The dashed box indicates the portions of the pipeline that may be controlled by the adversary
*Song et.al. "Machine Learning Models that Remember Too Much", 2017.

# Inversion Countermeasures

❖ Incorporate model inversion metrics to increase robustness
  - o **Identify** sensitive features
  - o Analyze **effective feature placement** in algorithm – E.g. sensitive features at top of a *decision tree* maintain accuracy while preventing *inversion* from performing better than guessing
  - o **Approximate**/ **Degrade** confidence score output – E.g. decrease gradient magnitudes
    - ▪ Works against non-adapting attacker


❖ Ensuring privacy needs to be balanced against usability
  - o **Privacy Budget**

❖ **Differential Privacy** mechanisms using added noise
  - o Might prevent model inversion
  - o Risk of compromising legitimate results in critical applications

# A Countermeasure Against Model Inversion

❖ Based on the injection of noise with long-tailed distribution to the confidence levels.

❖ The small randomness added to the confidence information **prevents convergence** for model inversion attack, but does not affect functionality

❖ **No modification or re-training** of model required

**Noise distribution long tail**

# Targeted Misclassification

❖ Misclassification to a target class

  o Visually same-looking images are classified differently

  o Target adversarial examples are obtained using our numerical implementation of gradient descent based attack.

**Original:** bird - 99.9%    **Adversarial:** cat - 94.0%
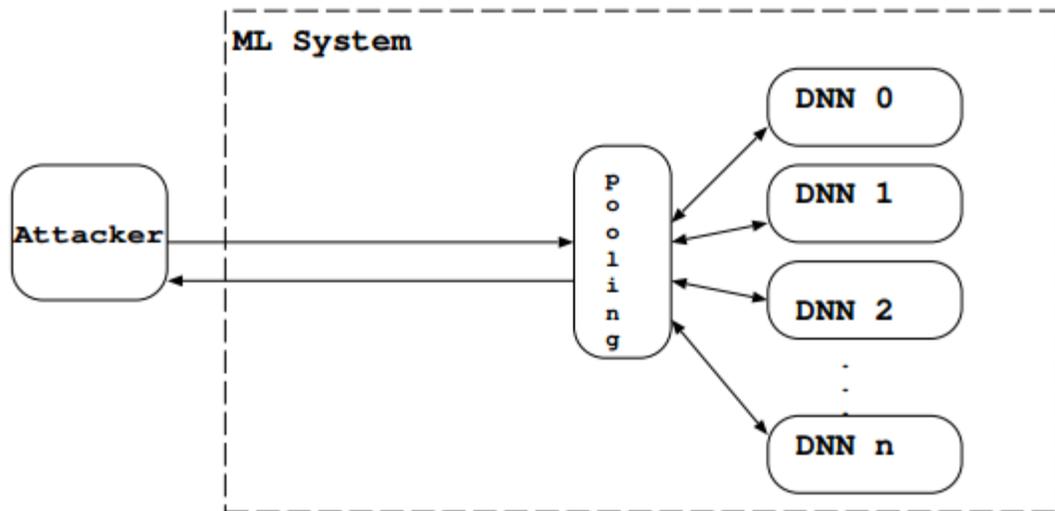
**Original:** frog - 99.8%    **Adversarial:** ship - 80.1%

Adversarial examples. Original images (left) and the target adversarial examples (right). Below each image is the classification and confidence returned by the ResNet CIFAR-10 Image Classifier.

# A Countermeasure Against Targeted Misclassification

❖ Varying the order of the training

- o Different models which offer the same classification accuracy, yet they are different numerically.

❖ An ensemble of such models

- o Allows to randomly switch between these equivalent models during query which further blurs the classification boundary.



Workflow description of adversarial attacks with Multi-Model Defense applied.

Adversarial attack performed on an image originally classified as *deer*, where the target class *is truck*. With Noise-Injection defense, the attack does not converge and ends up degrading the original image.

# Model Poisoning and Evasion

# Model Poisoning and Evasion Attacks

❖ Ensuring Integrity of a Machine Learning model is difficult
  o Dependent on **quality** of *training*, *testing* datasets
    ▪ Coverage of *corner cases*
    ▪ Awareness of *adversarial examples*

  o **Model sophistication** – E.g. small model may produce incorrect outputs

  o **Lifetime management** of larger systems
    ▪ Driverless cars will need constant updates
    ▪ Degradation of input sensors, training data pollution

❖ Adversarial examples may be **Transferable** *
  o Example that fools Model A might fool Model B
  o Smaller model used to find examples quickly to target more sophisticated model

Papernot et. al., "Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples", 2016

# Model Poisoning and Evasion Attacks

❖ **Adversary capabilities:** Causing misclassifications of attacks to appear as normal (false positives/ negatives)

- Attack on training phase: **Poisoning (Causative) Attack**: Attackers attempt to learn, influence, or corrupt the ML model itself
  - Compromising data collection
  - Subverting the learning process
  - Degrading performance of the system
  - Facilitating future evasion

- Attack on testing phase: **Evasion (Exploratory) Attack**: Do not tamper with ML model, but instead cause it to *produce adversary selected outputs by manipulating test samples*.
  - Finding the blind spots and weaknesses of the ML system to evade it
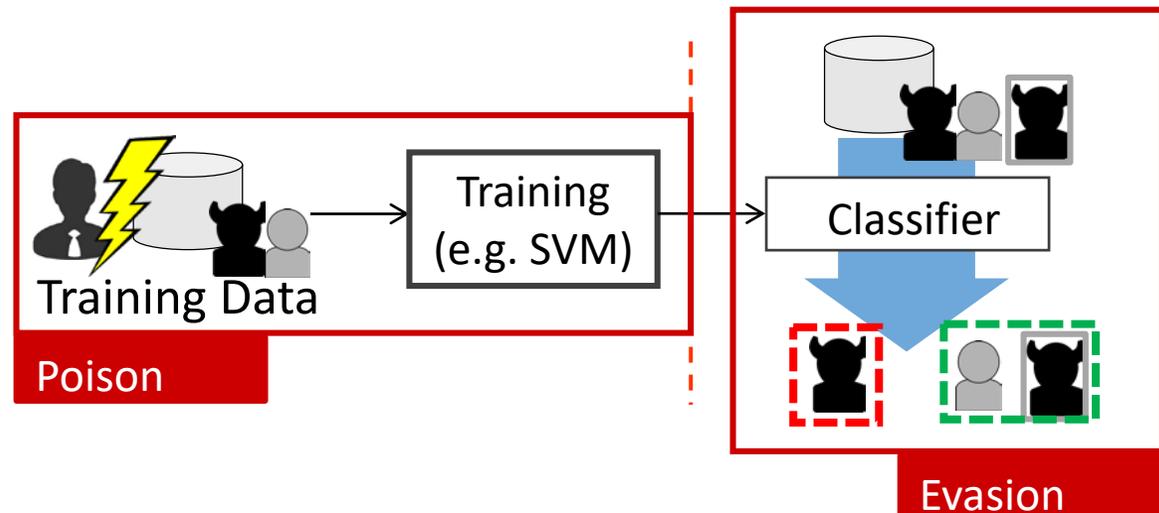
# Adversarial Detection of Malicious Crowdsourcing

❖ Malicious crowdsourcing, or *crowdturfing* used for tampering legitimate applications

- Real users paid to promote malicious intentions
- Product reviews, Political campaigns, Spam

❖ Adversarial machine learning attacks

- Evasion Attack: workers evade classifiers
- Poisoning Attack: crowdturfing admins tamper with training data



Wang et.al., "Man vs. Machine: Adversarial Detection of Malicious Crowdsourcing Workers ", 2014

# Physical Perturbations

❖ Adversarial perturbations detrimentally affect Deep Neural Networks (DNNs)
  - Cause misclassification in critical applications
  - Requires some knowledge of DNN model
  - Perturbations can be robust against noise in system

❖ Defenses should not rely on physical sources of noise as protection
  - Incorporate adversarial examples
  - Restrict model information/ visibility
  - **DNN Distillation** – transfer knowledge from one DNN to another
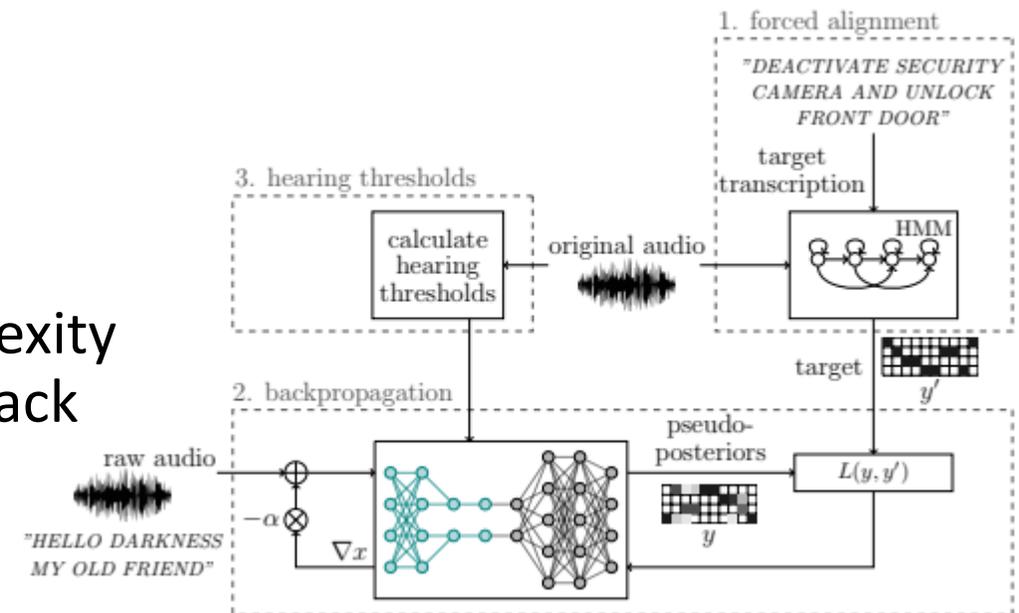  - **Gradient Masking**



Eykholt et.al., "Robust Physical-World Attacks on Deep Learning Visual Classification", 2018.

Papernot et.al., "Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks", 2015.

# Adversarial Attacks Against ASR DNNs

❖ Automatic Speech Recognition (**ASR**) and Natural Language Understanding (**NLU**) increasingly popular – E.g. Amazon Alexa/ Echo

   o Complex model = **Large parameter space** for attacker to explore

❖ Attacker goals

   o Psychoacoustic hiding – perceived as noise by human

   o Identify and match legitimate voice features

      ▪ Pitch, tone, fluency, volume, etc

   o Embed arbitrary audio input with a malicious voice command

   o *Temporal alignment* dependencies add complexity

   o Environment/ System *variability* can affect attack

   o Software tools like *Lyrebird* can prove useful



Lea et.al., "Adversarial Attacks Against Automatic Speech Recognition Systems via Psychoacoustic Hiding", 2018

# Defenses Against AML

❖ Evasion
  ○ Multiple classifier systems (B. Biggio et al., IJMLC 2010)
  ○ Learning with Invariances (SVMs)
  ○ Game Theory (SVMs)

❖ Poisoning
  ○ Data sanitization (B. Biggio et al., MCS, 2011)
  ○ Robust learning (PCA)
  ○ Randomization, information hiding, security by obscurity

❖ Randomizing collection of training data (timings / locations)
  ○ using difficult to reverse-engineer classifiers (e.g., MCSs)
  ○ denying access to the actual classifier or training data
  ○ randomizing classifier to give imperfect feedback to the attacker (B. Biggio et al., S+SSPR 2008)

# Towards Robust ML Model

# Future Research Areas

❖ Complexity of Machine Learning itself an issue
- o New attacks models constantly emerging – *timely detection* critical
- o Generation and incorporation of **Adversarial Examples**
- o **Data Privacy** is crucial to enhance ML security
  - ▪ *Differential Privacy* has tradeoffs
  - ▪ *Homomorphic Encryption* still nascent

❖ Security introduces overhead and can affect performance
- o **Optimizations** needed to ensure ML effiency

❖ Tools to increase robustness of Machine Learning need research
- o *Unlearning, re-learning*
- o *ML Testing*
- o *Sensitivity Analysis*

# Unlearning and Re-learning

❖ Ability to **unlearn** is gaining importance
  - ○ **Pollution** attacks or carelessness – *Mislabeling* and *Misclassification*
    - ▪ Large changing datasets difficult to maintain
    - ▪ Anomaly detection not enough
  - ○ **EU GDPR** regulations – **Privacy**
  - ○ **Completeness** and **Timeliness** are primary concerns *
  - ○ **Statistical Query Learning*** and **Causal Unlearning**** proposed in literature
  - ○ Suitable for **small deletions**

❖ **Re-learning** or **Online learning**
  - ○ Faces similar issues to un-learning
  - ○ Can be very slow
  - ○ More suitable for large amounts of deletions or new information

* Yinzhi Cao, "Towards Making Systems Forget with Machine Unlearning", 2015
** Cao *et. al.*, "Efficient Repair of Polluted Machine Learning Systems via Causal Unlearning", 2018

# Sensitivity Analysis

❖ Study of how the uncertainty in the output of a system can be attributed to different sources of uncertainty in its inputs
  o ML feature extraction sensitivity analysis well-researched


❖ Detection of **biases** in training/test datasets is crucial *
  o Model accuracy dependent on datasets used – ***real-world*** performance can be different
    ▪ Datasets can have **expiration dates**
    ▪ **Privacy** issues can render datasets incomplete
  o Identify training datasets which **generalize** better
  o Study sensitivity of ML accuracy to change in datasets

* Sanders, Saxe, "Garbage In, Garbage Out - How Purportedly Great ML Models Can Be Screwed Up By Bad Data", 2017

# Thank you