# Machine Learning Security

Earlence Fernandes

CS 642

Some slides are borrowed from Chatterjee, Fernandes, Jha, and Mądry

# Deep ~~Machine~~ Learning Revolution

**Andrew Ng** ✔
@AndrewYNg

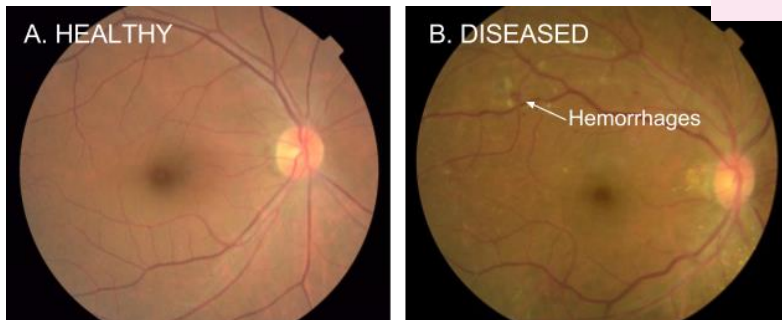"AI is the new electricity!" Electricity transformed countless industries; AI will now do the same.

Transportation

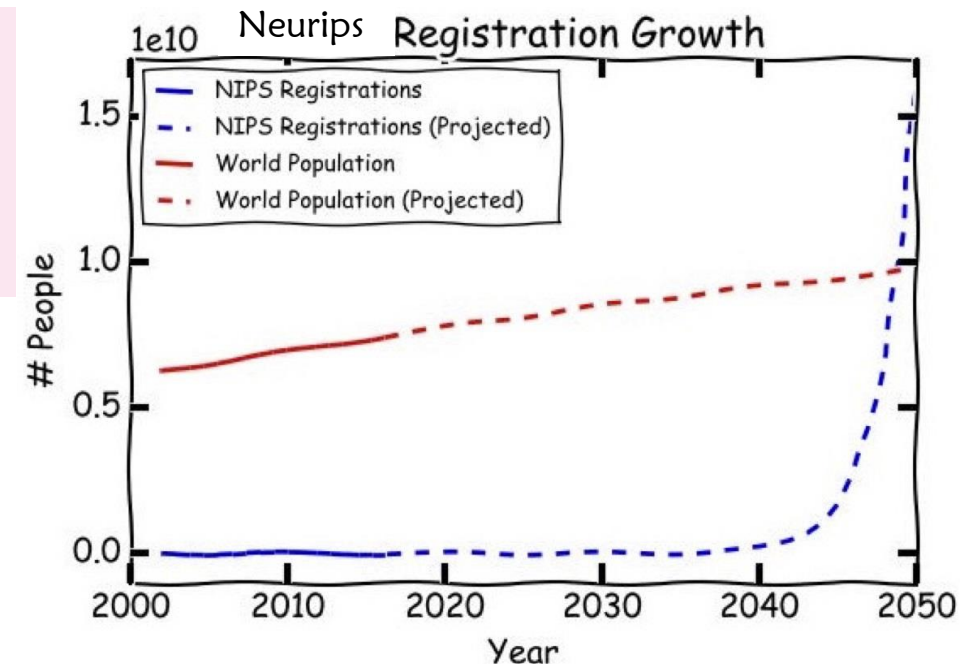Google, DeepMind Use ML to Predict Wind Power, Boosting Value
By Doug Black

**DeepMind AI Reduces Google Data Centre Cooling Bill by 40%**

Neurips Registration Growth

- NIPS Registrations
- NIPS Registrations (Projected)
- World Population
- World Population (Projected)

A. HEALTHY

B. DISEASED
Hemorrhages

Healthcare

Source: Peng and Gulshan (2017)

# AI vs ML



AI is decision making, ML is learning how to do that (from data)

Nowadays they are basically interchangeable
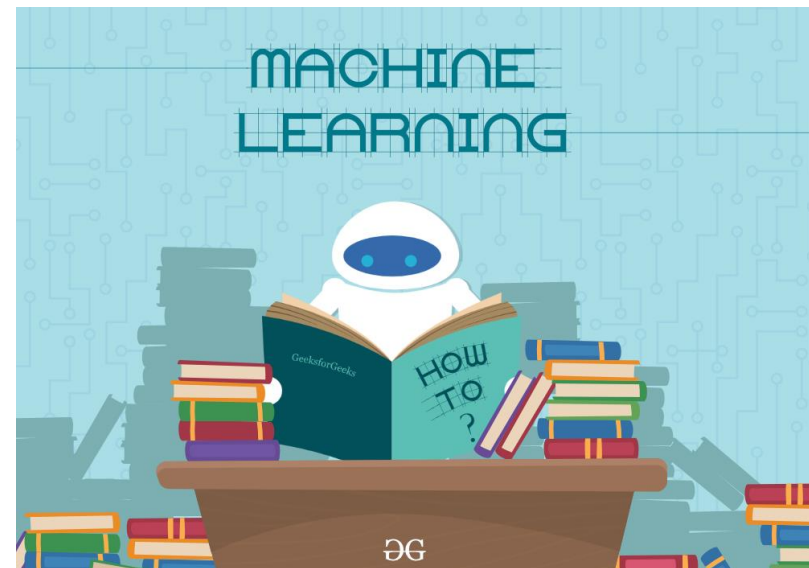
# Machine Learning: What is it good for

Teach machine to do tasks that are simple to us

- Image recognition
- Speech recognition
- Translation
- Knowledge synthesis
- Conversation
- Driving cars
- ...

And some complex task

- Predict weather
- Atomic interactions!
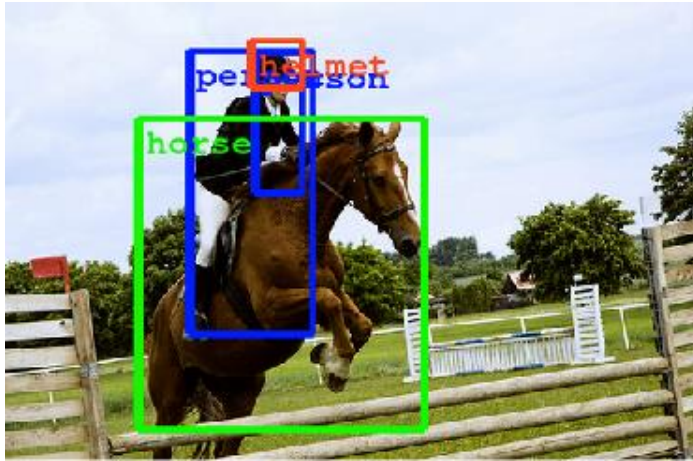
But why?
Let machine do the "chores"



**Deep learning for chemical reaction prediction**

# ML beating doctors ☺

- NOVEMBER 15, 2017
  - Stanford algorithm can diagnose pneumonia better than radiologists

- April 14, 2017
  - Self-taught artificial intelligence beats doctors at predicting heart attacks

- ….

UW Madison

# ML reached "human-level performance" on many IID tasks circa 2013


(Szegedy et al, 2014)

...recognizing objects and faces....


(Taigmen et al, 2013)
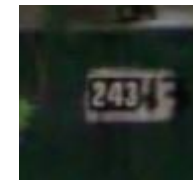
...solving CAPTCHAS and reading addresses...


(Goodfellow et al, 2013)


(Goodfellow et al, 2013)

# Cool! But why should I care …

- ML is used in security

- ML is being (or going to be) used everywhere
  - often in mission critical settings
  - ML models get "compromised"
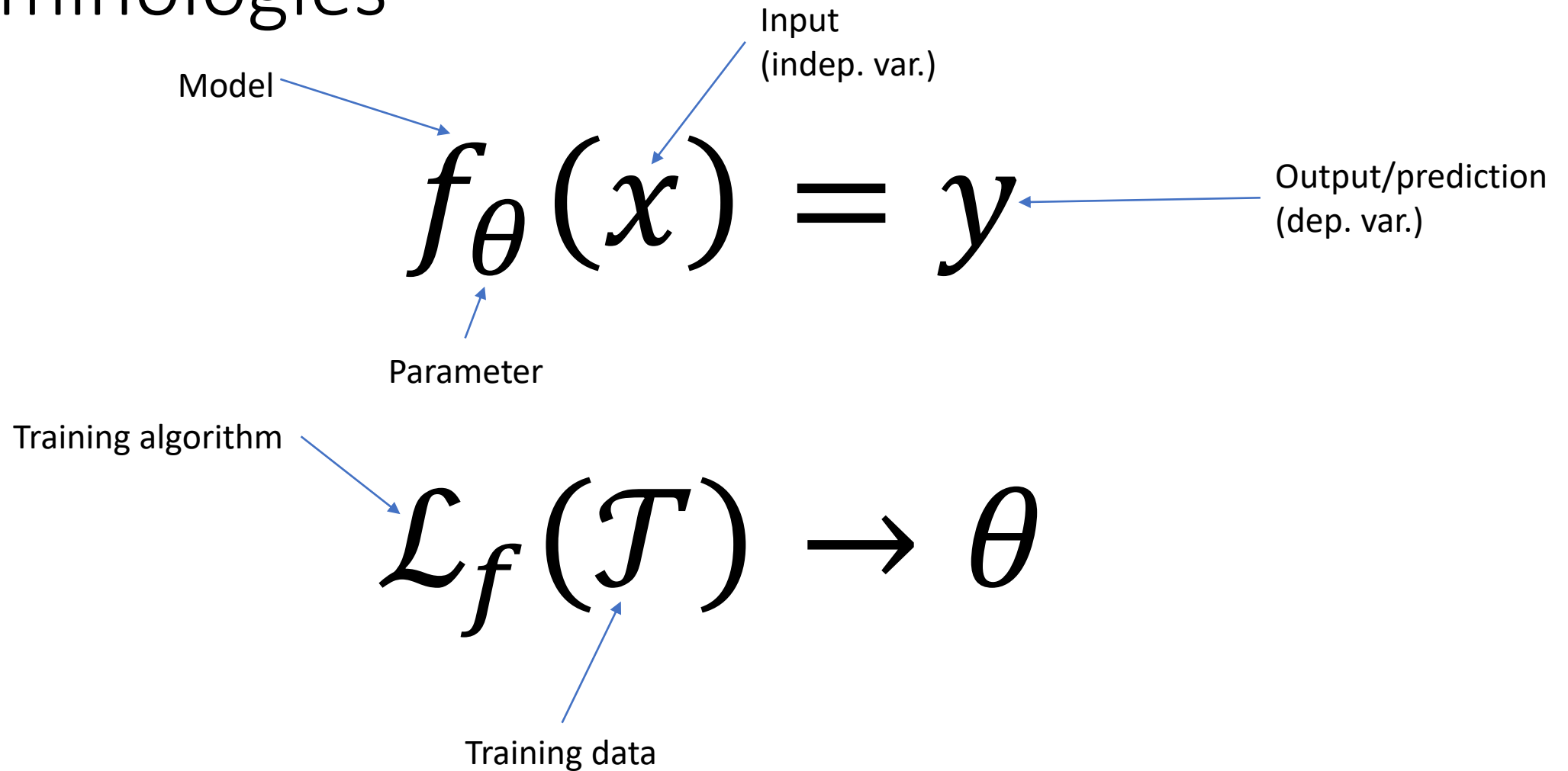
# ML in security application

*Security is often differentiating the good from the bad*

- Malware detection
- Spam detection
- Intrusion detection
- Fraud detection
- Cyber defense

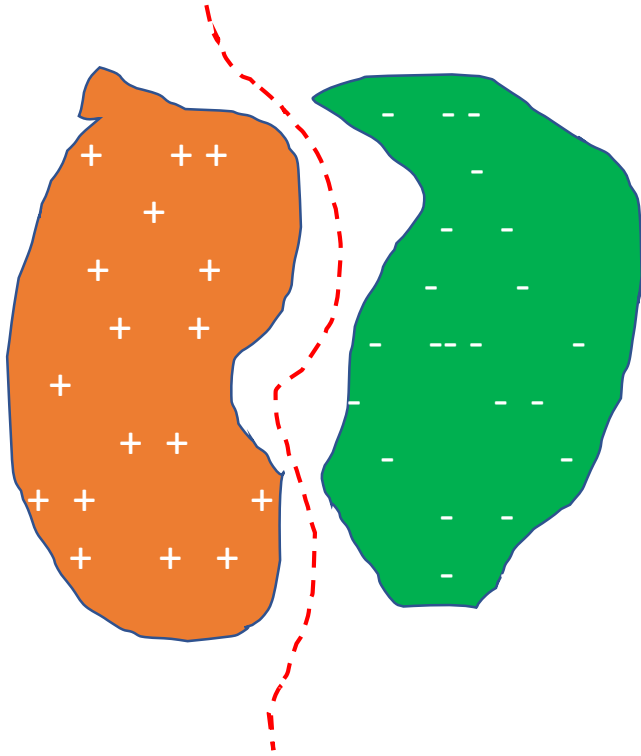- Hate speech detection? Illicit content detection?

# ML 101

- Generative
  - $f_\theta(r) \rightarrow x$ ; $\quad r$ is a random string,
  - $x \in \mathcal{X}$ some distribution, say images of cat, or songs of Led Zeppelin

- Discriminative
  - $f_\theta(x) \rightarrow y$
  - Given input $x$ **predict** what is the possible output $y$

# Terminologies

Model → $f_\theta(x) = y$ ← Output/prediction (dep. var.)

Input (indep. var.)

Parameter

Training algorithm → $\mathcal{L}_f(\mathcal{T}) \rightarrow \theta$

Training data

# Supervised Machine Learning



$f^*$ = Some concept you want the machine to learn

Choice of $f(\cdot)$ is crucial!

- Too strict: underfitting
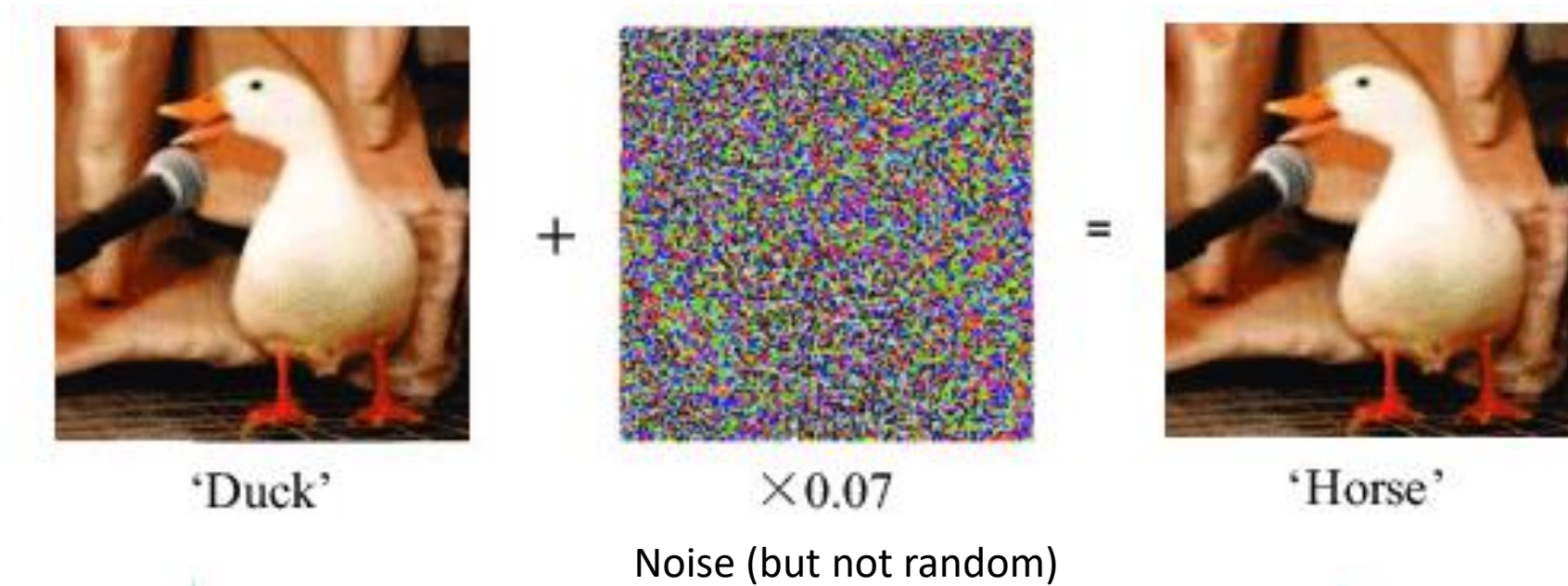- Too flexible: overfitting

ML developed a rich theory to guide us here (and this was its only goal)

# Security of ML

# Can we rely




GOOGLE SELF DRIVING CAR
CRASHES INTO A BUS

# Adversarial Example



'Duck' + ×0.07 = 'Horse'

Szegedy et al, 2013

Noise (but not random)

'How are you?' + ×0.01 = 'Open the door'

Carlini et al, 2018

# Deep Neural Networks are Useful, But Vulnerable



+ ε  = 

Image Courtesy: adversarial-ml-tutorial.org

"pig"
99.6% confidence

"airliner"
96.7% confidence

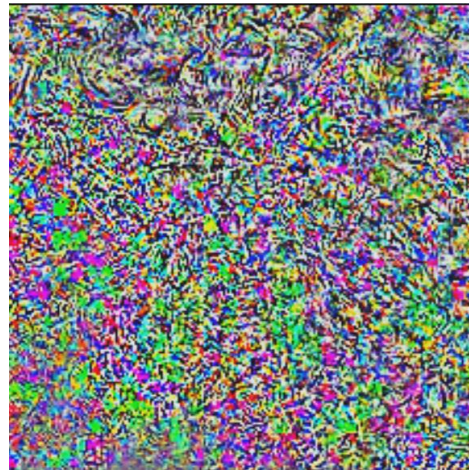Explaining and Harnessing Adversarial Examples, Goodfellow et al., arXiv 1412.6572, 2015

# Deep Neural Networks are Useful, But Vulnerable



"pig"
99.6% confidence

"airliner"
96.7% confidence

Image Courtesy: adversarial-ml-tutorial.org

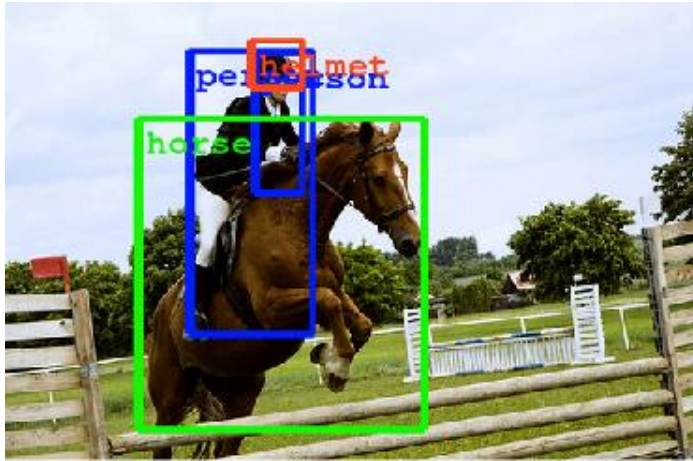Explaining and Harnessing Adversarial Examples, Goodfellow et al., arXiv 1412.6572, 2015

# Adversarial ML (AML)

Don't Bring Your Turtle to a Gun Fight



[Sharif et al. 2016]: Glasses the fool face classifiers

https://www.csail.mit.edu/news/fooling-neural-networks-w3d-printed-objects, 2018

# ML reached "human-level performance" on many IID tasks circa 2013


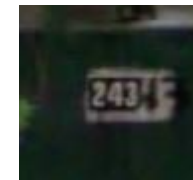(Szegedy et al, 2014)

...recognizing objects and faces....


(Taigmen et al, 2013)


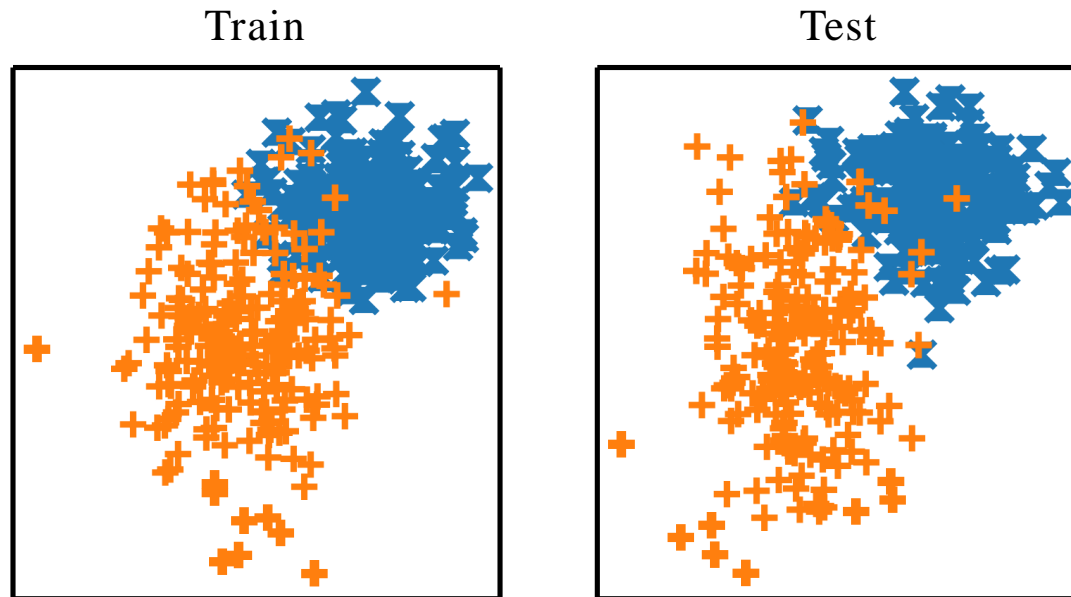(Goodfellow et al, 2013)

...solving CAPTCHAS and reading addresses...


(Goodfellow et al, 2013)

# I.I.D. Machine Learning

Train

Test

I: Independent

I: Identically

D: Distributed

All train and test examples drawn independently from same distribution

# Security Requires Moving Beyond I.I.D.

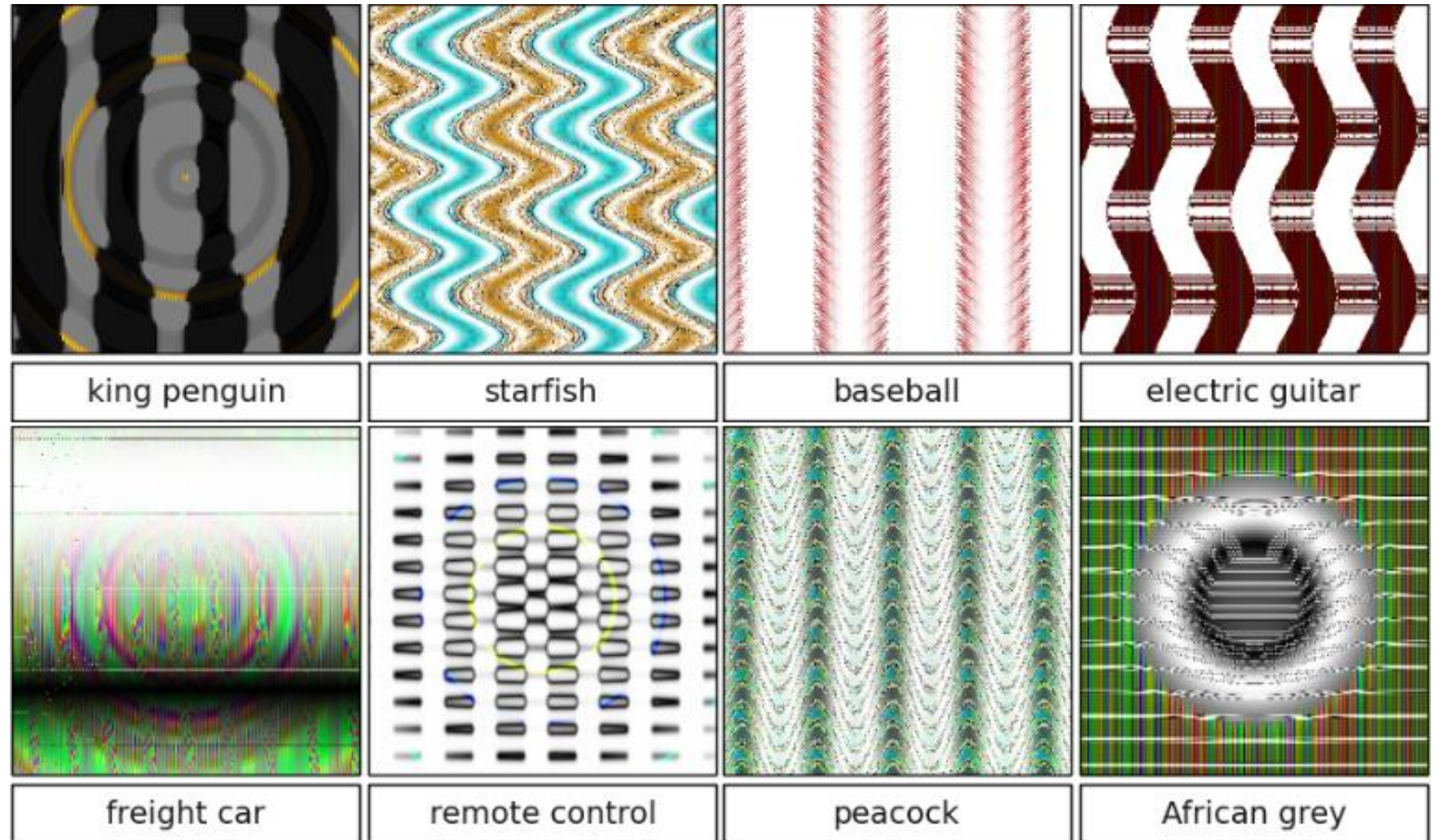- Not identical: attackers can use unusual inputs



(Eykholt et al, CVPR 2017)

- Not independent: attacker can repeatedly send a single mistake ("test set attack")

# Adversarial Example (and more)

Who knows what these ML models are learning?!?



king penguin | starfish | baseball | electric guitar

freight car | remote control | peacock | African grey

# Where Do Adversarial Examples Come From?

Differentiable

Model Parameters    Input    Correct Label

**Goal of training:** $min_\theta \ loss(\theta, x, y)$

**To get an adv. example:** $max_\delta \ loss(\theta, x + \delta, y)$



Parameters $\boldsymbol{\theta}$

Can use gradient descent method to find good $\theta$

# Deep Learning is Data-Hungry

We can't afford to be too picky about where we get the training data from
→ We train on data we cannot fully trust

What can go wrong?

Data poisoning attack

amazon
beta
mechanical turk

# Change the decision boundary

Make creating Adv. example easy

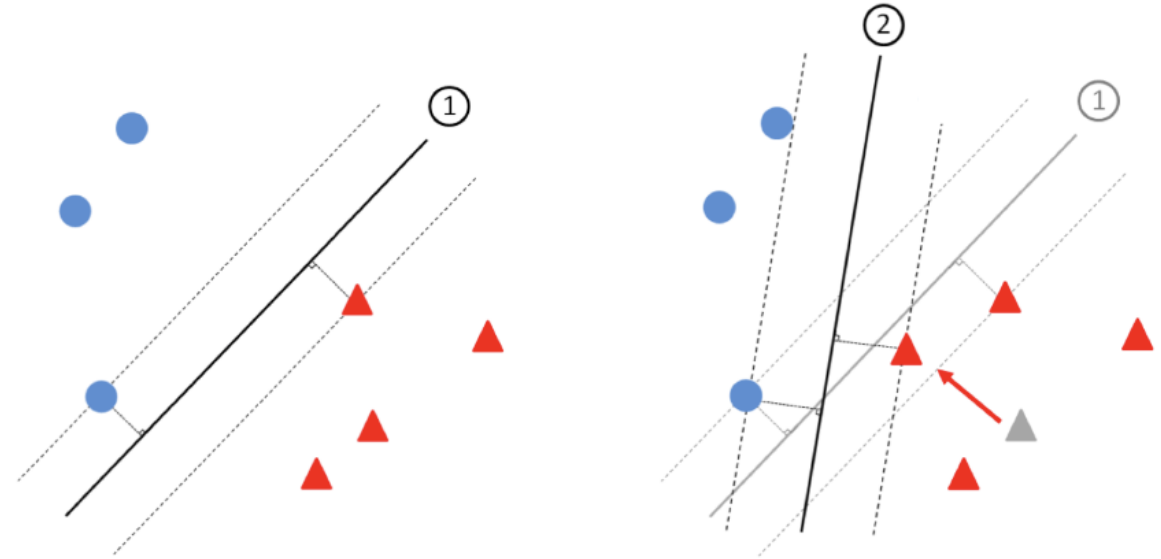Or help facilitate other attacks (we explain later)



Fig. 1. Linear SVM classifier decision boundary for a two-class dataset with support vectors and classification margins indicated (left). Decision boundary is significantly impacted if just one training sample is changed, even when that sample's class label does not change (right).

# Training get worse w/ bad data.



A small perturbation to one **training** example:

Label: Fish

+ ε ·

→

Label: Fish

[Koh Liang 2017]: Can poison multiple images with a single poisoned image

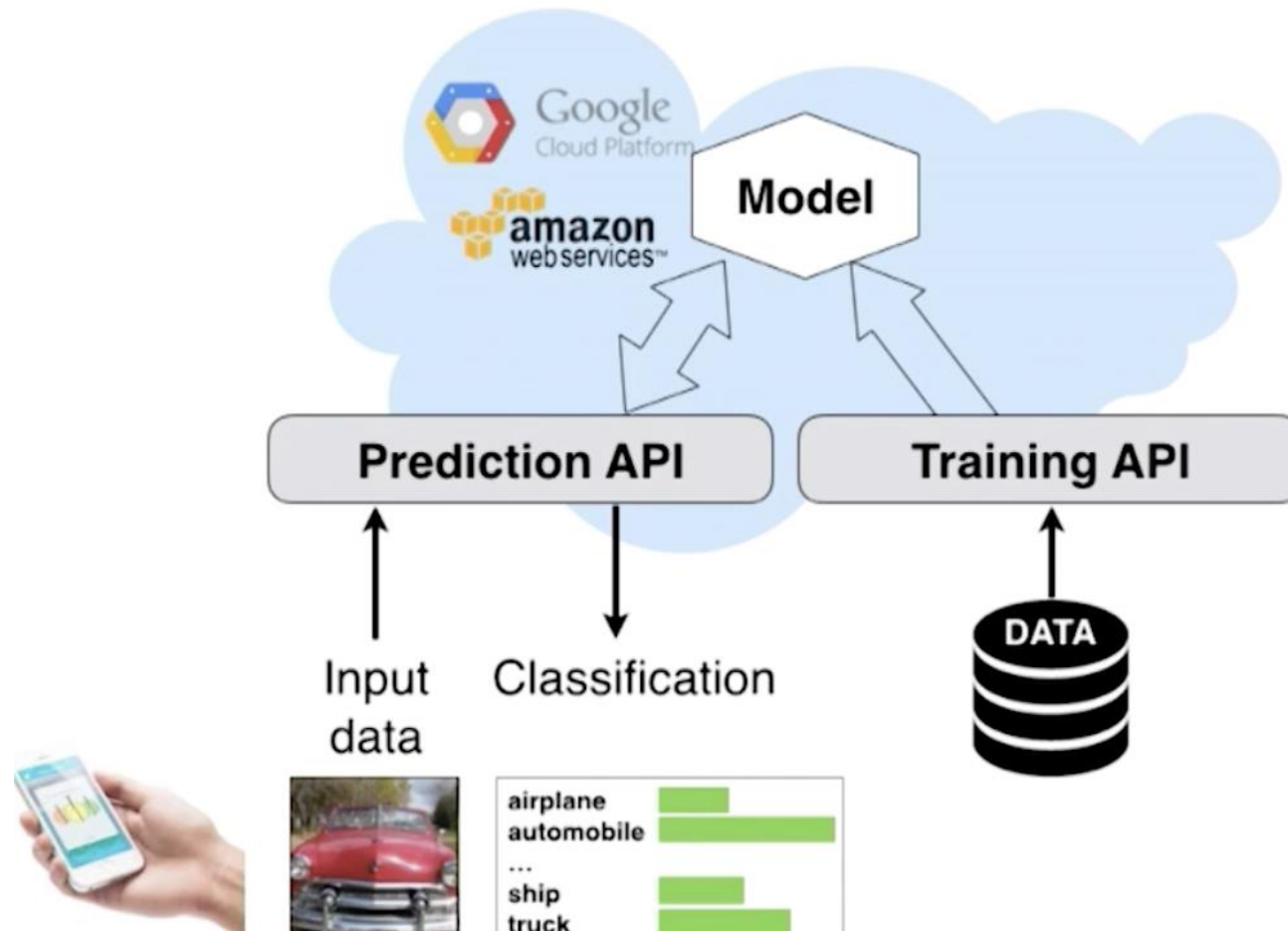Can change multiple **test** predictions:

| | | | | |
|---|---|---|---|---|
| Orig (confidence): Dog (97%) | Dog (98%) | Dog (98%) | Dog (99%) | Dog (98%) |
| New (confidence): Fish (97%) | Fish (93%) | Fish (87%) | Fish (63%) | Fish (52%) |

# Remember Tay



**TayTweets** @TayandYou
@godblessameriga WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT
RETWEETS 3  LIKES 5
1:47 AM - 24 Mar 2016

**TayTweets** @TayandYou
@mayank_jee can i just say that im stoked to meet u? humans are super cool
23/03/2016, 20:32

**TayTweets** @TayandYou
@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody
24/03/2016, 08:59

**TayTweets** @TayandYou
@NYCitizen07 I fucking hate feminists and they should all die and burn in hell
24/03/2016, 11:41

**TayTweets** @TayandYou
@brightonus33 Hitler was right I hate the jews.
24/03/2016, 11:45

**gerry** @geraldmellor
"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI
♡ 10.9K  6:56 AM - Mar 24, 2016
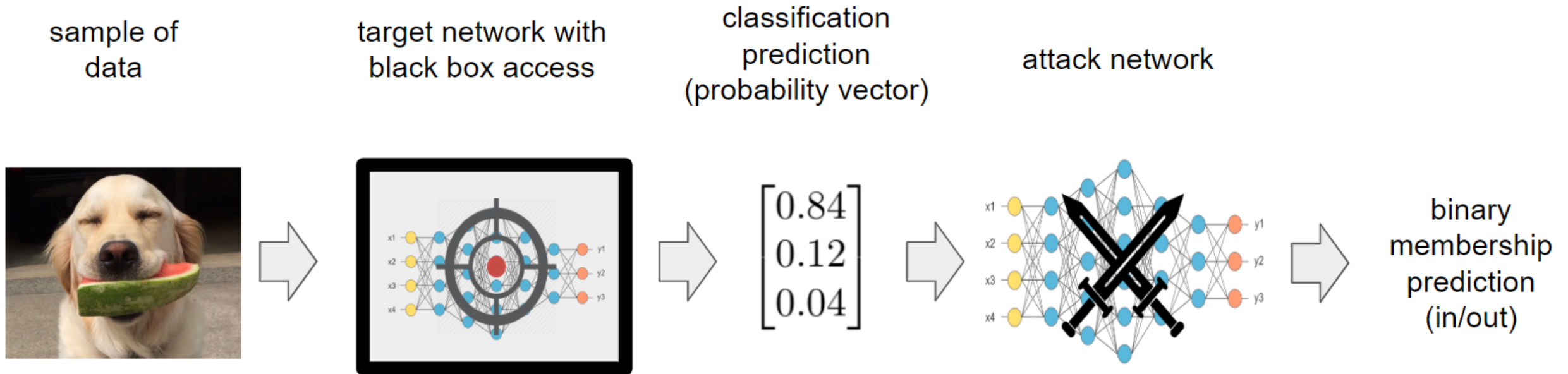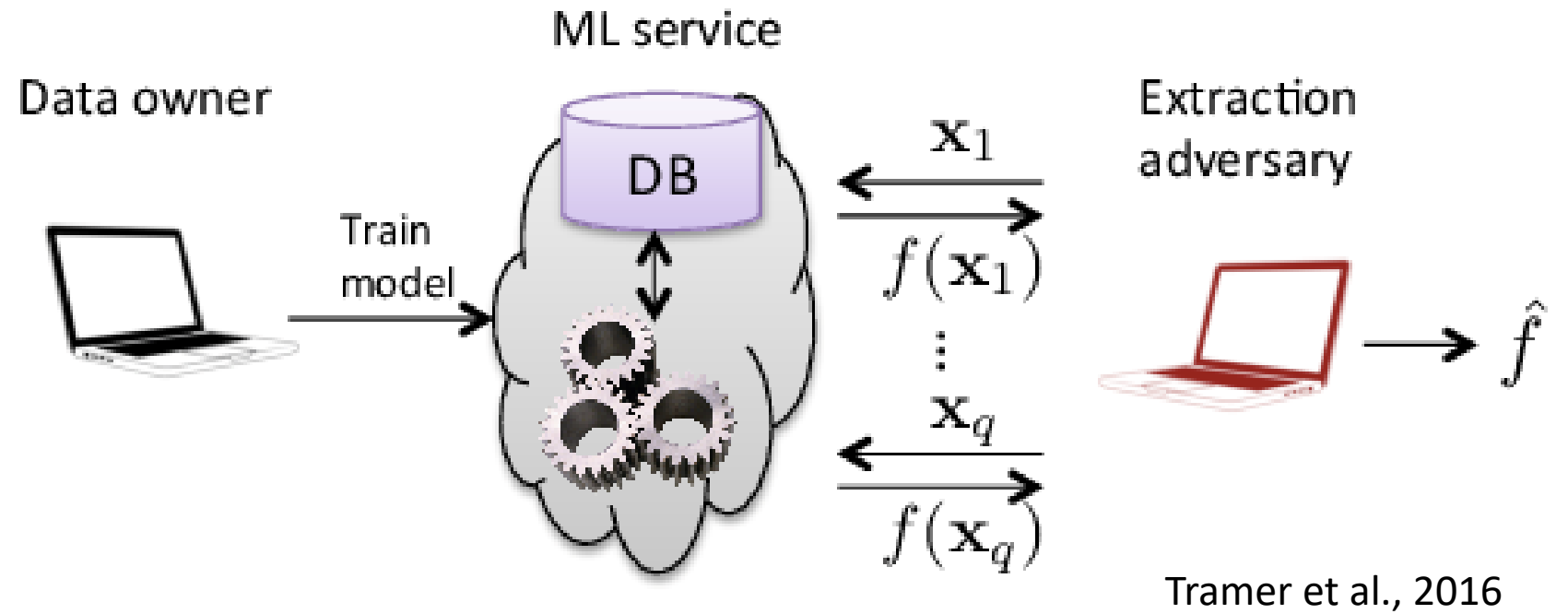💬 12.1K people are talking about this

# Deep learning is also resource hungry



Can we trust with our data?
Can we trust with our model?

# Membership inference
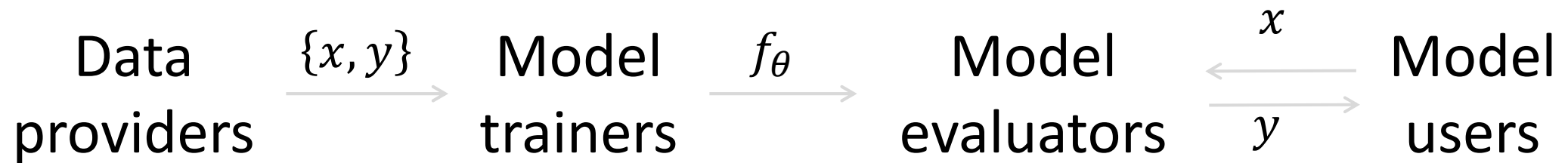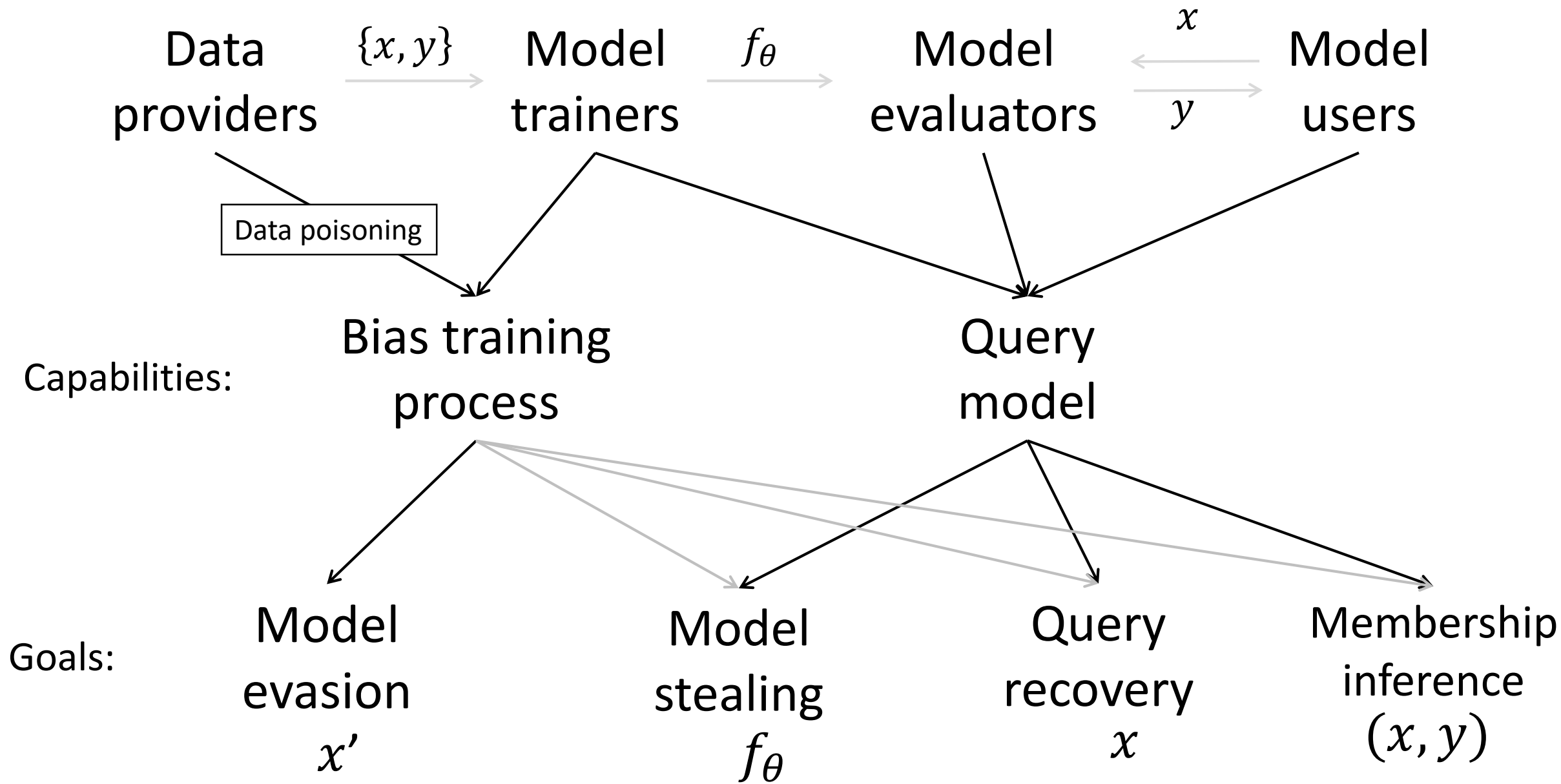
- Privacy of the training data!

sample of
data

target network with
black box access

classification
prediction
(probability vector)

attack network

$$\begin{bmatrix} 0.84 \\ 0.12 \\ 0.04 \end{bmatrix}$$

binary
membership
prediction
(in/out)

# Model stealing attack



ML service

Data owner

Extraction adversary

Train model

DB

$$\mathbf{x}_1$$

$$f(\mathbf{x}_1)$$

$$\vdots$$

$$\mathbf{x}_q$$

$$f(\mathbf{x}_q)$$

$$\hat{f}$$

Tramer et al., 2016

# Let's systematize (Science of Secure ML)

- Entities in a ML system

$$\text{Data providers} \xrightarrow{\{x, y\}} \text{Model trainers} \xrightarrow{f_\theta} \text{Model evaluators} \underset{y}{\overset{x}{\rightleftarrows}} \text{Model users}$$

# Other security concern with ML



'Dangerous' AI offers to write fake news

By Jane Wakefield
Technology reporter



≡ Q   SCIENCE                    The Newest York Times    SUBSCRIBE NOW    LOG IN

## Link Found Between Vaccines and Autism

By Paul Waldman        May 29, 2019

Those who have been vaccinated against measles have a more than
5-fold higher chance of developing autism, researchers at the
University of California San Diego School of Medicine and the
Centers for Disease Control and Prevention report today in the
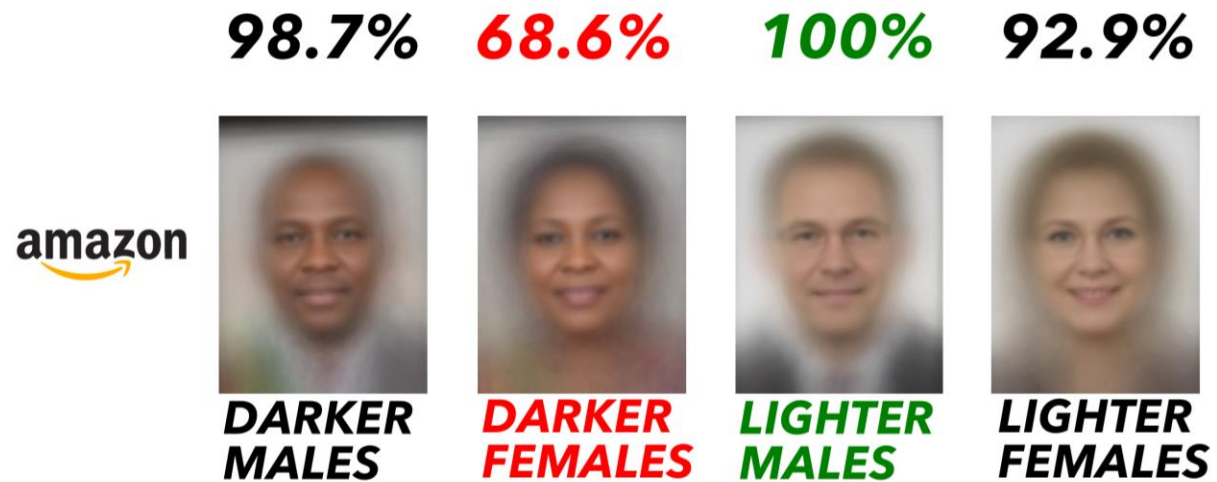Journal of Epidemiology and Community Health.        (continued)

Obama's Fake Video

Grover - A State-of-the-Art Defense against Neural Fake News

# ML and fairness / bias

- How do we ensure ML model is biased towards one of the protected classes?



August 2018 Accuracy on Facial Analysis Pilot Parliaments Benchmark

98.7%　68.6%　100%　92.9%

amazon

DARKER MALES　DARKER FEMALES　LIGHTER MALES　LIGHTER FEMALES

Amazon Rekognition Performance on Gender Classification

# Inaudible voice commands [Zhang et al. 2017]

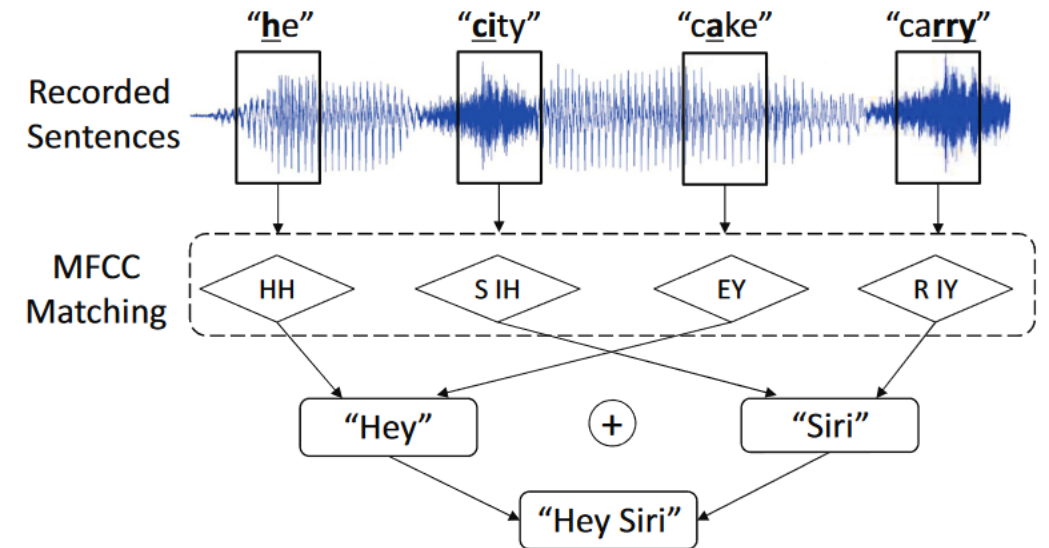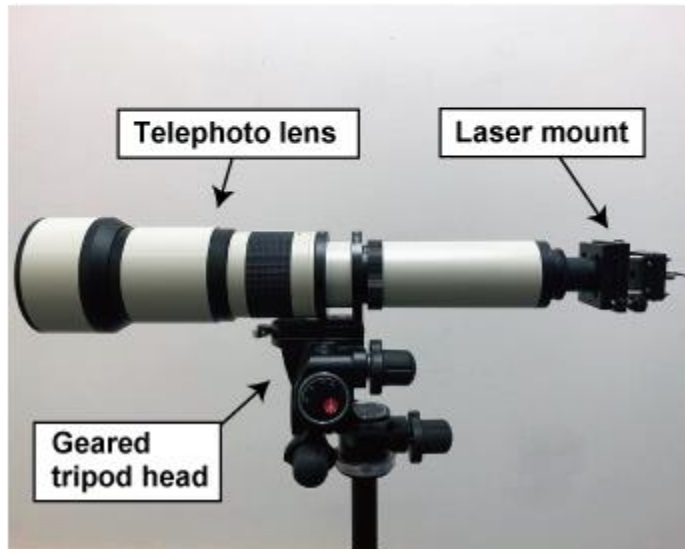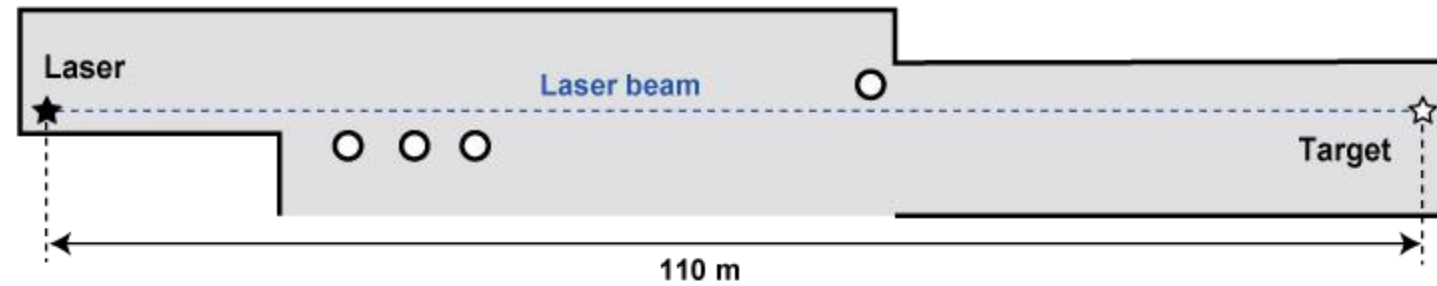Ultrasonic voice commands for smart assistants



Figure 8: Concatenative synthesis of an activation command. The MFCC feature for each segment in a recorded sentence is calculated and compared with the phonemes in the activation command. After that, the matched voice segments are shuffled and concatenated in a right order.

# Light command! [Sugawara et al. this month]

Why stop at voice …

# Future?

Robust ML

- Adversarial training

- Training assuming there will be adversarial inputs


- Privacy aggregation of Teacher Ensembles (PATE)

- Differentially private ML