

# NOTES

---

- Call for actions
  - 11/09 Lecture: Recording will be offered (on 11/16)
  - 11/20: Checkpoint II review deadline (on HotCRP)

# CS 499/579: TRUSTWORTHY ML

## PRELIMINARIES ON PRIVACY

Tu/Th 4:00 – 5:50 pm

Sanghyun Hong

[sanghyun.hong@oregonstate.edu](mailto:sanghyun.hong@oregonstate.edu)

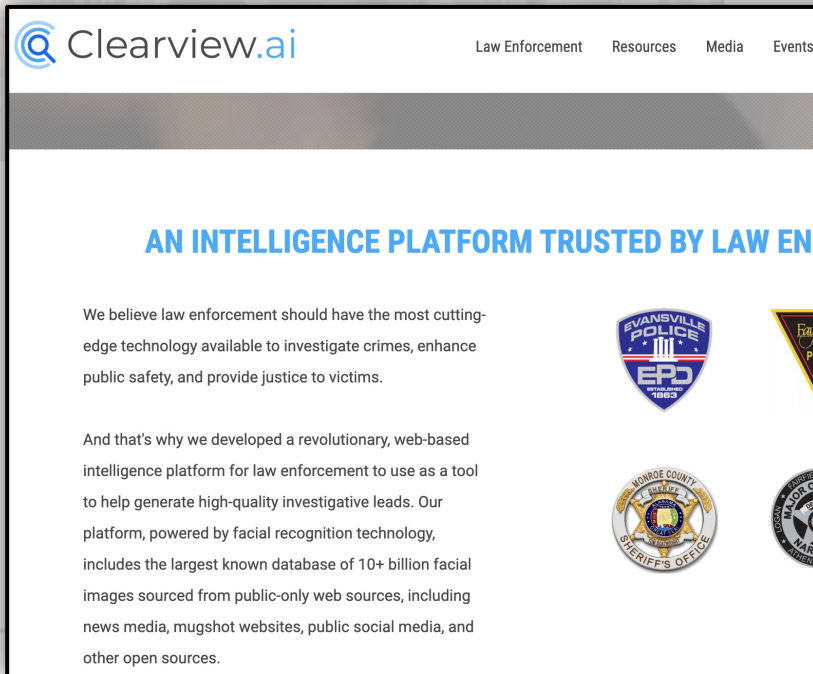


**Oregon State**  
University

**SAIL**  
Secure AI Systems Lab

**PRIVACY, PRIVACY, PRIVACY...**

# WHY PRIVACY MATTERS?




Clearview.ai

Law Enforcement Resources Media Events

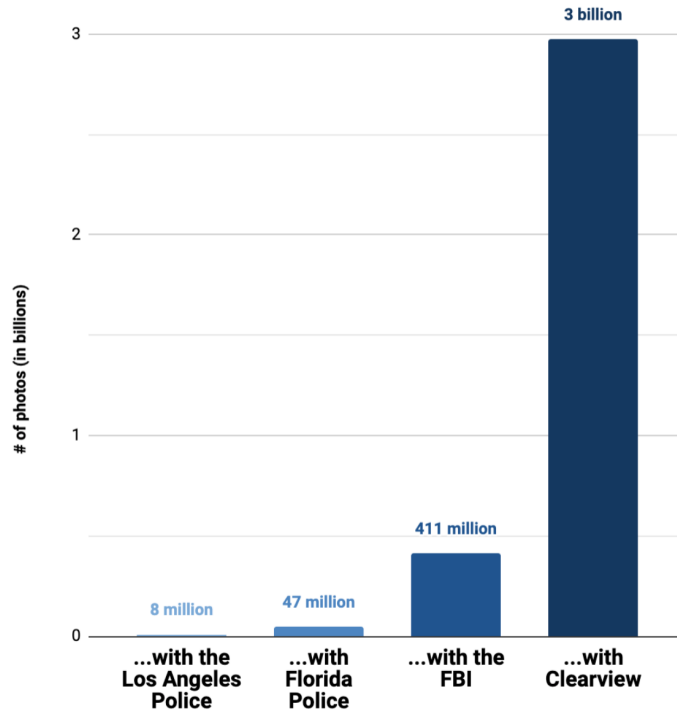
## AN INTELLIGENCE PLATFORM TRUSTED BY LAW ENFORCEMENT

We believe law enforcement should have the most cutting-edge technology available to investigate crimes, enhance public safety, and provide justice to victims.

And that's why we developed a revolutionary, web-based intelligence platform for law enforcement to use as a tool to help generate high-quality investigative leads. Our platform, powered by facial recognition technology, includes the largest known database of 10+ billion facial images sourced from public-only web sources, including news media, mugshot websites, public social media, and other open sources.



## This is how many photos you can search...



<sup>1</sup><https://www.nytimes.com/2020/01/18/technology/clearview-privacy-facial-recognition.html>

<sup>2</sup><https://www.muckrock.com/news/archives/2020/jan/18/clearview-ai-facial-recognition-records/>

# WHY PRIVACY MATTERS?

---

- Let's discuss
  - What is privacy?
  - What does privacy matter?
  - How is it different from security?

# HOW CAN WE MAKE IT PRIVATE?

---

- A perfect, yet not interesting solution:
  - No learning ... but this is *not* what we want
  - Hold-on, what if we anonymize some records?

# DE-ANONYMIZATION

---

- Setup
  - **Attacker:** de-anonymize anonymized records
  - **Victim** : anonymize sensitive data records
- Knowledge
  - Additional (or auxiliary information) about the data
- Capability
  - Query your data with some techniques
  - Perform post-processing computations on  $q$  (outputs)
  - ... (many more)

# DE-ANONYMIZATION

---

- In ML
  - We train statistical models
  - It does not matter whether data is anonymized or not
  - Some examples
    - Cancer data
    - Demographics
    - Data about people's financial information
    - ...
- Note:
  - “Anonymization of a data record might seem easy to implement. Unfortunately, it is *increasingly easy to defeat* anonymization by the very techniques that are being developed for many legitimate applications of big data.” [1]

[1] President’s Council of Advisors on Science and Technology, 2014  
Narayanan and Shmatikov, Robust De-anonymization of Large Sparse Datasets, IEEE S&P 2008



# HOW CAN WE MAKE IT PRIVATE?

---

- Shannon's perfect security
  - An adversary should not distinguish a message  $M$  from a random text  $R$



Claude Shannon (1916 ~ 2001)  
A Father of Information Theory  
and Modern Cryptography

# HOW CAN WE MAKE IT PRIVATE?

---

- Shannon's perfect security

- An adversary should not distinguish a message  $M$  from a random text  $R$

- Formally:

- $\Pr[M = m | C = c] = \Pr[M = m]$

- where

- $m$  is a message (from a set  $M$ )

- $c$  is a ciphertext (from a set of all ciphertexts  $C$ )

- $\Pr[C = c | M = m] = \Pr[C = c]$

- It means:

- Ciphertext provides no additional information

- Observing  $c$  does not help with guessing  $M = m$

- $c$  is independent of the message  $m$



Claude Shannon (1916 ~ 2001)  
A Father of Information Theory  
and Modern Cryptography

# HOW CAN WE MAKE IT PRIVATE?

---

- Perfect security in model training

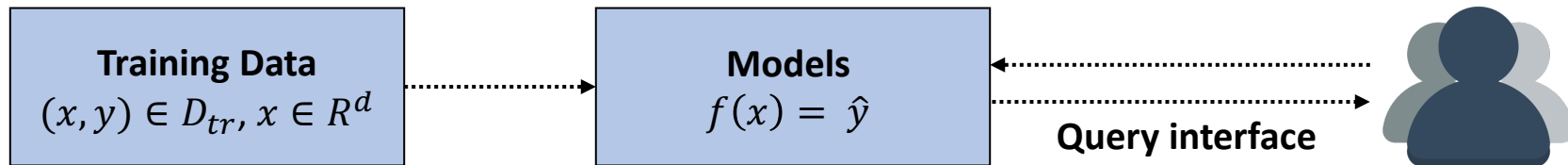


- Potential solutions:
  - **Encrypt-decrypt:** encrypt the training data and decrypt it to train a model
  - **Homomorphic encryption:** encrypt the training data and train a model on it
  - ...

# HOW CAN WE MAKE IT PRIVATE?

---

- Inferences with such model(s)

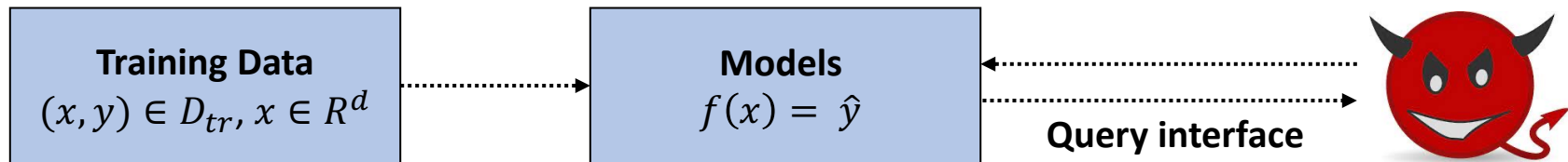


- Potential problems:
  - Perfect security-based solutions are computationally expensive (than vanilla training)
  - Only a limited number of users (who has a key) may use these models

# HOW CAN WE MAKE IT PRIVATE?

---

- Inferences with such model(s)



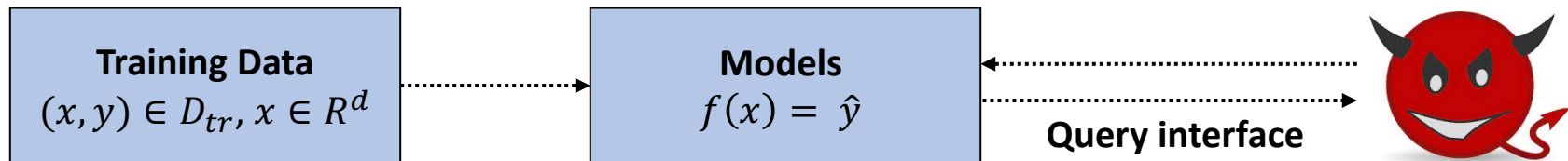
- Potential problems:
  - Perfect security-based solutions are computationally expensive (than vanilla training)
  - Only a limited number of users (who has a key) can use these models
  - Once a key is leaked, an adversary can query the model with any data

# WHAT AN ADVERSARY CAN DO WITH THE QUERY ACCESS?

# PRIVACY THREAT MODEL

---

- ML Pipeline

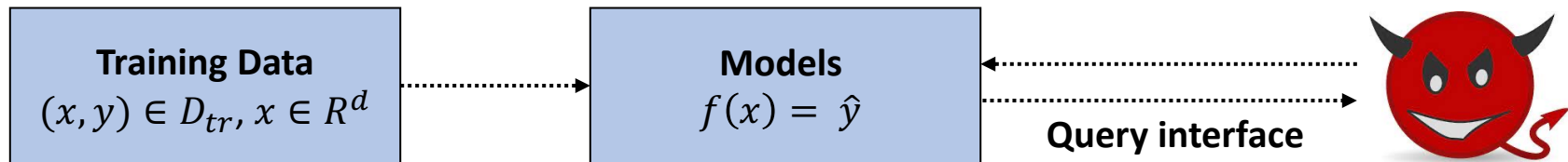


- Privacy risks

- Identify your membership in the training data
- Identify (sensitive) properties of your training data
- Identify (sensitive) attribute of a person that you know
- Reconstruct a sample completely
- Reconstruct a model behind the query interface
- ...

# PRIVACY THREAT MODEL

- ML Pipeline



- Privacy risks (from the view of the work by Dwork *et al.*)

- Tracing attack : Identify your membership in the training data
- Reconstruction : Identify (sensitive) properties of your training data
- De-anonymization: Identify (sensitive) attribute of a person that you know
- Reconstruction : Reconstruct a sample completely
- Reconstruction : Reconstruct a model behind the query interface
- ...



# PRIVACY THREAT MODEL

---

- The attack considers **non-trivial** cases
  - ex. Smoking causes cancer
  - Revealing this information is *not* a privacy attack
  - We know this is correlated without interacting with the target model
  
  - ex. A model trained on a dataset of lung cancer patients
  - ex. The model gets a patient information and returns the probability of getting the cancer
  - ex. We know the Person A is smoking
  - ex. We identify that A is in the dataset (defer the details to later on)
  - It's a *non-trivial* attack as we identify the information about an individual

# MEMBERSHIP INFERENCE: TRACING

---

- Setup

- **Victim:**

- Has a dataset  $x = \{x_1, \dots, x_n\}$  with  $n$ -i.i.d samples where each  $x_i$  is drawn from  $P$  over  $\{\pm 1\}^d$
    - For each query  $M$ , the victim returns the sample mean  $q$  over given sample  $x_i$ 's

- **Attacker:**

- Perform an attack  $A(y, q, z)$  that identify whether a target instance  $y \in \{\pm 1\}^d$  **IN** the dataset  $x$  or not (**OUT**) with  $m$ -i.i.d reference samples  $z = \{z_1, \dots, z_m\}$  and the sample mean  $q$

- **Procedure:**

# RECONSTRUCTION ATTACK I: ATTRIBUTE INFERENCE

---

- Setup
  - **Victim:**
    - For each  $i$ -th instance, the victim has  $(x_i, s_i)$  information
    - $x_i \in \{0, 1\}^d$ : public info. accessible by an adversary and  $s_i$ : is the one-bit secret
  - **Attacker:**
    - Perform an attack  $A$  that reconstructs  $s_i$  by exploiting query outputs  $\hat{q}$  and the public information  $A(x, M(x, s))$ , where the attacker knows  $k > 1$  public attributes
  - **Formally**

# RECONSTRUCTION ATTACK I: ATTRIBUTE INFERENCE

---

- Setup

- **Victim:**

- For each  $i$ -th instance, the victim has  $(x_i, s_i)$  information
    - $x_i \in \{0, 1\}^d$ : public info. accessible by an adversary and  $s_i$ : is the one-bit secret

- **Attacker:**

- Perform an attack  $A$  that reconstructs  $s_i$  by exploiting query outputs  $\hat{q}$  and the public information  $A(x, M(x, s))$ , where the attacker knows  $k > 1$  public attributes

- **Approximation:**

- Linear statistics (*e.g.*, linear SVM, linear regression, ...)
    - Practical constraints (# Queries)
      - Ideally  $2^n$  queries to solve the subset-sum problem
      - Practically, considering the tradeoff btw error and accuracy, we can do it in polynomial time

# RECONSTRUCTION ATTACK II: MODEL EXTRACTION

---

- Setup

- **Victim:**

- Has a model  $f(x) = y$  trained on a confidential data
    - For each query  $M$ , the victim returns the output  $y_i$  over given sample  $x_i$ 's

- **Attacker:**

- Perform an attack (i.e., trains a surrogate model  $f'$  that is functionally equivalent to  $f$ )

# Thank You!

Tu/Th 4:00 – 5:50 pm

Sanghyun Hong

<https://secure-ai.systems/courses/MLSec/F23>



**Oregon State**  
University

**SAIL**  
Secure AI Systems Lab