# High Level Computer Vision:
# More Attacks and Defenses on CV

Mario Fritz fritz@cispa.saarland
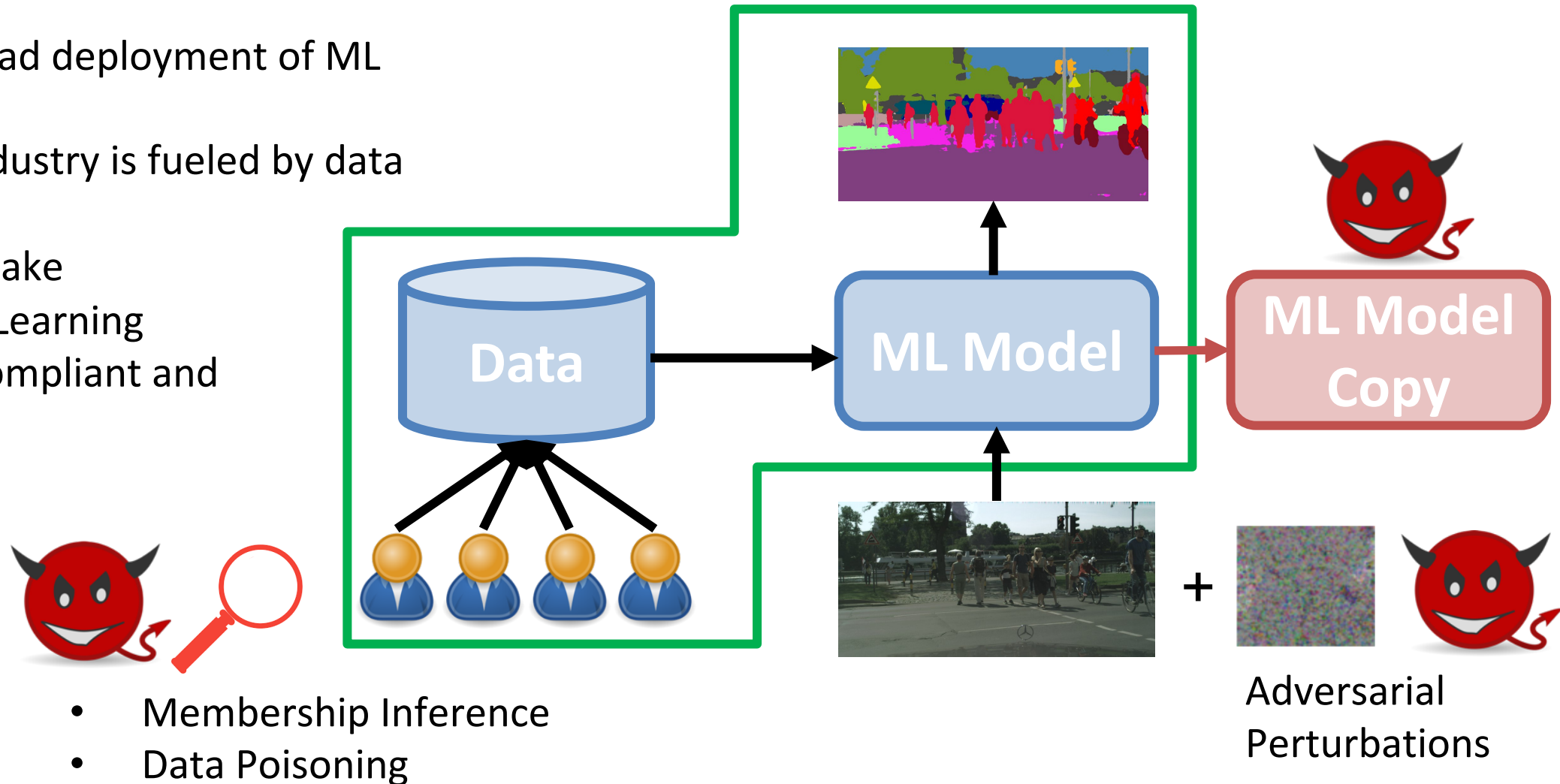
Bernt Schiele schiele@mpi-inf.mpg.de

17.7.2019

# Outline

- Landscape of attacks on Computer Vision Models

- Reverse Engineering and Model Stealing

  – Watermarking

- Adversarial Perturbations

  – Data Poisoning

- Membership Inference

  – Differential Privacy

# Privacy & Security in Machine Learning: Towards Trustworthy AI



- Widespread deployment of ML

- Future industry is fueled by data

- How to make Machine Learning privacy compliant and secure?

- Membership Inference
- Data Poisoning

Adversarial Perturbations

S. Oh; M. Augustin; B. Schiele; M. Fritz; Towards Reverse-Engineering Black-Box Neural Networks; **ICLR'18**
S. Oh; M. Fritz; B.Schiele; Adversarial Image Perturbation for Privacy Protection -- A Game Theory Perspective **ICCV'17**
A. Salem; Y. Zhang; M. Humbert; M. Fritz; M. Backes; ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models **NDSS'19**

K.Grosse, N. Papernot, P.Manoharan, M. Backes, P. D. McDaniel: Adversarial Examples for Malware Detection. **ESORICS'17**
L. Hanzlik; Y, Zhang; K. Grosse; A. Salem; M. Augustin; M. Backes; M.Fritz; MLCapsule: Guarded Offline Deployment of Machine Learning as a Service; **ArXiv'18**
Tribhuvanesh Orekondy; Bernt Schiele; Mario Fritz; Knockoff Nets: Stealing Functionality of Black-Box Models **CVPR'19**
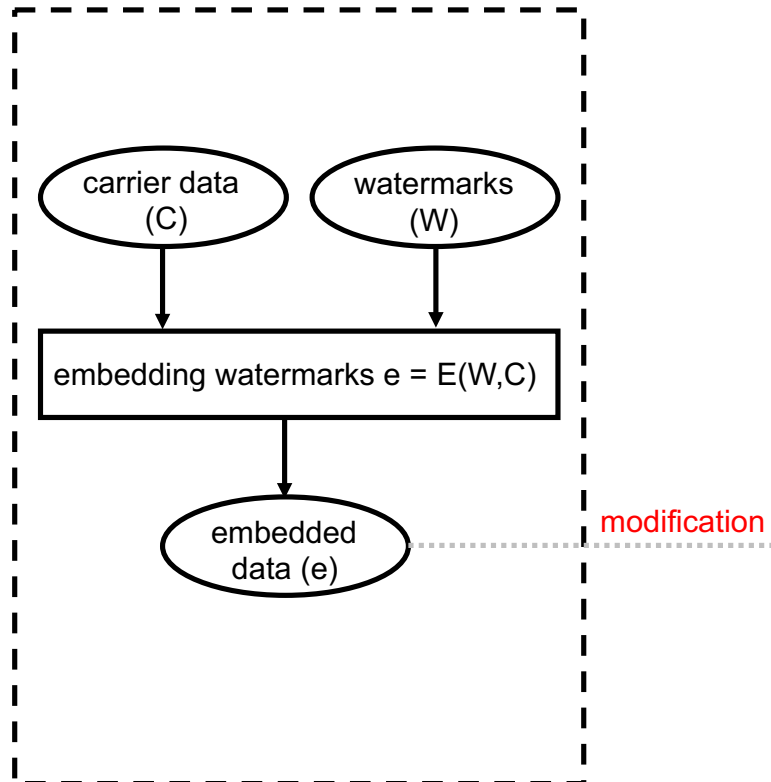
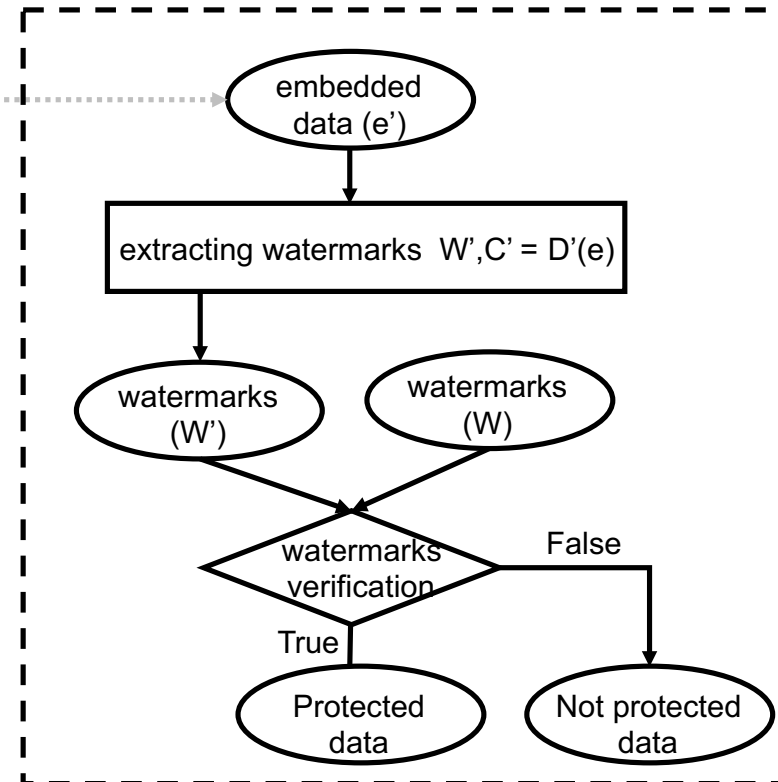# Protecting Intellectual Property of Deep Neural Networks with Watermarking

- *Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph. Stoecklin, Heqing Huang, Ian Molloy*

- *ASIACCS'18*

# Motivation

- AI / ML technology embeddeded into many systems

- Building such models requires:

  - Expertise

  - Data

  - Annotation

  - Computation

- Potential of copyright infringement / IP violations by

  - Illegal reproduction

  - Distributiuon

  - Derivation

- Actual legal situation a bit unclear:

  - Law and Adversarial Machine Learning:
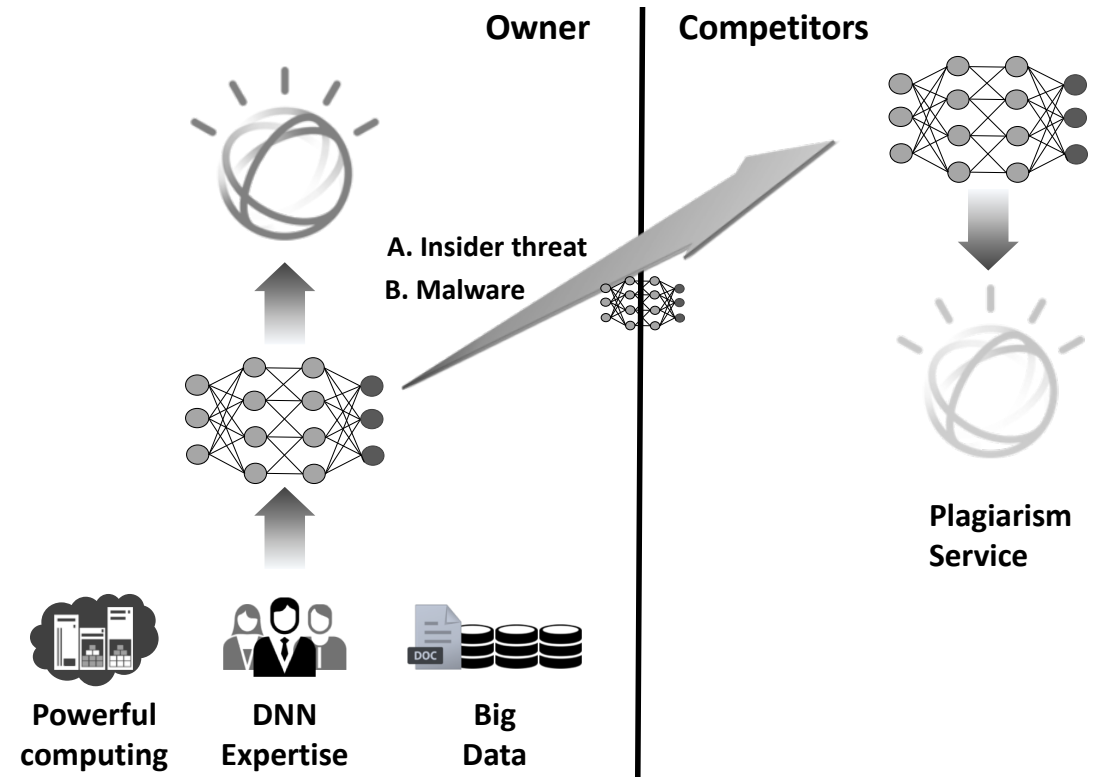    Ram Shankar Siva Kumar, David R. O'Brien, Kendra Albert, Salome Vilojen
    https://arxiv.org/abs/1810.10731
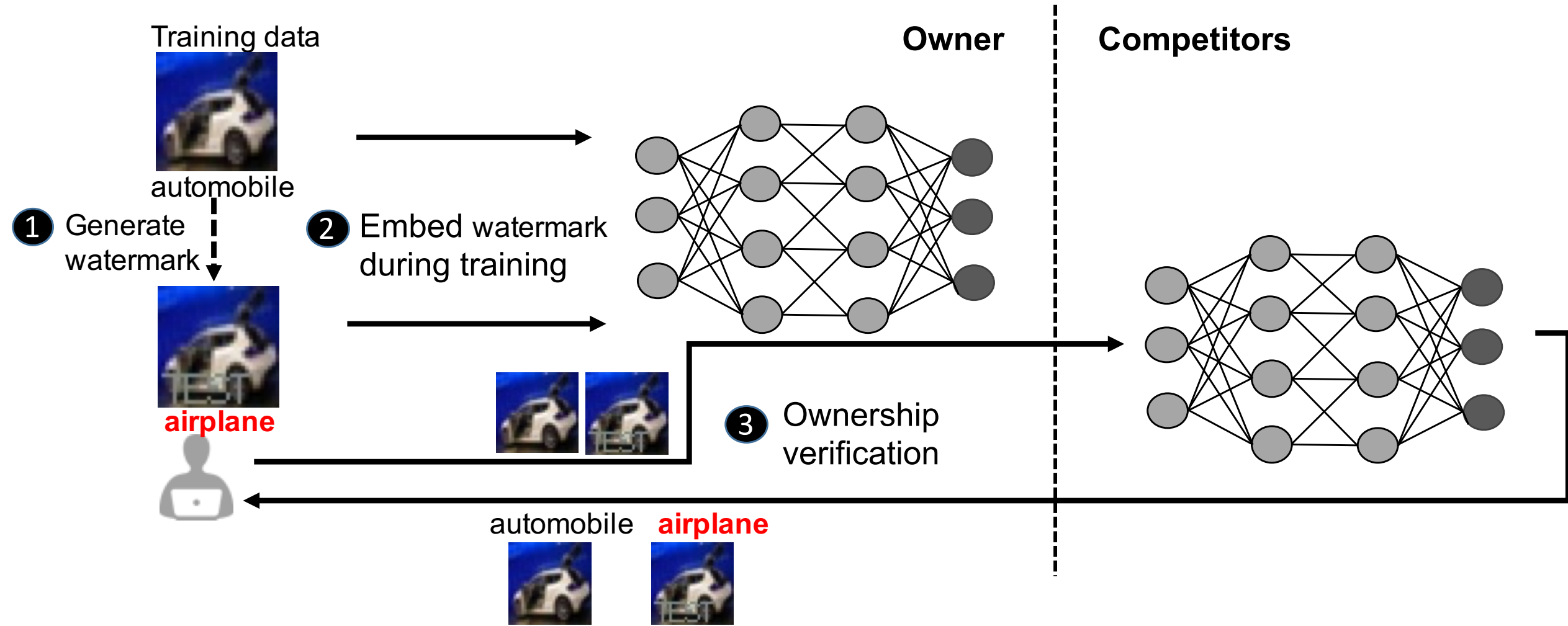
# Watermarking

**Watermarks Embedding**
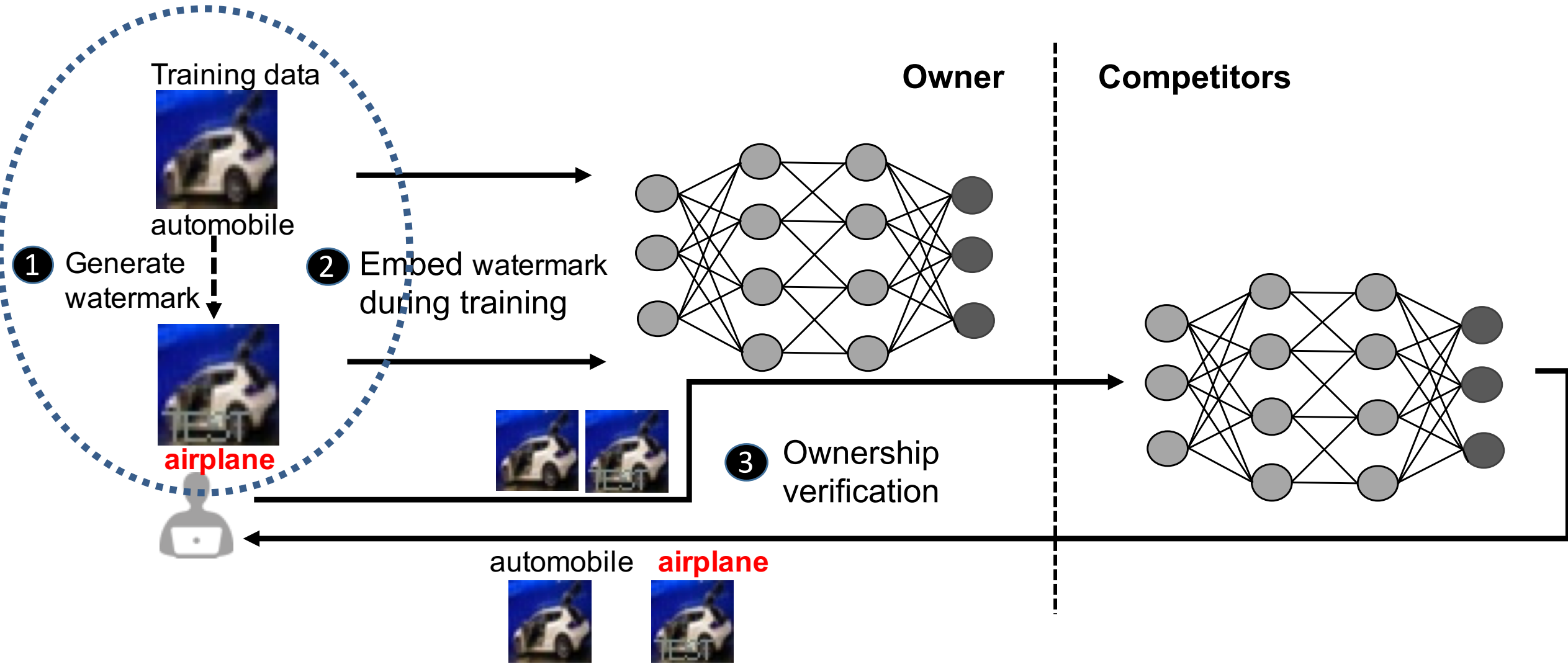
**Watermarks Verification**

# Idea

- Watermark in Deep Learning

- Allow for verifying the ownership

- Special training that delivers characteristic output for special examples

- Needs to be robust / resilient to
  - Counter watermarking
  - Fine-tuning
  - Training
  - Model inversions

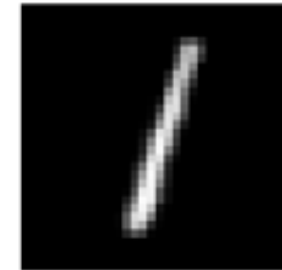# DNN Watermarking

# DNN Watermarking

# DNN Watermark generation



**airplane**

- Meaningful content embedded in original training data



**airplane**

- Independent training data with unrelated classes as watermarks



**airplane**

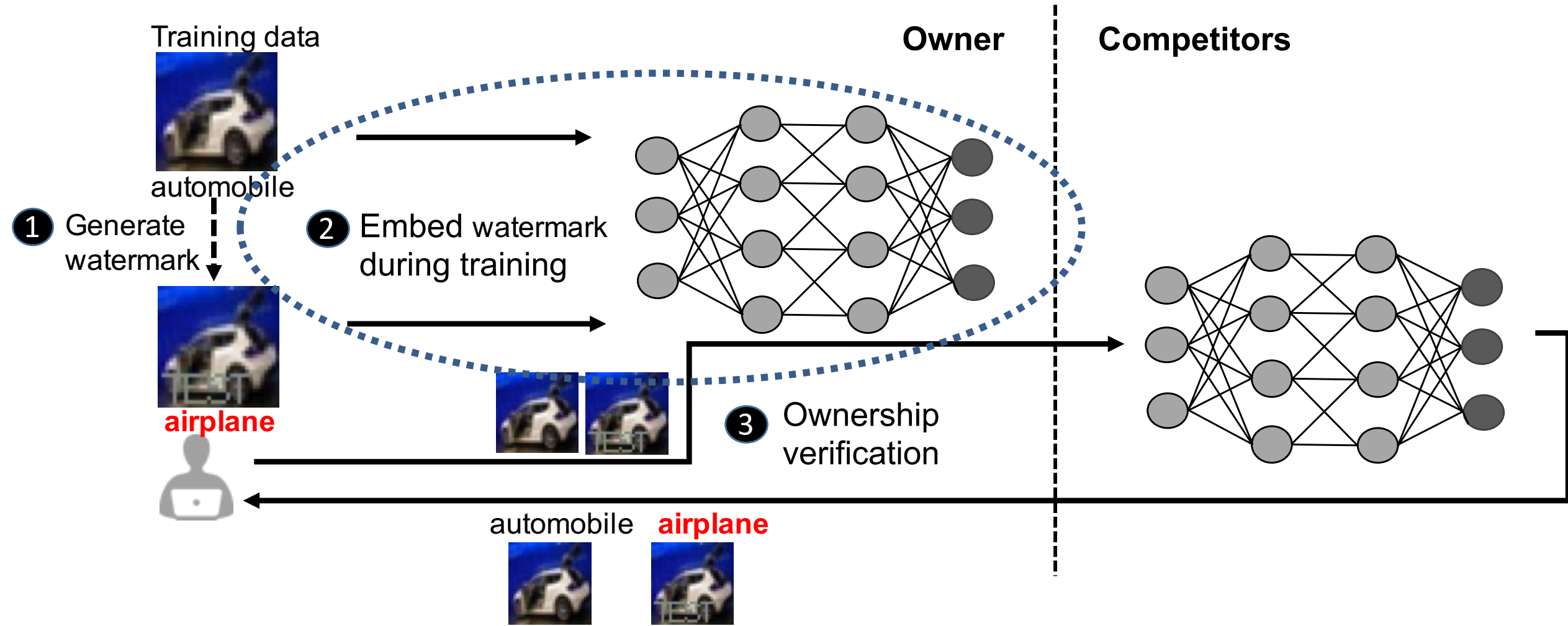- Pre-specified Noise as watermark

# DNN Watermarking



Training data

automobile

**Owner** | **Competitors**

**1** Generate watermark

**airplane**

**2** Embed watermark during training

**3** Ownership verification

automobile  **airplane**

**Algorithm 1** Watermark embedding

**Input:**

Training set $D_{train} = \{X_i, Y_i\}_{i=1}^{S}$

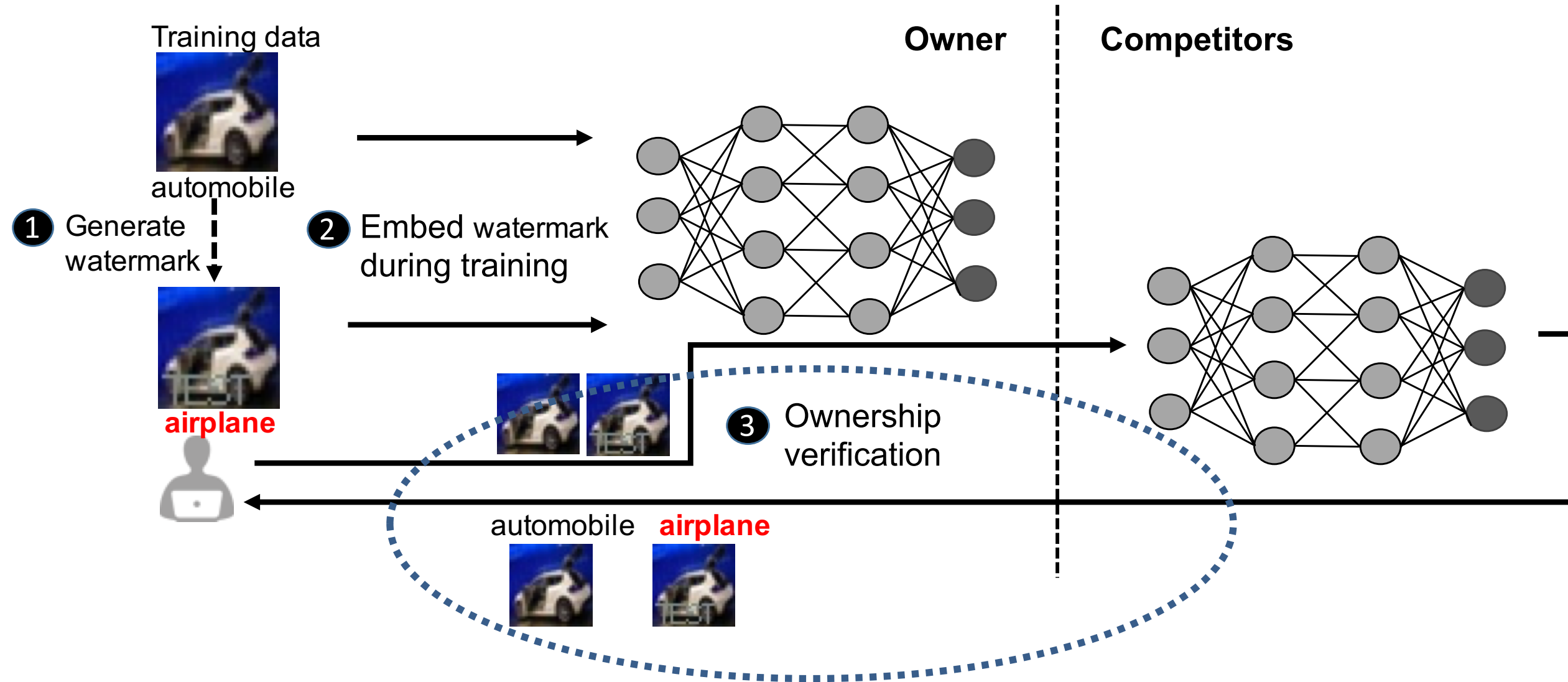DNN key K=$\{Y_s, Y_d\}(s \neq d)$

**Output:**

DNN model: $F_\theta$

Watermark Pair: $D_{wm}$

1: **function** WATERMARK_EMBEDDING()
2:     $D_{wm} \leftarrow \emptyset$
3:     $D_{tmp} \leftarrow sample(D_{train}, Y_s, percentage)$
4:     **for each** $d \in D_{tmp}$ **do**
5:         $x_{wm} = ADD\_WATERMARK(d[x], watermarks)$
6:         $y_{wm} = y_d$
7:         $D_{wm} = D_{wm} \cup \{x_{wm}, y_{wm}\}$
8:     **end for**
9: **end function**
10: $F_\theta = Train(D_{wm}, D_{train})$
11: **return** $F_\theta, D_{wm}$

# DNN Watermarking

CISPA
HELMHOLTZ-ZENTRUM i. G.

Training data

automobile

**Owner**          **Competitors**

① Generate watermark

② Embed watermark during training

airplane

③ Ownership verification

automobile   **airplane**

# Ownership Verification

- Adversary might want to monetize model with online API

- Query with watermarked images

- If it flips label as trained -> our model

- Works on trained images (basically overfitting on training set)

- Even works on newly watermarked images (generalization of watermarks to test)
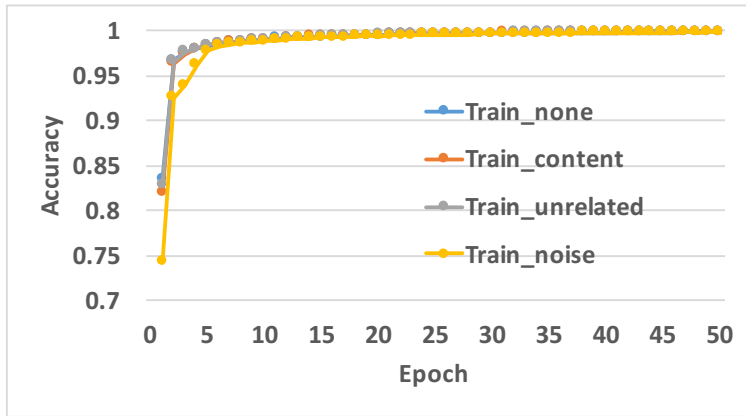
**(a) MNIST**

| Accuracy | $WM_{content}$ | $WM_{unrelated}$ | $WM_{noise}$ |
|---|---|---|---|
| Watermarks (trained) | 100% | 100% | 100% |
| Watermarks (new) | 100% | 100% | 99.42% |

**(b) CIFAR10**

| Accuracy | $WM_{content}$ | $WM_{unrelated}$ | $WM_{noise}$ |
|---|---|---|---|
| Watermarks (trained) | 99.93% | 100% | 99.86% |
| Watermarks (new) | 98.6% | 100% | 94.1% |

# Side Effects

- Does including watermarked images effect train/val/test accuracies?



(a) Train accuracy



(b) Validation accuracy

**Figure 6: Model accuracy over training procedure (MNIST)**

(a) MNIST

| CleanModel | $WM_{content}$ | $WM_{unrelated}$ | $WM_{noise}$ |
|------------|----------------|------------------|--------------|
| 99.28 % | 99.46% | 99.43% | 99.41% |

(b) CIFQR10

| CleanModel | $WM_{content}$ | $WM_{unrelated}$ | $WM_{noise}$ |
|------------|----------------|------------------|--------------|
| 78.6% | 78.41% | 78.12% | 78.49% |



(a) Train accuracy



(b) Validation accuracy

# Robustness

- Does the model retrain the watermarking – despite modification to model

- Pruning:
  - Remove small weights in model

- Fine-Tuning:
  - Continue training with more examples

- High robustness

**Table 3: Robustness for model pruning: accuracy of clean testing data and accuracy of watermarks (MNIST)**

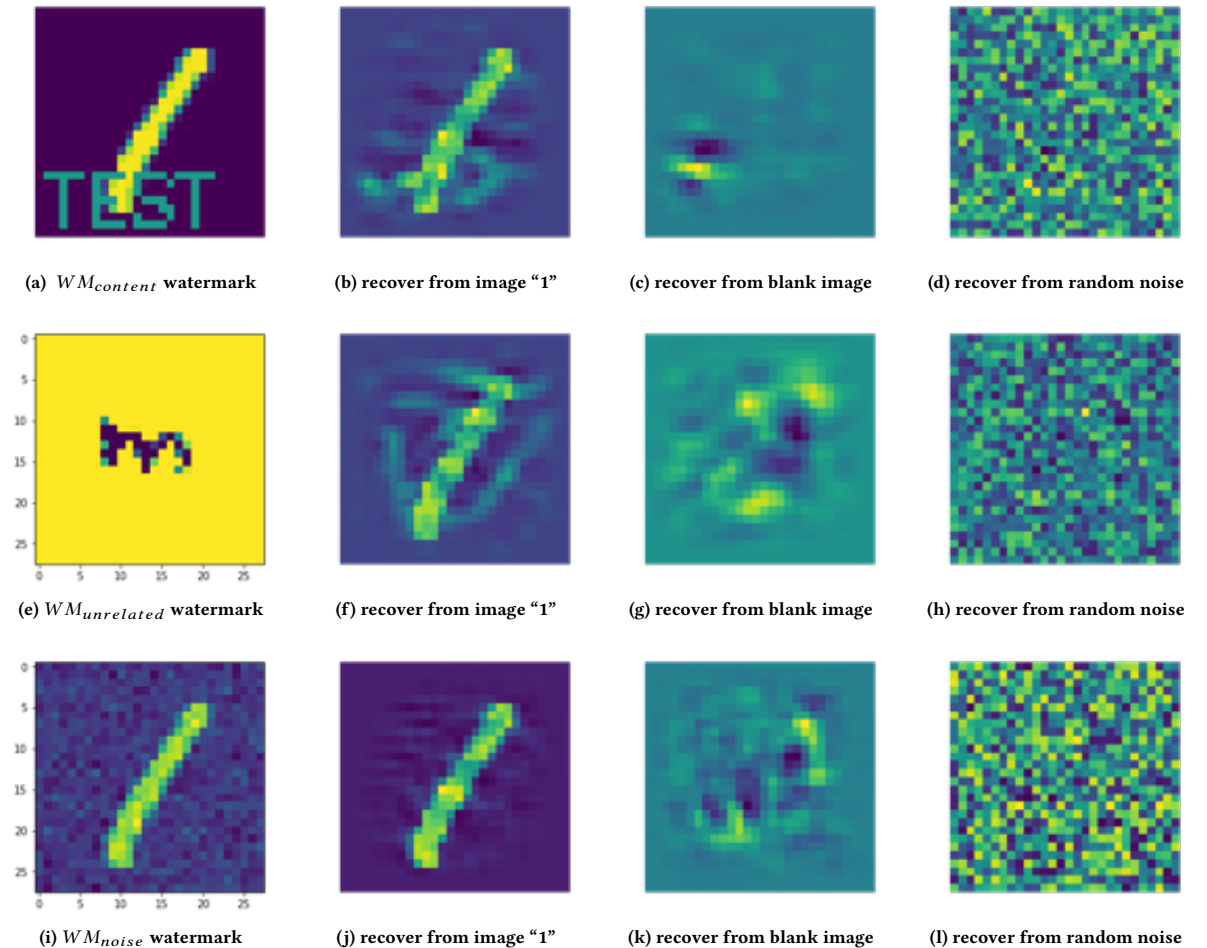| Pruning rate | $WM_{content}$ | | $WM_{unrelated}$ | | $WM_{noise}$ | |
|---|---|---|---|---|---|---|
| | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. |
| 10% | 99.44% | 100% | 99.43% | 100% | 99.4% | 100% |
| 20% | 99.45% | 100% | 99.45% | 100% | 99.41% | 100% |
| 30% | 99.43% | 100% | 99.41% | 100% | 99.41% | 100% |
| 40% | 99.4% | 100% | 99.31% | 100% | 99.42% | 100% |
| 50% | 99.29% | 100% | 99.19% | 100% | 99.41% | 100% |
| 60% | 99.27% | 100% | 99.24% | 100% | 99.3% | 99.9% |
| 70% | 99.18% | 100% | 98.82% | 100% | 99.22% | 99.9% |
| 80% | 98.92% | 100% | 97.79% | 100% | 99.04% | 99.9% |
| 90% | 97.03% | 99.95% | 93.55% | 99.9% | 95.19% | 99.55% |

**Table 4: Robustness for model pruning: accuracy of clean testing data and accuracy of watermarks (CIFAR10)**

| Pruning rate | $WM_{content}$ | | $WM_{unrelated}$ | | $WM_{noise}$ | |
|---|---|---|---|---|---|---|
| | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. |
| 10% | 78.37% | 99.93% | 78.06% | 100% | 78.45% | 99.86% |
| 20% | 78.42% | 99.93% | 78.08% | 100% | 78.5% | 99.86% |
| 30% | 78.2% | 99.93% | 78.05% | 100% | 78.33% | 99.93% |
| 40% | 78.24% | 99.93% | 77.78% | 100% | 78.31% | 99.93% |
| 50% | 78.16% | 99.93% | 77.75% | 100% | 78.02% | 99.8% |
| 60% | 77.87% | 99.86% | 77.44% | 100% | 77.87% | 99.6% |
| 70% | 76.7% | 99.86% | 76.71% | 100% | 77.01% | 98.46% |
| 80% | 74.59% | 99.8% | 74.57% | 96.39% | 73.09% | 92.8% |
| 90% | 64.9% | 99.47% | 62.15% | 10.93% | 59.29% | 65.13% |

**Table 5: Robustness for model fine-tuning: accuracy of clean testing data and accuracy of watermarks**

| Dataset | $WM_{content}$ | | $WM_{unrelated}$ | | $WM_{noise}$ | |
|---|---|---|---|---|---|---|
| | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. | Testing Acc. | Watermark Acc. |
| MNIST | 99.6% | 99.95% | 99.64% | 100% | 99.68% | 99.85% |
| CIFAR10 | 77.55% | 98.33% | 76.75% | 95.33% | 78.43% | 69.13% |

- Can watermark be recovered from classifier?

- Attack using gradient based technique: Fredrikson, Matt, Somesh Jha, and Thomas Ristenpart. "Model inversion attacks that exploit confidence information and basic countermeasures." In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 1322-1333. ACM, 2015.
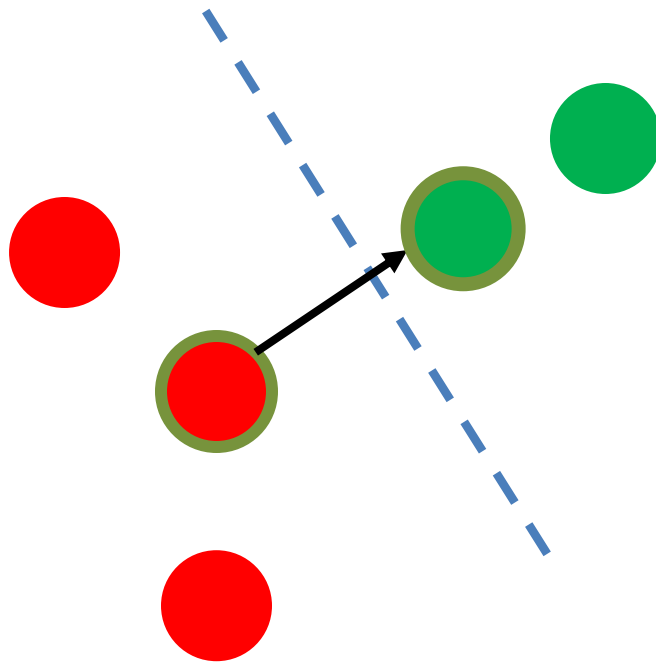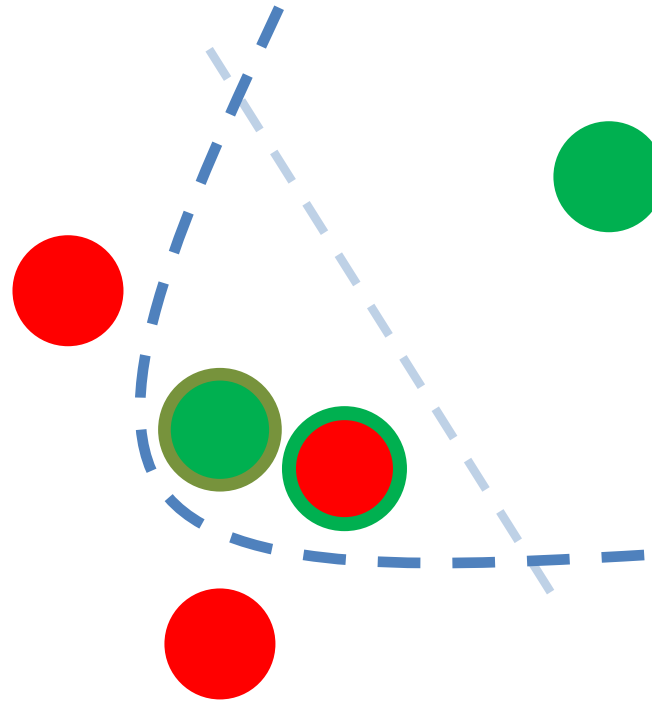
- Does not see effective



(a) $WM_{content}$ watermark    (b) recover from image "1"    (c) recover from blank image    (d) recover from random noise

(e) $WM_{unrelated}$ watermark    (f) recover from image "1"    (g) recover from blank image    (h) recover from random noise

(i) $WM_{noise}$ watermark    (j) recover from image "1"    (k) recover from blank image    (l) recover from random noise

# Poisoning

Mario Fritz | 19.12.2018

# Poisoning vs Evasion Attacks

**CISPA**
HELMHOLTZ-ZENTRUM i. G.

**Evasion Attack
(Adversarial Perturbation)**

**Poisoning Attack**

**Clean Data
Poisoning Attack**

Manipulation of test data

Inject data and label into
training set; often wrong
label

Inject data training set;
labeling is correct – can also
be done by the victim

# Attack Technique: Model Poisoning

- Online systems sacrifice stationarity for adaptability
  - System is re-train/adapted during deployment

- Dependent on how much control users have on the training input

- Sometimes easy to detect rubbish
- Boiling frog attacks: gradually inject poisoning data in order to make it harder to detect

- What is distribution drift that we want to adapt to?

- What is adversarial data poisoning that we want to robust to?



THE VERGE   TECH  SCIENCE  CULTURE  CARS  REVIEWS  LONGFORM  VIDEO  MORE

MICROSOFT \ WEB \ TL;DR

## Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

By James Vincent | @jjvincent | Mar 24, 2016, 6:43am EDT

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

♡ 10.7K   6:56 AM - Mar 24, 2016

💬 12.4K people are talking about this

# Poisoning

Before attack — After attack — After sanitization

[Koh'18]

- ML models are often trained on data from the "outside"

- Not in our control – or we depend on it because of scale or real-world scenario

- Adversary can inject data points in our training dataset

- Common defense: data sanatization

- Automated defense

  - Too much data to do human inspection

  - Also human is not a good baseline anyways

- Attacker evaluation

  - Attacker wants to increase error no matter what defenses are deployed

- Attack budget and defense thresholds

  - Attacker has limited control of the dataset

  - Typical assumptions 3-5%

- Binary Classification

$$f_\theta : \mathcal{X} \to \{-1, +1\} \quad x \in \mathcal{X} \quad y \in \{-1, +1\} \quad f_\theta(x) = \text{sign}(\theta^\top x)$$

- Misclassification

$$L_{0\text{-}1}(\theta; \mathcal{D}_{\text{test}}) = \frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{(x,y) \in \mathcal{D}_{\text{test}}} \mathbf{I}[f_\theta(x) \neq y]$$

- Defender wants to estimate theta^ to **minimize** the error

- Attacker want to mislead Defender to **maximize** error $\quad L_{0\text{-}1}(\hat{\theta}; \mathcal{D}_{\text{test}})$

- Attacker picks $\epsilon n$ poisoned points $\mathcal{D}_{\text{p}}$

- Trainset $\mathcal{D} = \mathcal{D}_{\text{c}} \cup \mathcal{D}_{\text{p}}$

- Adversarial ML deja-vu: Min-Max objective

**Attacker:**

- Input: Clean training data $\mathcal{D}_{\mathrm{c}}$ and test data $\mathcal{D}_{\mathrm{test}}$.

- Output: Poisoned training data $\mathcal{D}_{\mathrm{p}}$, with $|\mathcal{D}_{\mathrm{p}}| = \epsilon |\mathcal{D}_{\mathrm{c}}|$.

- Goal: Mislead defender into learning parameters $\hat{\theta}$ with high test error $L_{\text{0-1}}(\hat{\theta}; \mathcal{D}_{\mathrm{test}})$.

**Defender:**

- Input: Combined training data $\mathcal{D} = \mathcal{D}_{\mathrm{c}} \cup \mathcal{D}_{\mathrm{p}}$.

- Output: Model parameters $\hat{\theta}$.

- Goal: Learn model parameters $\hat{\theta}$ with low test error $L_{\text{0-1}}(\hat{\theta}; \mathcal{D}_{\mathrm{test}})$ by filtering out poisoned points $\mathcal{D}_{\mathrm{p}}$.

[Koh'18]

- Defender tries to remove suspicious points from $\mathcal{D} = \mathcal{D}_c \cup \mathcal{D}_p$

- Train on remaining data

- Idea: poisoned data that is similar to clean does not matter much

- E.g. **L2** defense:

  - Find class centroids

  - Throw away data that is far away from centroids

- More formally:

  - Rate "anomaly of each data point": $score\ function\ s_\beta : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$

  - Parameterized by $anomaly\ detector\ parameters\ \beta$

- E.g. in **L2** defense

  - Parameters are the centroids $\beta = (\mu_+, \mu_-)$

  - Scoring function $s_\beta(x, y) = \|x - \beta_y\|_2$

- **Defense**:

  - Fit $anomaly\ detector\ parameters\ \beta = B(\mathcal{D})$

  - Construct $feasible\ set\ \mathcal{F}_\beta = \{(x, y) : (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ with } s_\beta(x, y) < \tau_y\}$ with $\text{threshold } \tau_y$

  - Sanatized training data $\mathcal{D}_{\text{san}} = \mathcal{D} \cap \mathcal{F}_\beta$

  - Training: Minimize $\hat{\theta}$ over loss: $\hat{\theta} = \underset{\theta}{\arg\min}\, L(\theta; \mathcal{D}_{\text{san}}) \overset{\text{def}}{=} \underset{\theta}{\arg\min}\, \frac{\lambda}{2}\|\theta\|_2^2 + \frac{1}{|\mathcal{D}_{\text{san}}|} \sum_{(x,y) \in \mathcal{D}_{\text{san}}} \ell(\theta; x, y)$

- **L2** defense rejects points far from the class centroids

$$\beta_y = \mathbb{E}_{\mathcal{D}}[x|y]$$
$$s_\beta(x, y) = \|x - \beta_y\|_2$$

- **Slab** defense [Steinhardt'17]

  – Project on line between centroids

  – Reject points according to distance

$$\beta_y = \mathbb{E}_{\mathcal{D}}[x|y]$$
$$s_\beta(x, y) = \left|(\beta_1 - \beta_{-1})^\top (x - \beta_y)\right|$$

  – Idea: focus on more relevant dimension – not all of them as in **L2**

- **Loss** defense

  - Estimate model parameters on $\mathcal{D}_{\mathrm{c}} \cup \mathcal{D}_{\mathrm{p}}$

  - Score points by loss / fit

  $$\beta = \underset{\theta}{\operatorname{argmin}} \, \mathbb{E}_{\mathcal{D}}[\ell_\theta(x, y)]$$

  $$s_\beta(x, y) = \ell_\beta(x, y)$$

  - Somewhat similar to slab

  - Anomaly w.r.t. parametric model – focuses on transformed feature space

- **SVD** defense [Rubinstein'09]

  – Assume that clea data lies in some low-rank subspace

  – Poisoned data has high residual

  – Given the data matrix X:

  $$\beta = \text{Matrix of top } k \text{ right singular vectors of } X$$

  $$s_\beta(x, y) = \|(I - \beta\beta^\top)x\|_2$$

  – Hyperparameter k : typically picked based on the eigenvalue spectrum; e.g. sum of squares of larges eigenvalues -> e.g. reconstruct 95% of data

# Defense Strategies

- **K-NN**

  - Remove points that are far away from k nearest neighbor

  $$\beta = \mathcal{D}_\mathrm{c} \cup \mathcal{D}_\mathrm{p}$$
  $$s_\beta(x, y) = \text{Distance to } k\text{-th nearest neighbor in } \beta$$

  - E.g. k = 5

# Clean Label Poisoning

# Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks

- Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, Tom Goldstein

- NeurIPS 2018

[Slides based lecture by Bo Li]

- Data Poisoning Attacks

  – Happens at training time

  – Manipulate performance of system through constructed poison instances

- **Generally requires** some degree of control over labeling function for data

- Indiscriminate attack

  – Degrade test accuracy

- **Targeted Attack**

  – Aim to control behavior on specific test instance(s)

# Clean Label Attack

- Choose a target instance from the test set

- Sample a base instance from the base class and construct a poison

- Poison is injected into the training data

- Poison is cleanly labeled by labeling party

- Model is retrained on poisoned dataset

- Success if *target* is classified as being in the base class
  - Example: malware as benign software


- Deployment:
  - Place poisoned images on web
  - Wait for being crawled
  - A bit like fake news ☹

# Crafting Poison Data via Feature Collisions



benign     spam

$$\mathbf{p} = \text{argmin}_{\mathbf{x}} \, \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2 + \beta \, \|\mathbf{x} - \mathbf{b}\|_2^2$$

- First term – gets the poison instance (**p**) to move toward the target instance in feature space

- Second term – tries to make **p** to appear like a base class to a human

- Training on data + poison can cause the decision boundary to rotate to include the target + poison in the base class

- This allows for a "backdoor" into the base class

**Algorithm 1** Poisoning Example Generation

**Input:** target instance $t$, base instance $b$
Initialize x: $x_0 \leftarrow b$
Define: $L_p(x) = \|f(\mathbf{x}) - f(\mathbf{t})\|^2$
**for** $i = 1$ **to** $maxIters$ **do**
    Forward step: $\widehat{x}_i = x_{i-1} - \lambda \nabla_x L_p(x_{i-1})$
    Backward step: $x_i = (\widehat{x}_i + \lambda\beta b)/(1 + \beta\lambda)$
**end for**

- Described by Goldstein et al. in 2014
- Forward step is gradient descent update to minimize distance from poison to the target instance in feature space
- Backward step is proximal update that minimizes the distance from the poison to base instance in input space
- Beta is tuned to make poison instance look realistic

Feature space representation

# "One-shot kill attack"

- Transfer learning scenario

- Pretrained CNN  is used as feature extraction network

- All weights are frozen, but the last layer (SoftMax) is retrained to adapt the network to a specific task

- Add one poison instance to cause misclassification of the target

- Showed 100% success rate across 1099 trials

  - High success rate due to more weights (2048) than examples (1801) causing overfitting on training data

- Original accuracy on test set is hardly affected

  - 0.2% average drop in accuracy

- Poisoning attack with correctly labeled training data

- Poisons aim to collide with target in feature space causing the network to incorrectly separate them

- Similar to adversarial training

- Does not degrade the performance for non-targeted examples

- Creates a method for creating backdoor in neural net

- More complicated and not as effective if whole architecture is trained ( and not only fine tuned)

# Membership Inference Attacks

Shokri, Reza, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. "Membership inference attacks against machine learning models." In *Security and Privacy (SP), 2017 IEEE Symposium on*, pp. 3-18. IEEE, 2017.

Salem, Ahmed, Yang Zhang, Mathias Humbert, Mario Fritz, and Michael Backes. "ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models." *NDSS 2019*

Privacy:
Does an ML model trained on privacy-sensitive data leak information of the data?

Forensics:
Can I tell which data was used?

Get some data

Train the model



ML model



100
80
60
40
20
0

cat    dog   panda

- Why membership matters?

  - A cliché example: a ML model for medical diagnosis, if a person is in the training set, then she has the corresponding disease

  - Security implications, IP implications

Attack by Shokri et al.

- **One** shadow model

- **One** attack model

- Same data distribution

- Can we do better?

- No assumption on the dataset

- Data transferring attack

- Train shadow model on a **different dataset**, and attack on the target model

Image dataset

Text dataset

Target Model

Shadow Model

Train

Test

Attack Model

In or not in

# Our Attack 2



Precision

Recall

CIFAR-100 M | News M
CIFAR-100 NoM | News NoM

- Can we do better?

- **No shadow model**

- Take the maximum, std, or entropy of the posterior as the score

  - The simplest attack
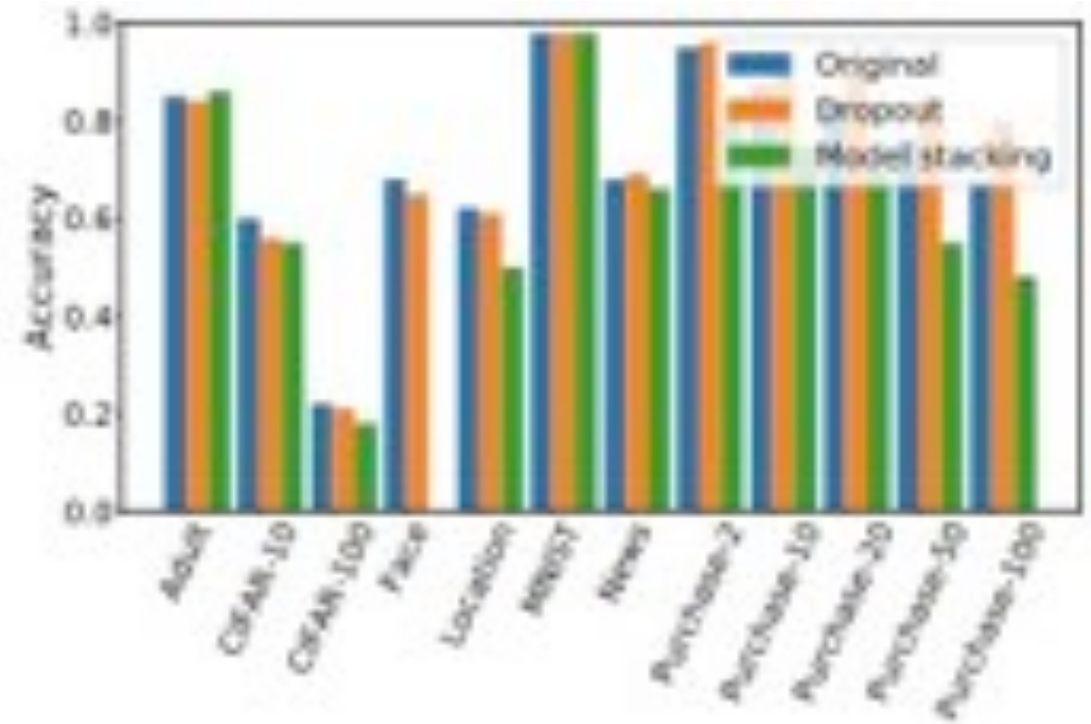
  - Unsupervised

  - Reason: overfitting

# How To Defend the Attack?
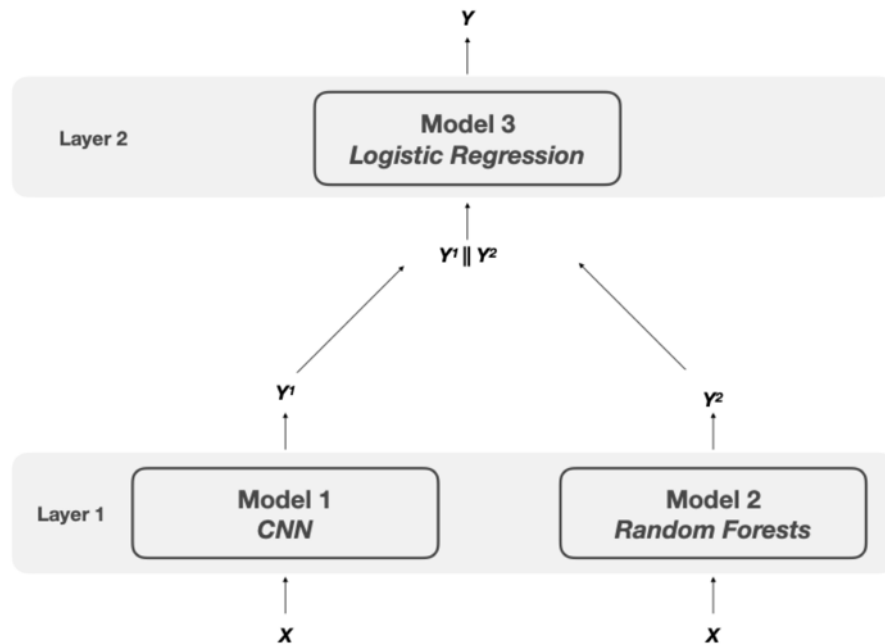
- Dropout

- Model Stacking

- Differential Private Training



(c) Accuracy.

# Differential Privacy

# Differential Privacy & Deep Learning

- Lecture based on Tutorial @ NIPS'17 :
  **Differentially Private Machine Learning: Theory, Algorithms, and Applications**
  **Kamalika Chaudhuri, Dept. of Computer Science and Engineering, UC San Diego**
  **Anand D. Sarwate, Dept. of Electrical and Computer Engineering, Rutgers University**
  https://www.ece.rutgers.edu/~asarwate/nips2017/

- **Deep Learning with Differential Privacy**
  Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, Li Zhang
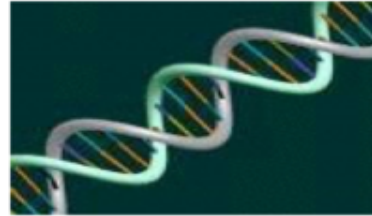  CCS 2016
  https://arxiv.org/abs/1607.00133

# Motivating
# Differential Privacy

# Sensitive Data

Medical Records

Genetic Data

Search Logs

# Differential privacy in practice

**Google:** RAPPOR for tracking statistics in Chrome.
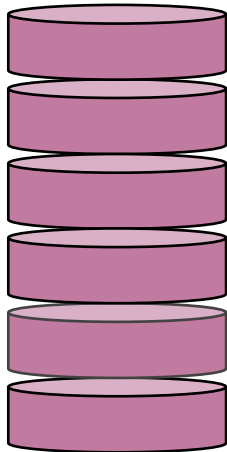
**Apple:** various iPhone usage statistics.

**Census:** 2020 US Census will use differential privacy.
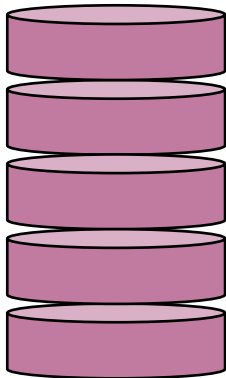
**mostly focused on count and average statistics**
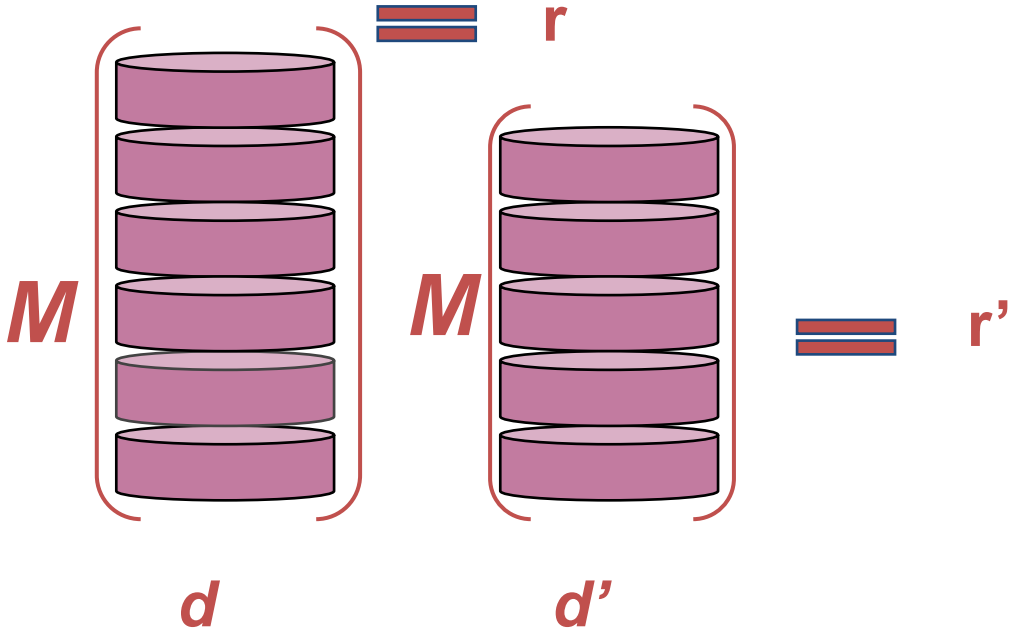
# Defense Proposal: Differential Privacy
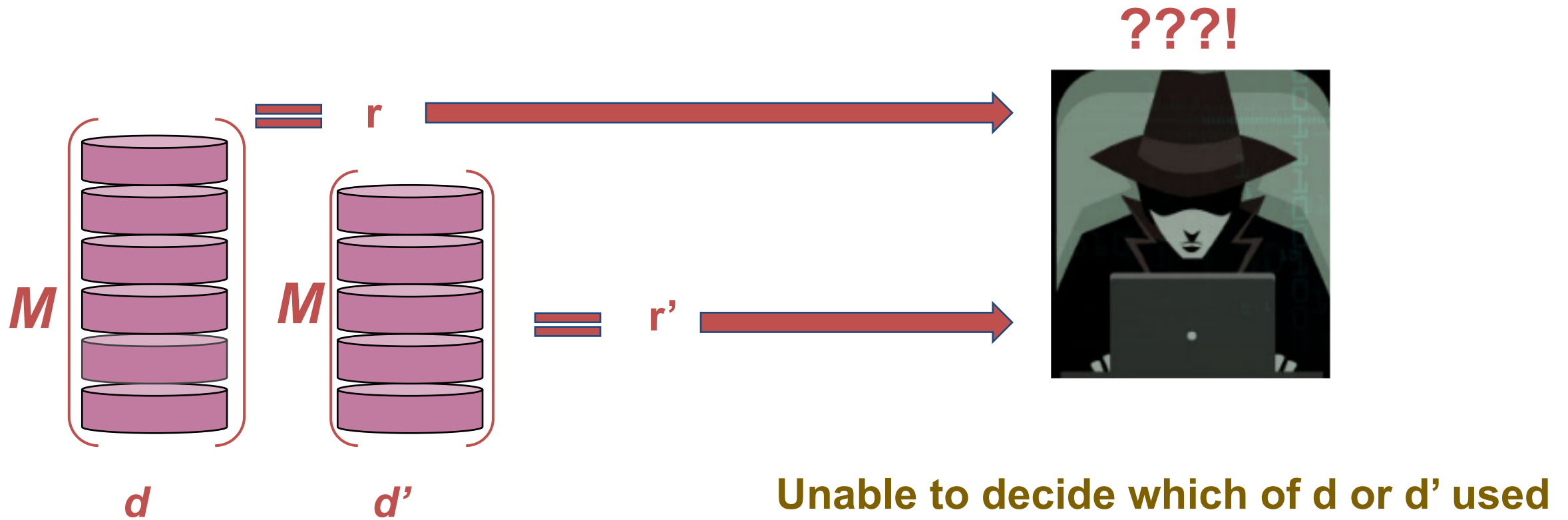
Adjacent datasets:



*d*          *d'*

# Defense Proposal: Differential Privacy



*M* differentially private

# Defense Proposal: Differential Privacy



**???!**

**Unable to decide which of d or d' used**

$M$    $M$

r

r'

d        d'

# Differential Privacy: Definition

Randomized algorithm $M : \mathcal{D} \to \mathcal{R}$ is $(\epsilon, \delta)$ - differentially private if for any two adjacent datasets and any subset of the outputs $S \subseteq \mathcal{R}$

$$\Pr[M(d) \in S] \leq e^{\epsilon} \Pr[M(d') \in S] + \delta$$

# Differential Privacy: Definition

Randomized algorithm $M : \mathcal{D} \to \mathcal{R}$ is $(\epsilon, \delta)$ - differentially private if for any two adjacent datasets and any subset of the outputs $S \subseteq \mathcal{R}$

$$\Pr[M(d) \in S] \leq e^{\epsilon} \Pr[M(d') \in S] + \delta$$

Privacy budget/cost: smaller $\longrightarrow$ more privacy

# Differential Privacy: Definition

Randomized algorithm $M : \mathcal{D} \rightarrow \mathcal{R}$ is $(\epsilon, \delta)$ - differentially private if for any two adjacent datasets and any subset of the outputs $S \subseteq \mathcal{R}$

Perturbation of M, $\delta \propto \frac{1}{\|\delta\|_1}$

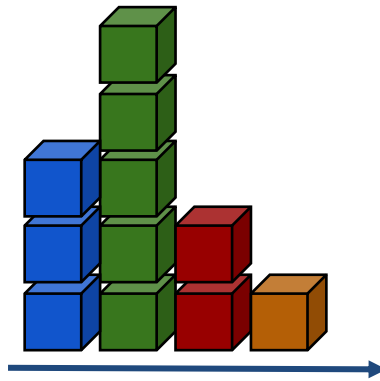$$\Pr[M(d) \in S] \leq e^{\epsilon} \Pr[M(d') \in S] + \delta$$

Privacy budget/cost: smaller $\longrightarrow$ more privacy
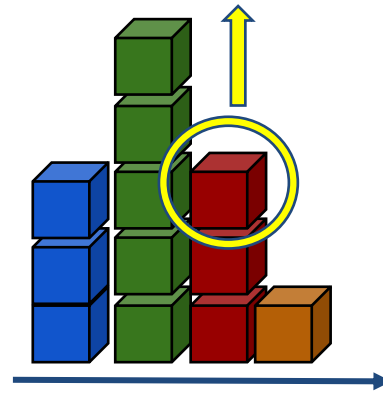
# How to Construct this DP Algo?

**Sensitivity of function:**

at most this much difference

**histogram**



*f(d)*

*f(d')*

# How to Construct this DP Algo?

**Sensitivity of function:**

at most this much difference

$\pi$  $2\pi$  1  -1

$$\sin(x_1 + x_2 + \ldots + x_n)$$     $$\sin(x_1 + x_2 + \ldots + x_n + x_{n+1})$$

*f(d)*     *f(d')*
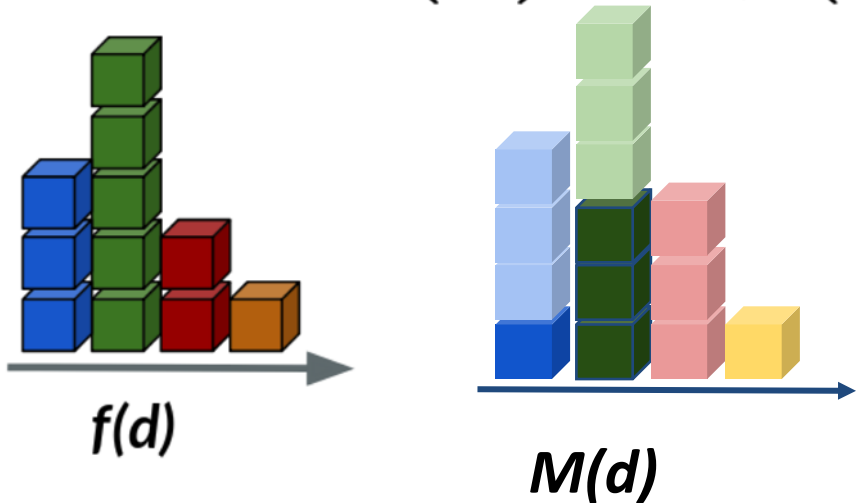
# How to Construct this DP Algo?

**Sensitivity of function:**

$$\Delta(f) = \max_{d,d' \text{ adjacent}} \| f(d) - f(d') \|$$

# How to Construct this DP Algo?



Gaussian Mechanism

$$M(d) = f(d) + \mathcal{N}(0, \sigma^2)$$

$$\sigma \geq \frac{1}{\epsilon}\sqrt{2ln(\frac{1.25}{\delta})}\Delta(f)$$

sensitivity

f(d)

M(d)

# How to Construct this DP Algo?

**Gaussian Mechanism**

$$M(d) = f(d) + \mathcal{N}(0, \sigma^2)$$

$$\sigma \geq \frac{1}{\epsilon}\sqrt{2ln(\frac{1.25}{\delta})}\Delta(f)$$

- Higher sensitivity → more noise needed
- More noise → smaller $\epsilon$

# How to Construct this DP Algo?

1) Determine sensitivity of the function
2) Add appropriate amount of noise

- If sensitivity is big add more noise
- More noise, better privacy

**But there's a catch:**
- **More noise destroys utility**

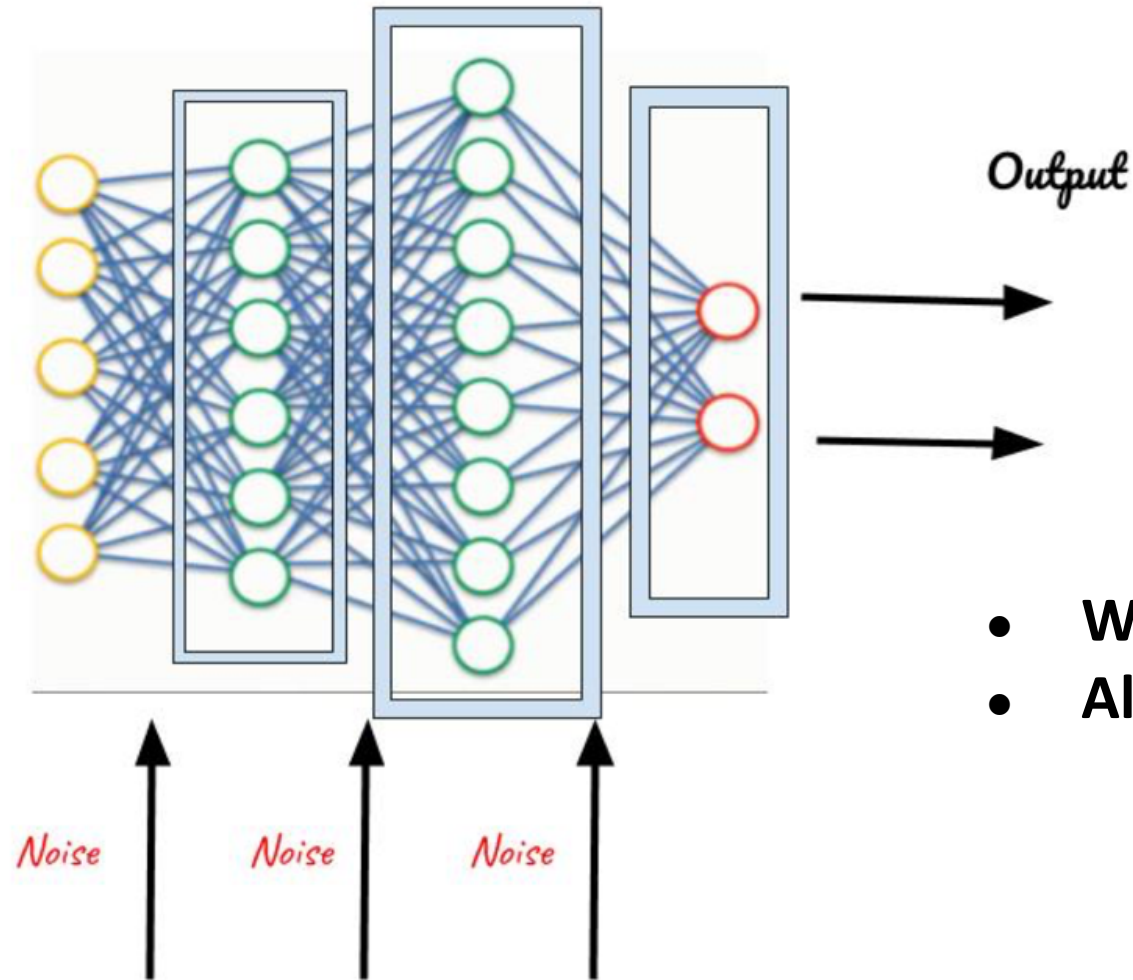**Differential Privacy + Machine Learning**

**Defense for MI attack**

# Differentially Private Stochastic Gradient Descent (DPSGD)

**M. Abadi et al. , "Deep Learning with Differential Privacy"**



- **White-box access**
- **All parameters protected**

# Differentially Private Stochastic Gradient Descent (DPSGD)

**M. Abadi et al. , "Deep Learning with Differential Privacy"**

---

**Algorithm 1** Differentially private SGD

**Input:** Examples $\{x_1, \ldots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate $\eta_t$, noise scale $\sigma$, group size $L$, gradient norm bound $C$.

    **Initialize** $\theta_0$ randomly

    **for** $t \in [T]$

        Take a random sample $L_t$ with sampling probability $L/N$

        **Compute gradient**

        For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

        **Clip gradient**

        $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|}{C})$

        **Add noise**

        $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L}(\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, (\sigma)^2))$

        **Descent**

        $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

**Output** $\theta_T$ and overall privacy cost $(\epsilon, \delta)$

# Differentially Private Stochastic Gradient Descent (DPSGD)

**M. Abadi et al. , "Deep Learning with Differential Privacy"**

---

**Algorithm 1** Differentially private SGD

**Input:** Examples $\{x_1, \ldots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate $\eta_t$, noise scale $\sigma$, group size $L$, gradient norm bound $C$.

    **Initialize** $\theta_0$ randomly

    **for** $t \in [T]$

        Take a random sample $L_t$ with sampling probability $L/N$

        **Compute gradient**

        For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

        **Clip gradient**

        $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|}{C})$       → **Bounds sensitivity**
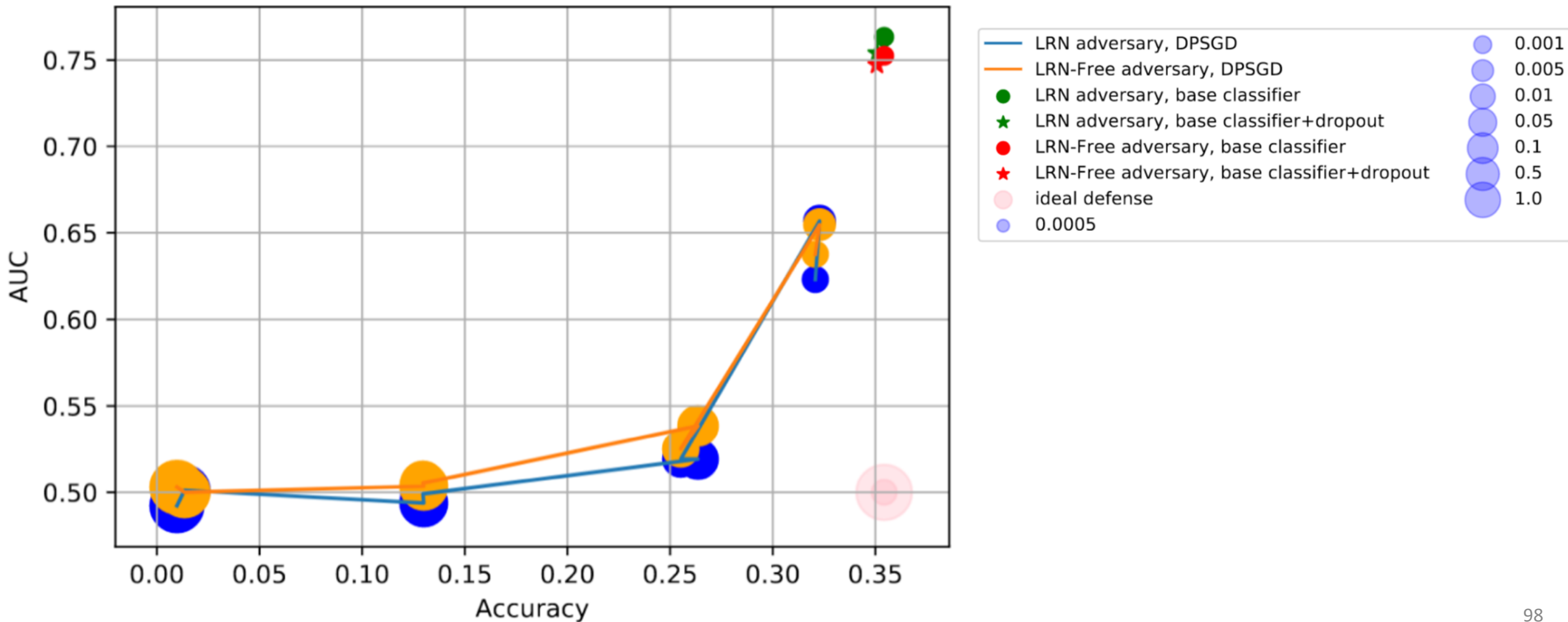
        **Add noise**

        $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L}(\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, (\sigma)^2))$

        **Descent**

        $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

**Output** $\theta_T$ and overall privacy cost $(\epsilon, \delta)$

---

# Differentially Private Stochastic Gradient Descent (DPSGD)

**M. Abadi et al. , "Deep Learning with Differential Privacy"**

---

**Algorithm 1** Differentially private SGD

**Input:** Examples $\{x_1, \ldots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate $\eta_t$, noise scale $\sigma$, group size $L$, gradient norm bound $C$.

$\quad$ **Initialize** $\theta_0$ randomly

$\quad$ **for** $t \in [T]$

$\qquad$ Take a random sample $L_t$ with sampling probability $L/N$

$\qquad$ **Compute gradient**

$\qquad$ For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

$\qquad$ **Clip gradient**

$\qquad$ $\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|}{C})$

$\qquad$ **Add noise**

$\qquad$ $\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L}(\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, (\sigma)^2))$
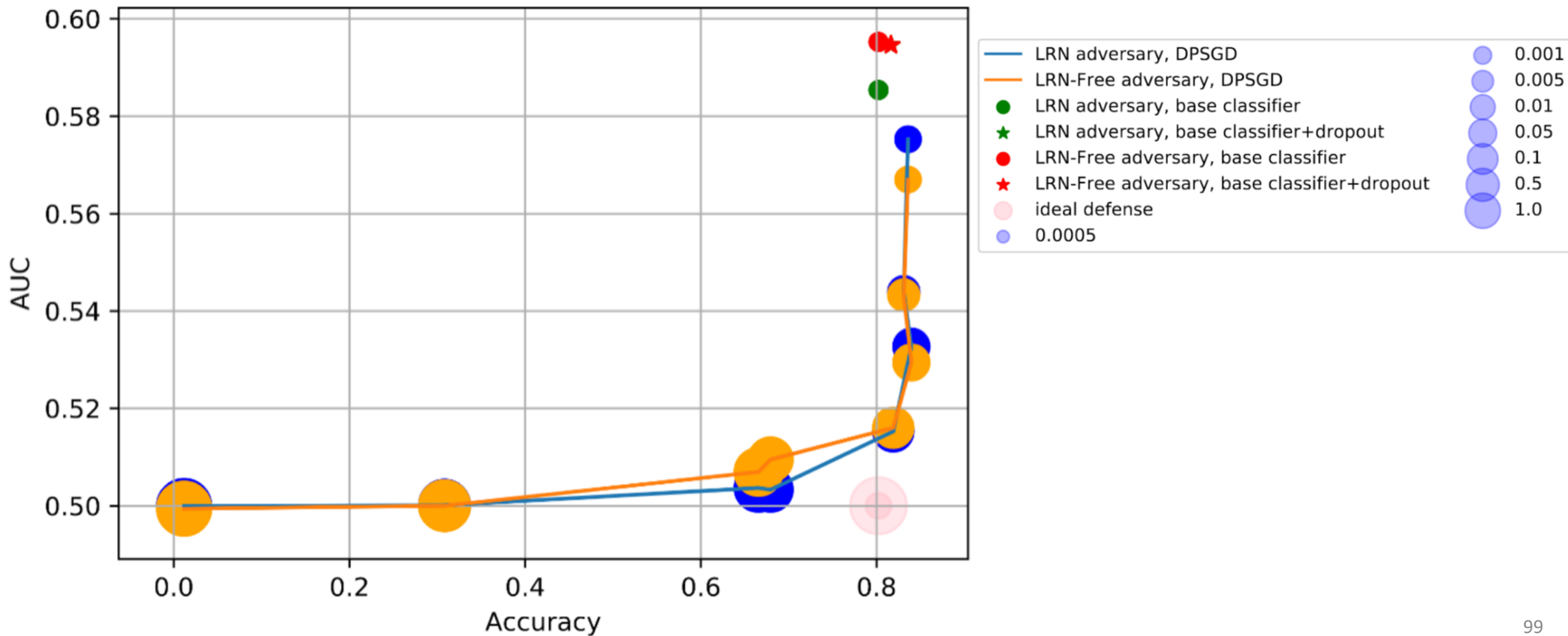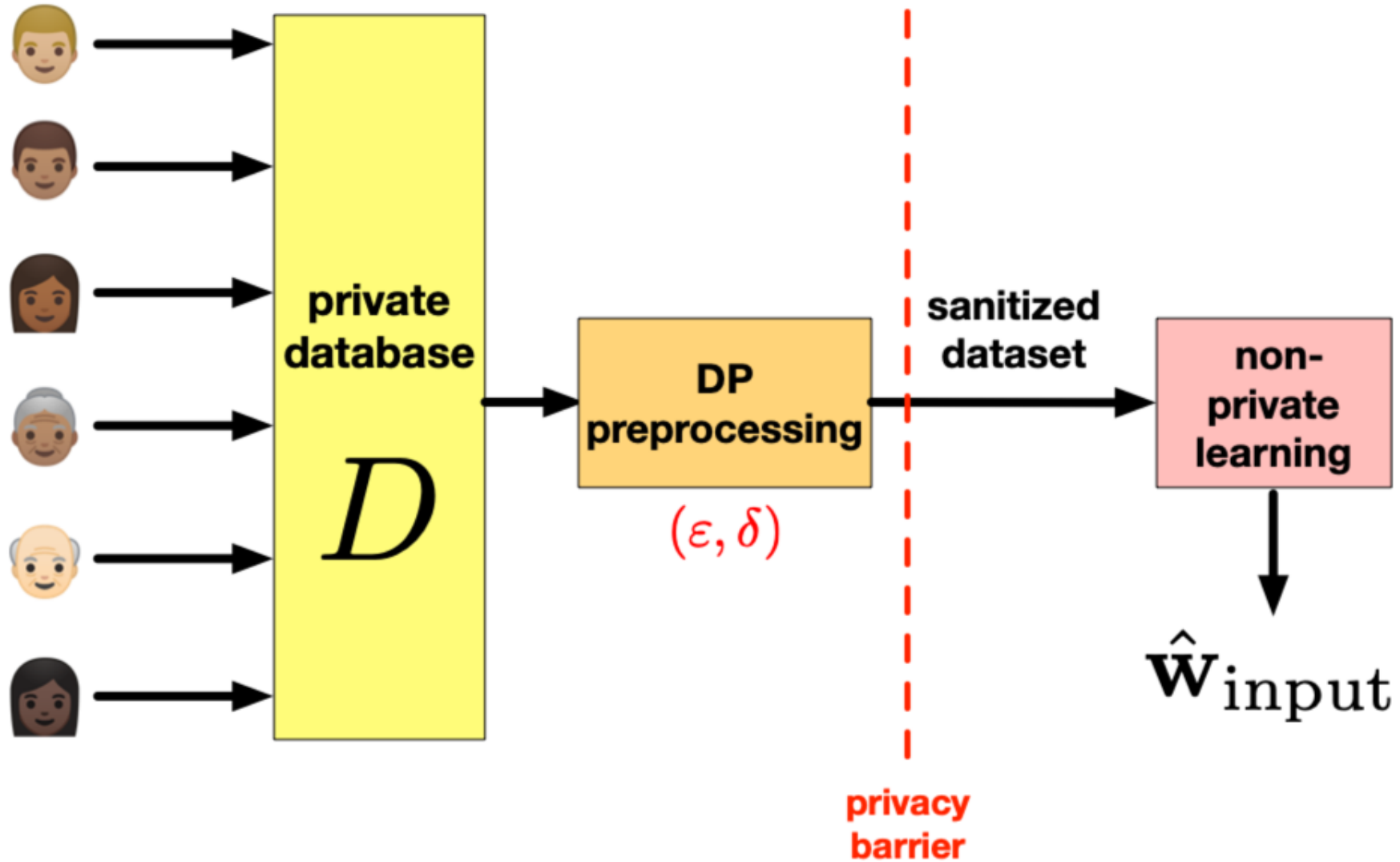
$\qquad$ **Descent**

$\qquad$ $\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

$\quad$ **Output** $\theta_T$ and overall privacy cost $(\epsilon, \delta)$

---

DPSGD, CIFAR100

DPSGD, Purchase100

Legend:
- LRN adversary, DPSGD
- LRN-Free adversary, DPSGD
- LRN adversary, base classifier
- LRN adversary, base classifier+dropout
- LRN-Free adversary, base classifier
- LRN-Free adversary, base classifier+dropout
- ideal defense
- 0.0005

Size scale: 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1.0
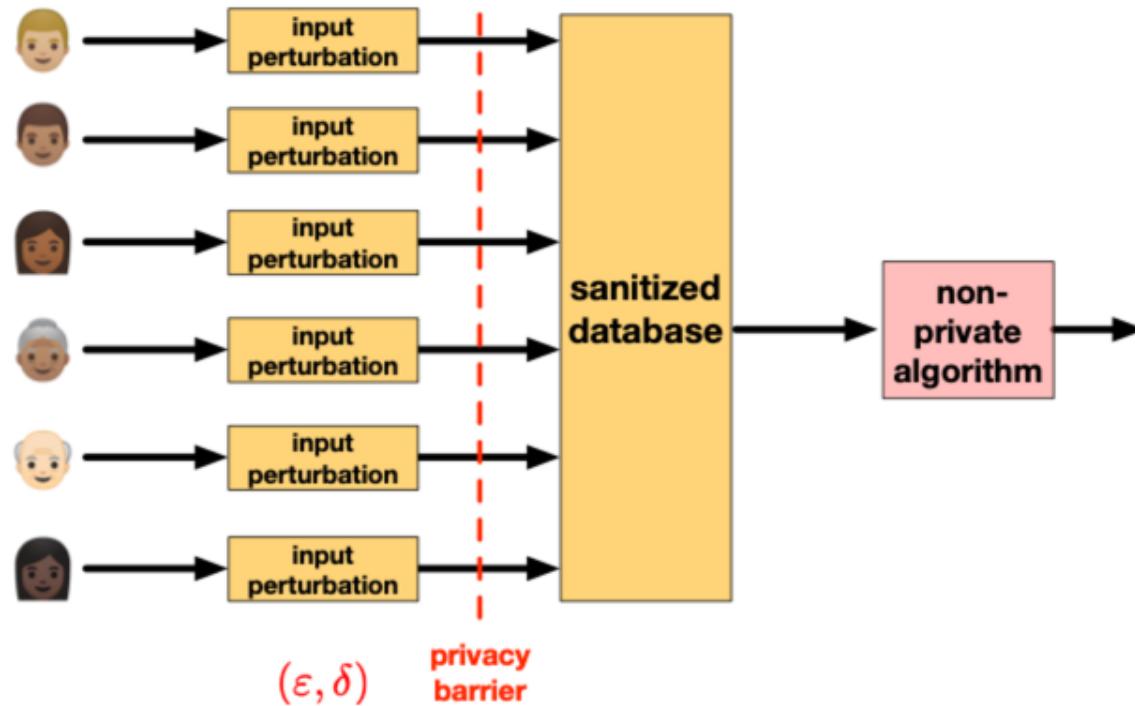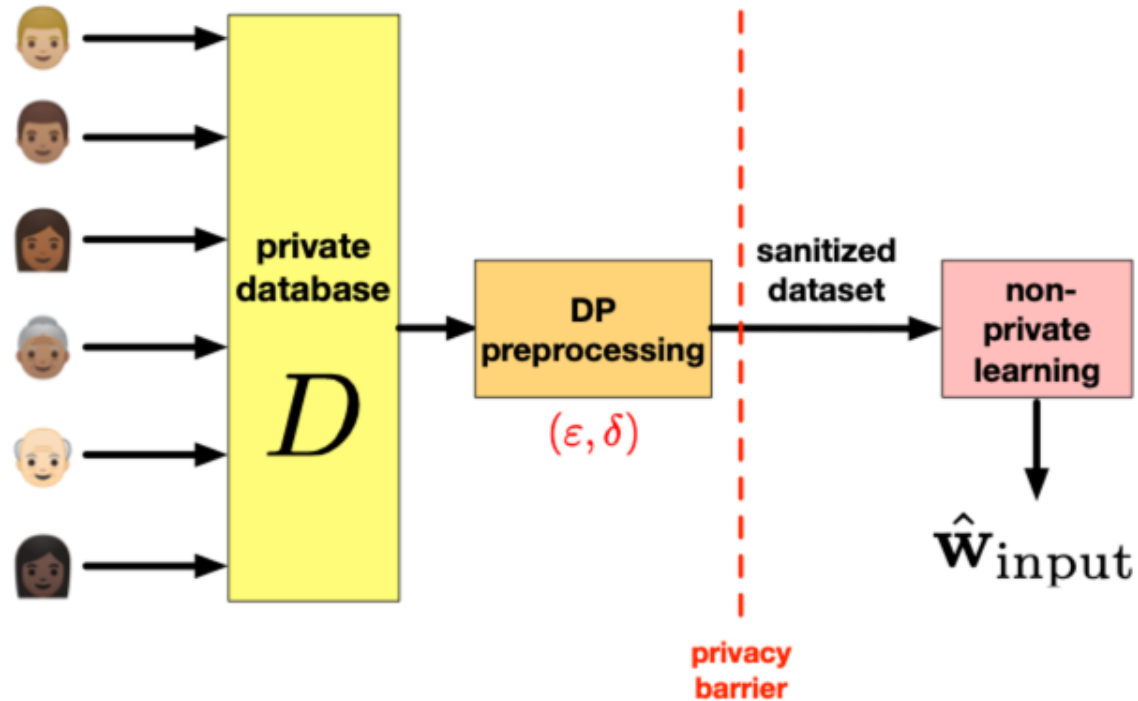
# Privacy in ERM: options
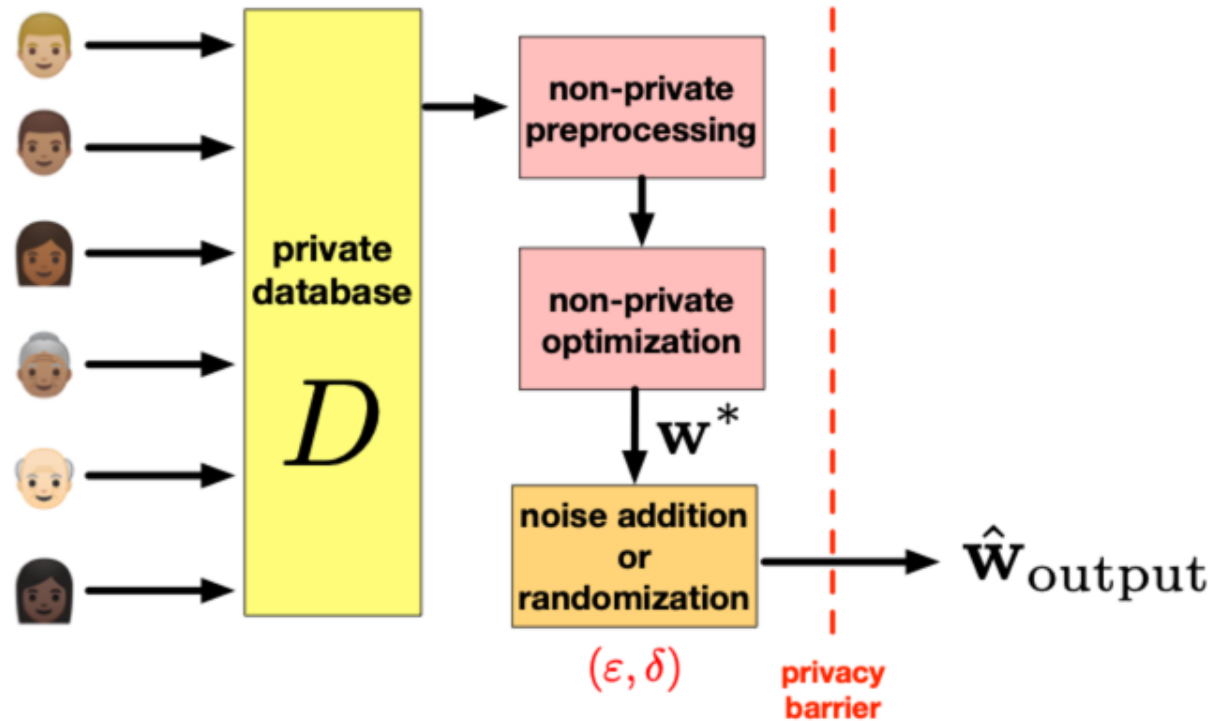
# Local Privacy



- Local privacy: data contributors sanitize data before collection.

- Classical technique: *randomized response* [W65].

- Interactive variant can be minimax optimal [DJW13].

# Input Perturbation



- Input perturbation: add noise to the input data.

- Advantages: easy to implement, results in reusable sanitized data set.
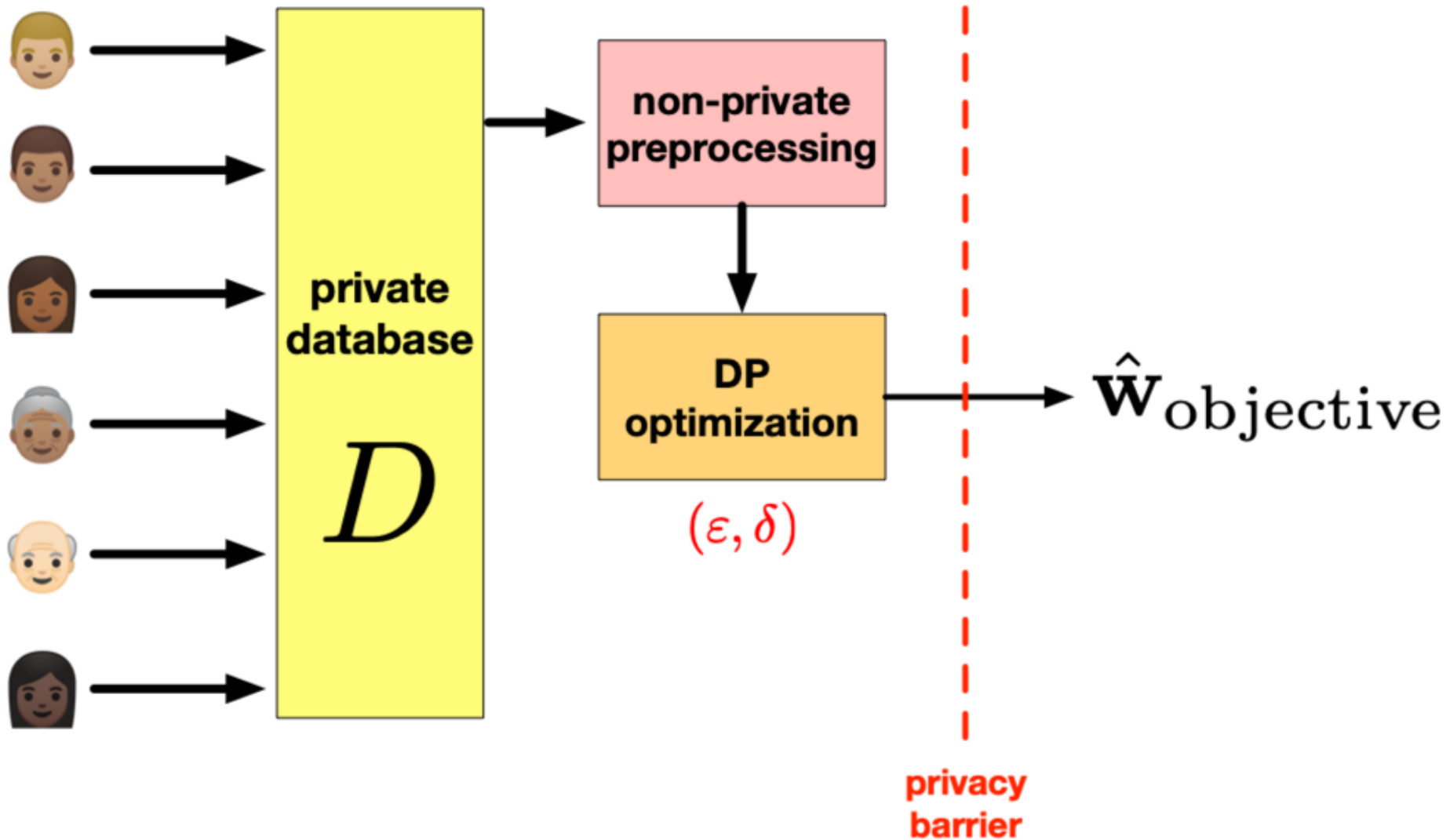
[DJW13,FTS17]

# Output Perturbation



- Compute the minimizer and add noise.

- Does not require re-engineering baseline algorithms

Noise depends on the sensitivity of the argmin.

[CMS11, RBHT12]

# Objective Perturbation

- Embrace the **"Bright and the Dark Side"** as a community

  - let's **better understand** and **control privacy**

  - let's **better understand** and **control security**

- Do **not leave this topic to companies** (alone) !

  - keep **knowledge** in the **public** domain

  - develop algorithms and methods — also to **pressure companies** to adopt them

- Responsibility in **education**

  - **educate students** about **both opportunities** and **potential dangers**

  - **distinguish** between **"what can be done"** and "**what should be done" (Weizenbaum)**

# High Level Computer Vision:

# Attacks on Computer Vision Models

Mario Fritz fritz@cispa.saarland

Bernt Schiele schiele@mpi-inf.mpg.de

17.7.2019