

MACHINE LEARNING AND PRIVACY

Pei-Yuan Wu

National Taiwan University

Dept. Electrical Engineering

peiyuanwu@ntu.edu.tw

1

OUTLINE

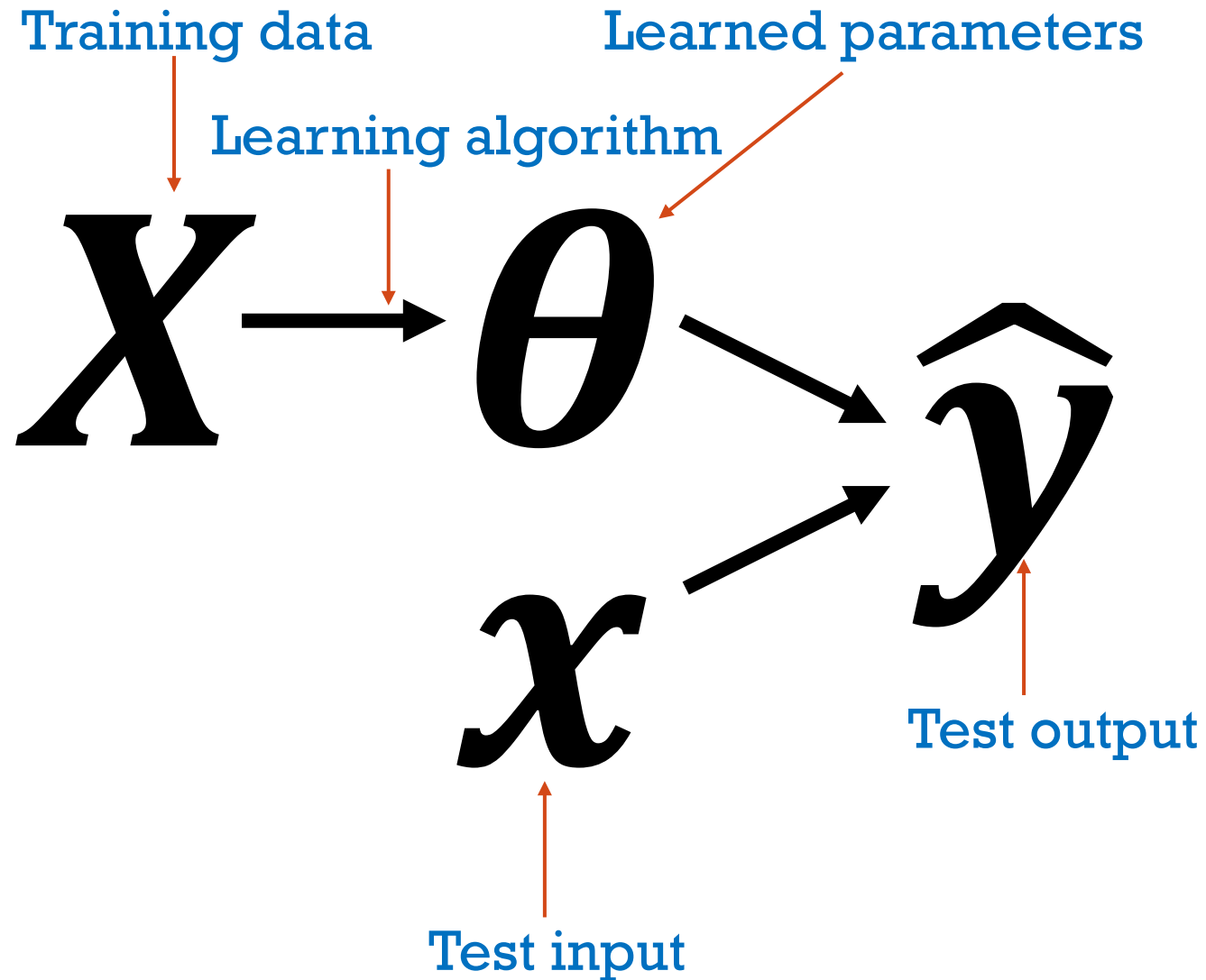
- Security against Attacks that use Machine Learning
 - Model Inversion Attack
 - Membership Inference Attack
 - Adversarial examples attack
 - Differential Privacy
 - Homomorphic Encryption
- Homomorphic encryption example
 - Privacy-Preserving Principle Component Analysis
- Secure Multi Party Computation
- Compressive Privacy Example
- Provable defense against adversarial example attack



SECURITY AGAINST ATTACKS THAT USE MACHINE LEARNING



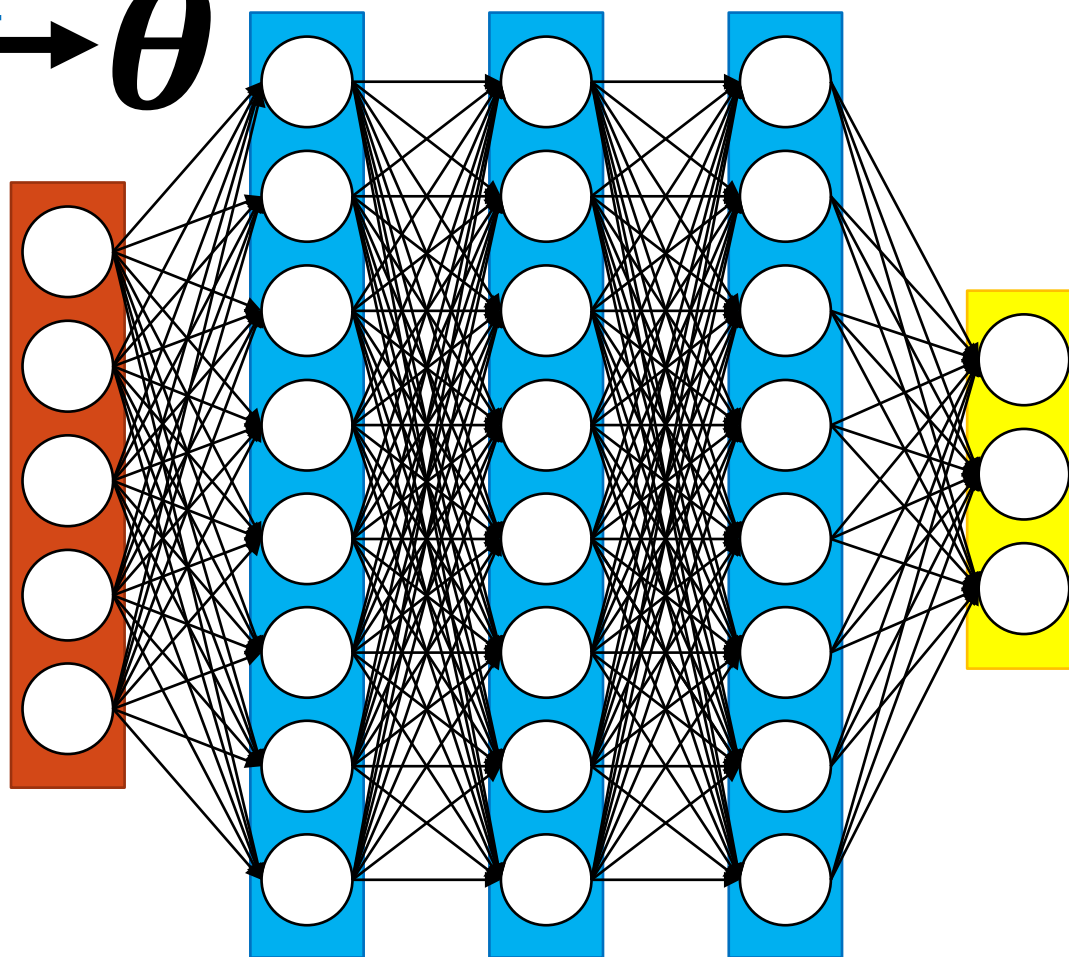
MACHINE LEARNING PIPELINE



X $\xrightarrow{\text{Learning algorithm}}$ θ

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training data



Test input

- 0 2%
- 1 10%
- 2 3%
- 3 2%
- 4 50%
- 5 4%
- 6 3%
- 7 10%
- 8 1%
- 9 15%

\hat{y}

Test output

5

MODEL INVERSION ATTACK

■ Adversary Target

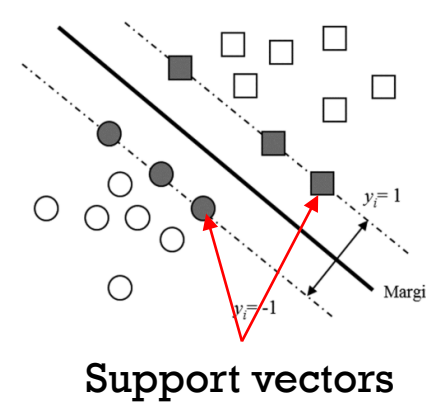
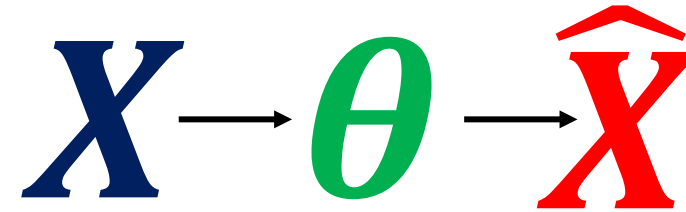
- Create dataset \hat{X} resemble X used to create model θ

■ Attack Scheme

- Support Vector Machine model reveals training data as support vectors.
- Exploit confidence information:
 - ✓ Many APIs reveal confidence values along with class predictions.
 - ✓ Find training samples for a peculiar class as the input yielding highest confidence on that peculiar class.

■ Remedy

- Only allow black-box access to the model.
- Confidence values not revealed or rounded.



Face in training data



Reconstructed Face by model inversion attack



Black-box face reconstruction with rounding confidence

MEMBERSHIP INFERENCE ATTACK

■ Adversary Target

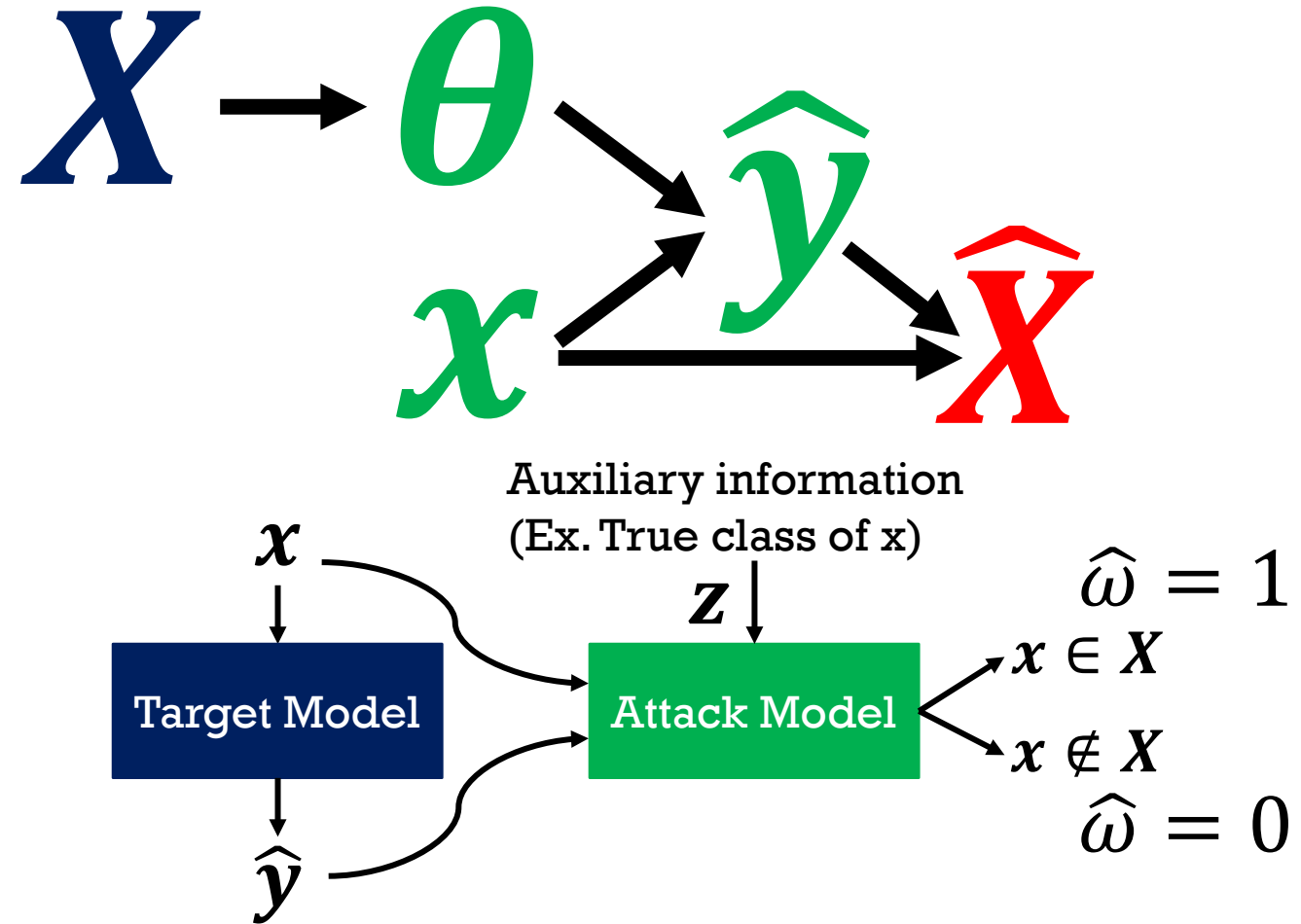
- Determine if a sample is in X used to create model θ .
- ✓ Example: Was Bob's record used to train ML model associated with AIDS?

■ Attack Scheme

- Exploits the difference between predictions made on training samples versus unseen samples.

■ Remedy

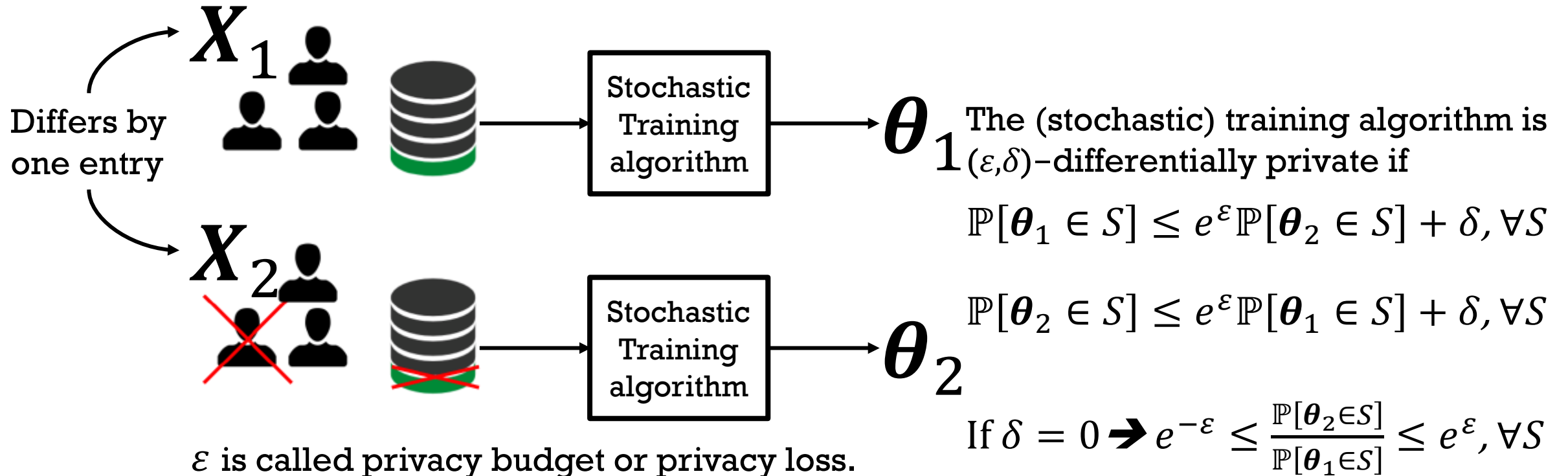
- Coarser precision of confidence values, only reveal top k confidence values.
- Differential privacy.



R. Shokri, M. Stronati, C. Song and V. Shmatikov, "Membership Inference Attacks Against Machine Learning Models," *2017 IEEE Symposium on Security and Privacy (SP)*, San Jose, CA, 2017, pp. 3-18.

(ϵ, δ) -DIFFERENTIAL PRIVACY

- The presence/absence of an entry in the training data has little effect on the trained parameters \rightarrow Difficult to perform membership inference attack



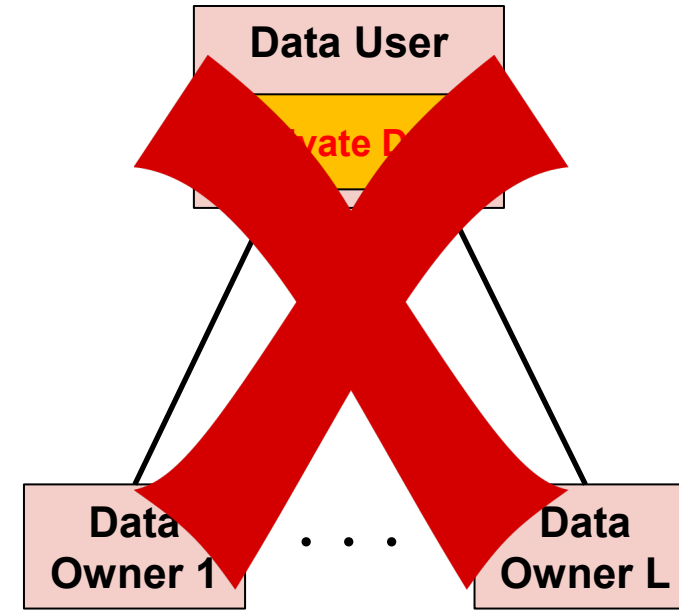
ϵ is called privacy budget or privacy loss.

Small $\epsilon \rightarrow \theta_1$ and θ_2 has similar probability distributions

\rightarrow Difficult to infer \mathbf{X} from θ

CRYPTOGRAPHIC APPROACHES TO DISTRIBUTED MACHINE LEARNING

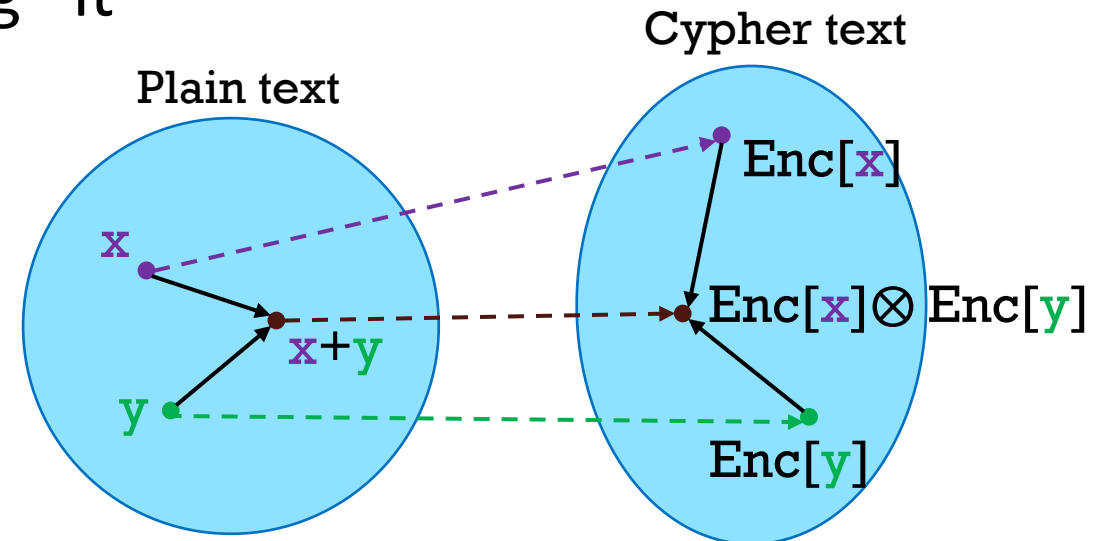
- In collaborative learning involving **multiple data owners**, we need a privacy-aware distributed approach:
 - The data user should build machine learning models, while the data owners keep their data **private**.
- Special cryptographic approaches allow computing the data without “seeing” it



- **Additive homomorphic encryption:**

Addition on original data = Modular multiplication on encrypted data

$$Enc[x + y] = Enc[x] \otimes Enc[y]$$



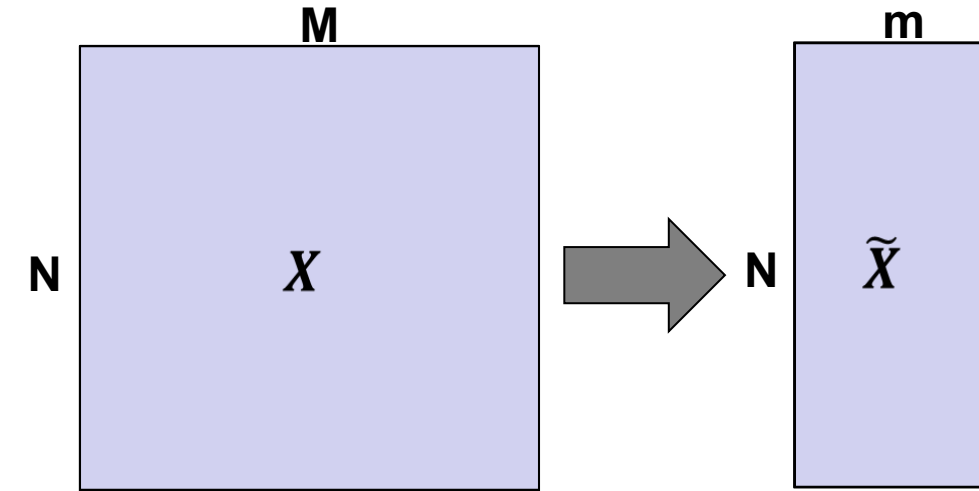


HOMOMORPHIC ENCRYPTION EXAMPLE: PRIVACY-PRESERVING PRINCIPLE COMPONENT ANALYSIS

PRIVACY-PRESERVING PRINCIPLE COMPONENT ANALYSIS

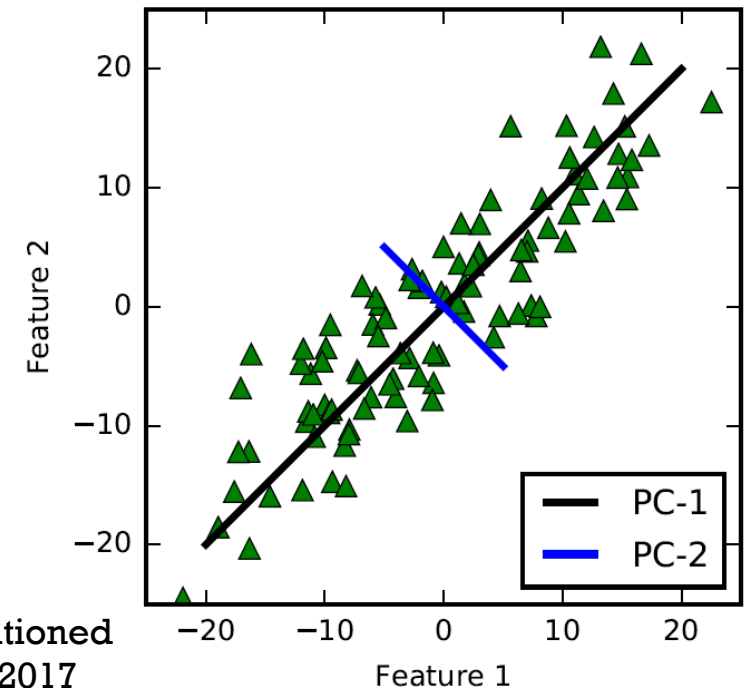
■ Dimensionality Reduction:

- Map the original high-dimensional data onto a lower-dimensional subspace.
- Reduces the noise, compresses the data, decreases the computational cost, and prevents over-fitting.



■ Principle Component Analysis:

- Target: Find the best sub-space that preserves most of the variance in the original data
- The principal axes are orthogonal, uncorrelated and ordered by how much variability they retain.



CENTRALIZED PCA

- PCA performs Eigen-decomposition of the center-adjusted scatter matrix

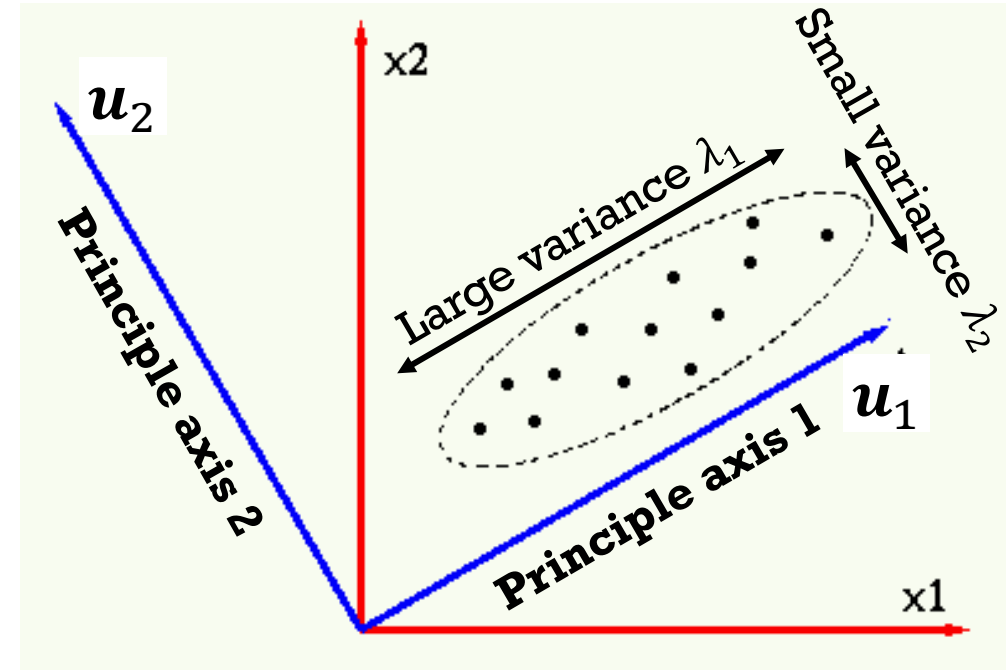
$$S = \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T = U\Lambda U^T$$

$U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M]$ is a unitary matrix

Principle axes

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \lambda_M \end{bmatrix} \text{ is a diagonal matrix}$$

Eigen Values



- In centralized PCA, the original data is needed to compute the scatter matrix (suitable only for **single data owner**)
- Privacy-aware distributed approach is needed for **multiple data owners**.

DISTRIBUTED SCATTER MATRIX COMPUTATION

- Rewrite scatter matrix as follows:

$$\mathbf{S} = \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T - N \boldsymbol{\mu} \boldsymbol{\mu}^T = \mathbf{R} - \frac{1}{N} \mathbf{v} \mathbf{v}^T$$

$$\text{where } \mathbf{R} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T, \mathbf{v} = \sum_{i=1}^N \mathbf{x}_i$$

- Each data owner computes the share: $DS^\ell = \{\mathbf{R}_\ell, \mathbf{v}_\ell, N_\ell\}$

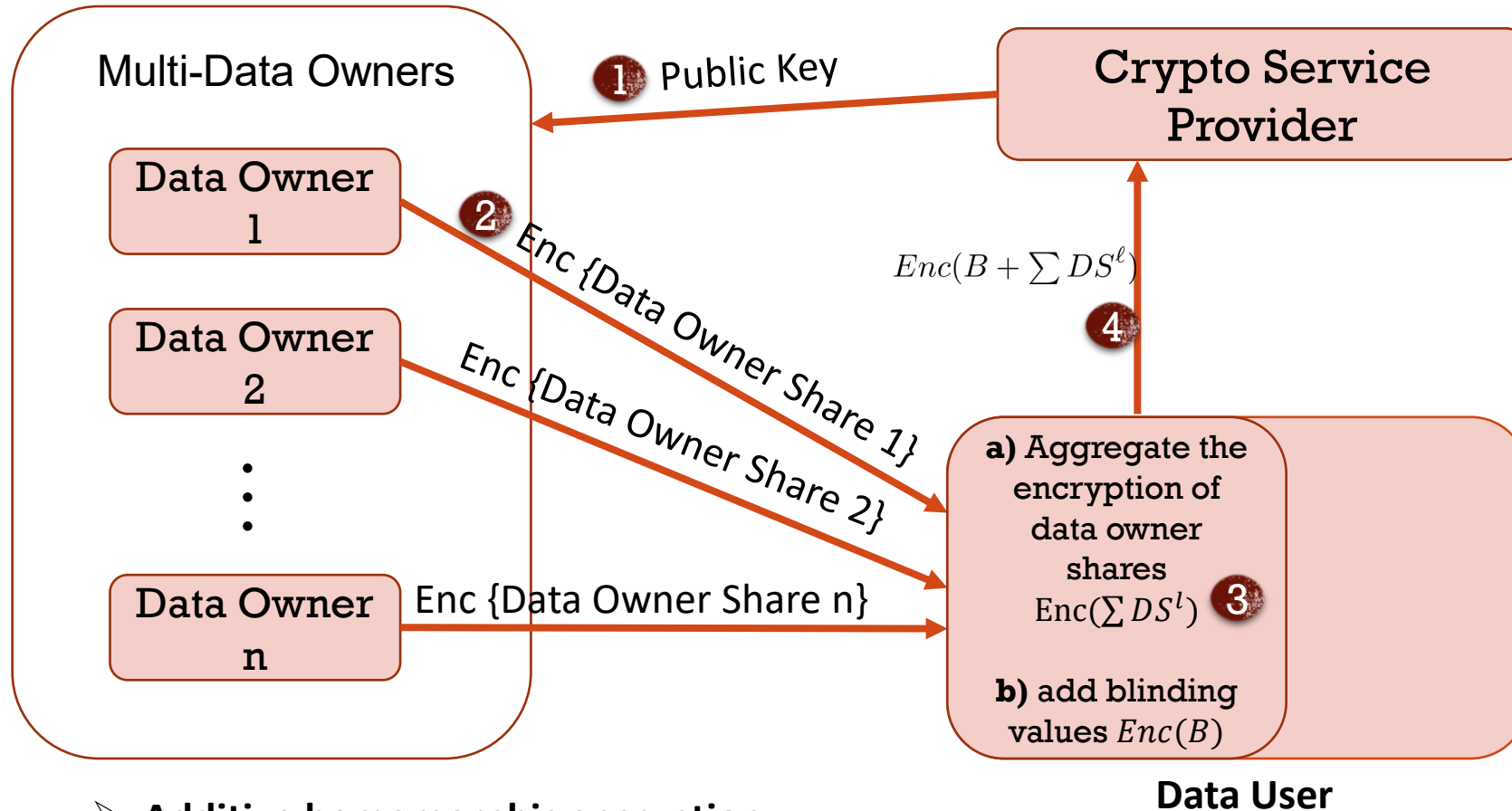
$$\mathbf{R}_\ell = \sum_{i \in P_\ell} \mathbf{x}_i \mathbf{x}_i^T, \mathbf{v}_\ell = \sum_{i \in P_\ell} \mathbf{x}_i, N_\ell = |P_\ell|$$

P_ℓ is the set of training samples from data owner ℓ

- The scatter matrix is aggregated from all data owners:

$$\mathbf{S} = \mathbf{R} - \frac{1}{N} \mathbf{v} \mathbf{v}^T, \mathbf{R} = \sum_{\ell} \mathbf{R}_\ell, \mathbf{v} = \sum_{\ell} \mathbf{v}_\ell, N = \sum_{\ell} N_\ell$$

ARCHITECTURE-SCATTER MATRIX COMPUTATION

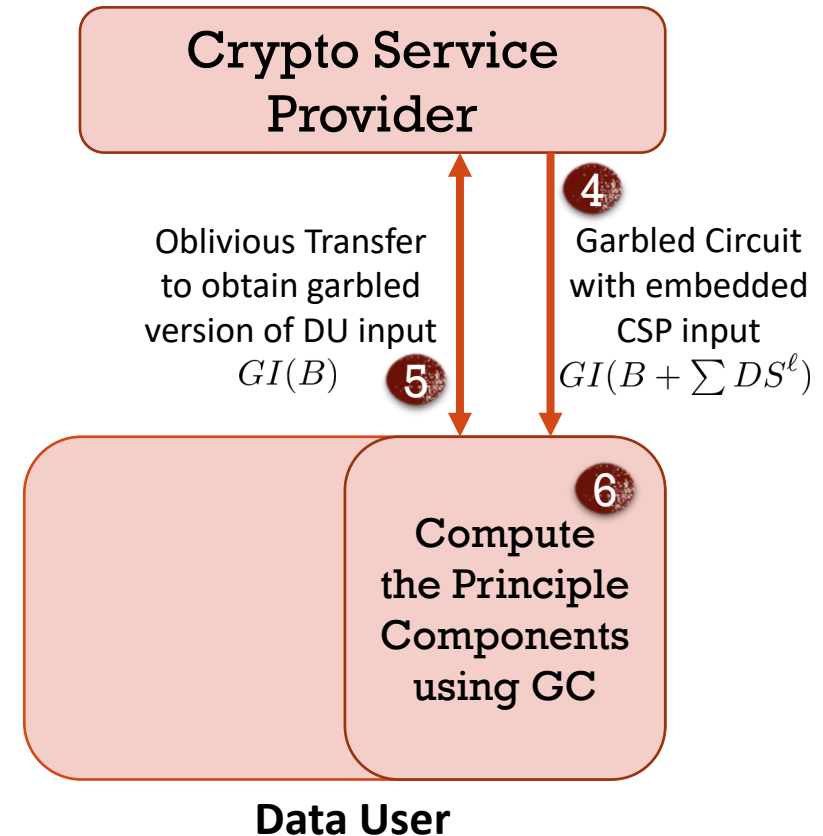
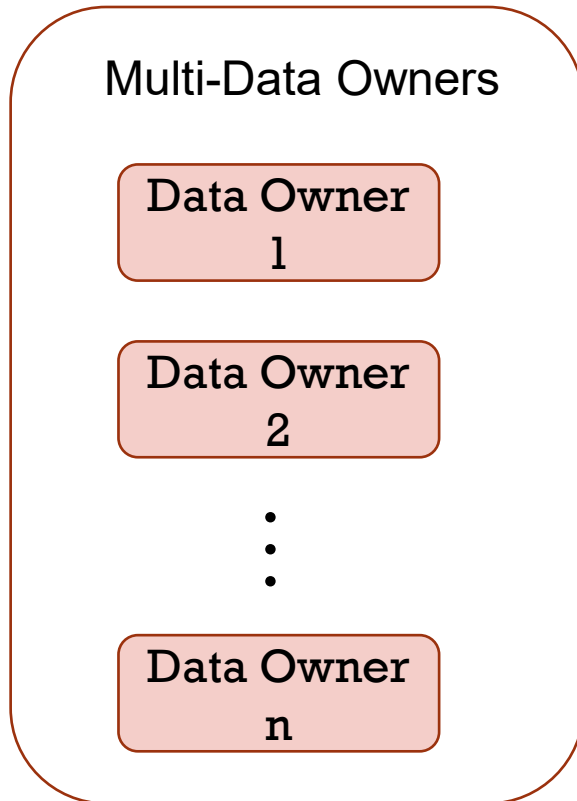


➤ **Additive homomorphic encryption:**

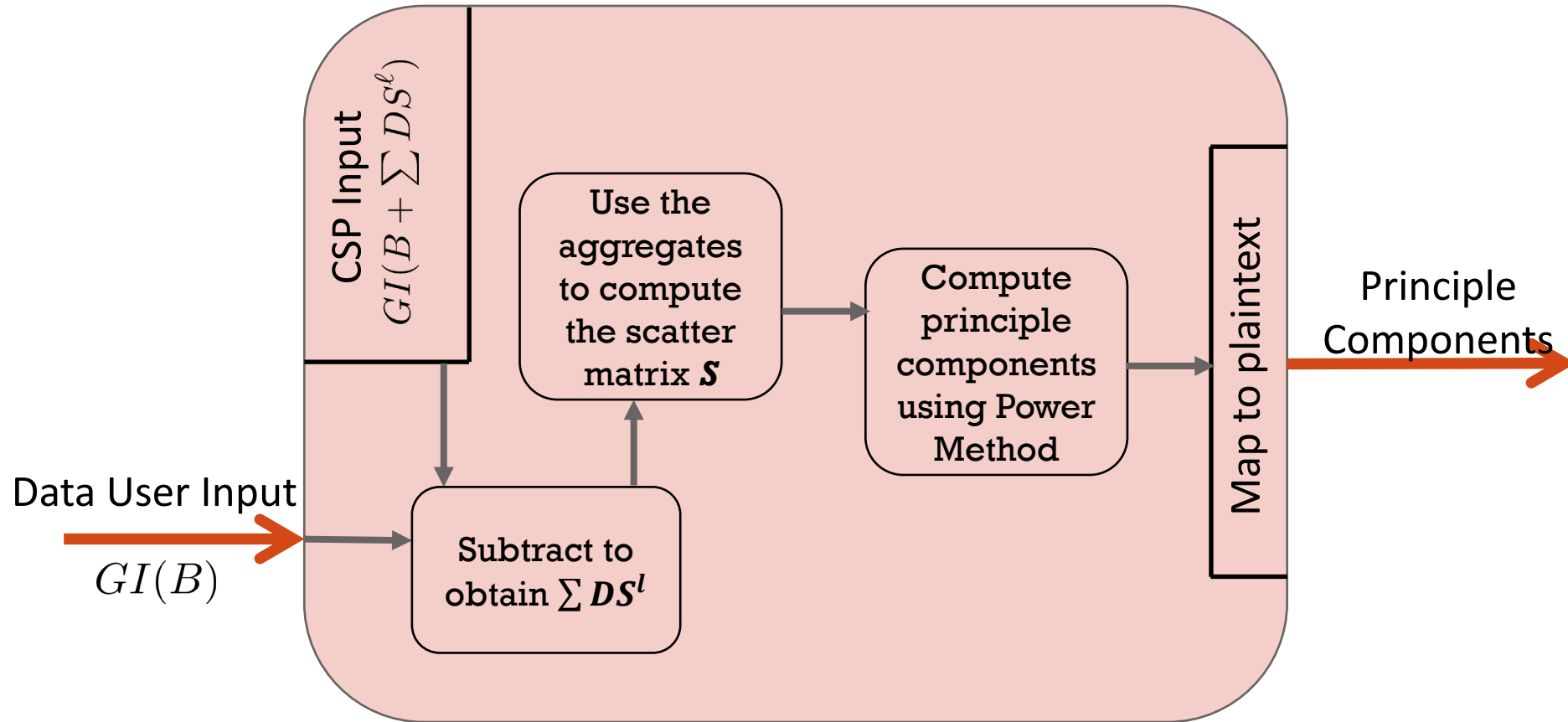
Addition on original data = Modular multiplication on encrypted data

$$Enc[x + y] = Enc[x] \otimes Enc[y]$$

ARCHITECTURE-PRINCIPLE AXIS COMPUTATION

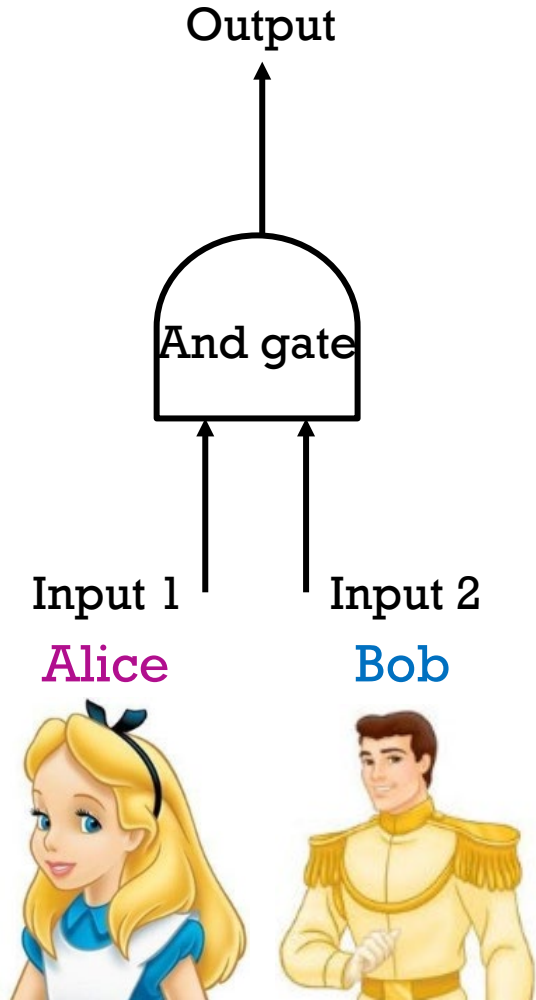


ARCHITECTURE-PRINCIPLE AXIS COMPUTATION (CONT'D)















GARbled CIRCUIT

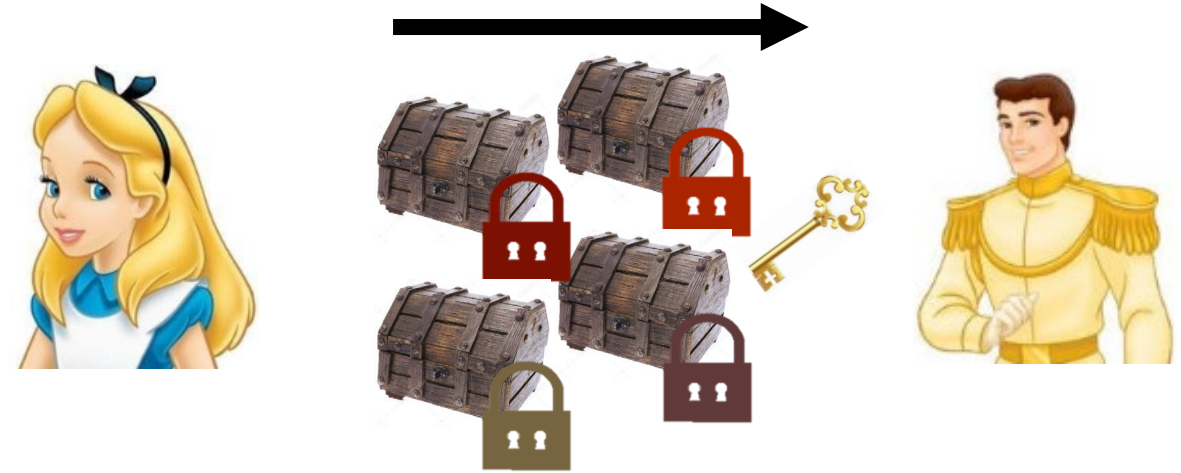
Alice and Bob want to compute AND gate together, but not revealing their own inputs



Alice makes the keys and locks, and lock the cases

Input 1	Input 2	Output
 0	 0	 0
 0	 1	 0
 1	 0	 0
 1	 1	 1

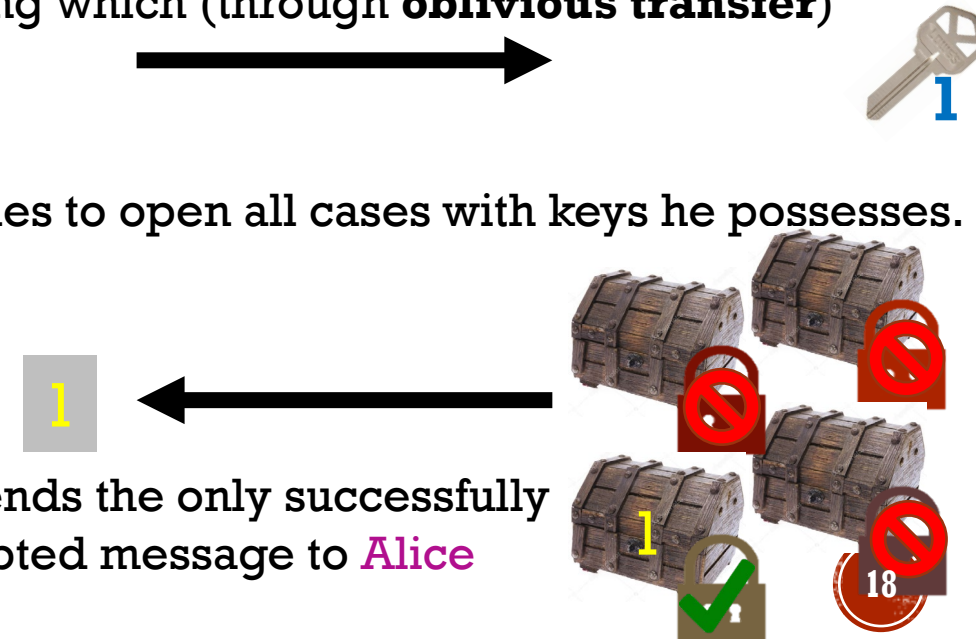
Alice shuffles the cases and gives them to Bob the locked cases, as well as her selected key



Bob picks his key from Alice without her knowing which (through oblivious transfer)

Bob tries to open all cases with keys he possesses.

Bob sends the only successfully decrypted message to Alice



COMPUTATION COSTS

$Enc[DS^\ell]$

PCA Eigen Decomposition
with Garbled circuit

Dataset	Features	Classes	Avg. DO time	Avg. DU Coll. / Add time	CSP Dec. time	DU PCA Comp. time
Diabetes	8	2	0.63 sec	10 ms	0.67 sec	28.3 sec (8)
Breast Cancer	10	2	0.93 sec	11 ms	1 sec	49.6 sec (8)
Australian	14	2	1.7 sec	12 ms	1.8 sec	119.1 sec (8)
German	24	2	5 sec	17 ms	5 sec	16.3 min (15)
Ionosphere	34	2	9.8 sec	24 ms	9.9 sec	43.2 min (15)
SensIT Acoustic	50	3	22.5 sec	40 ms	22.7 sec	126.7 min (15)

$$Enc \left[B + \sum_{\ell} DS^{\ell} \right] = Enc[B] \otimes_{\ell} Enc[DS^{\ell}]$$

$$Dec \left[Enc \left[B + \sum_{\ell} DS^{\ell} \right] \right] = B + \sum_{\ell} DS^{\ell}$$

CPU: i5-6600K @ 3.5GHz

RAM: 8 GB

Paillier's Cryptosystem with 1024 bits key length

No multi-threading was used

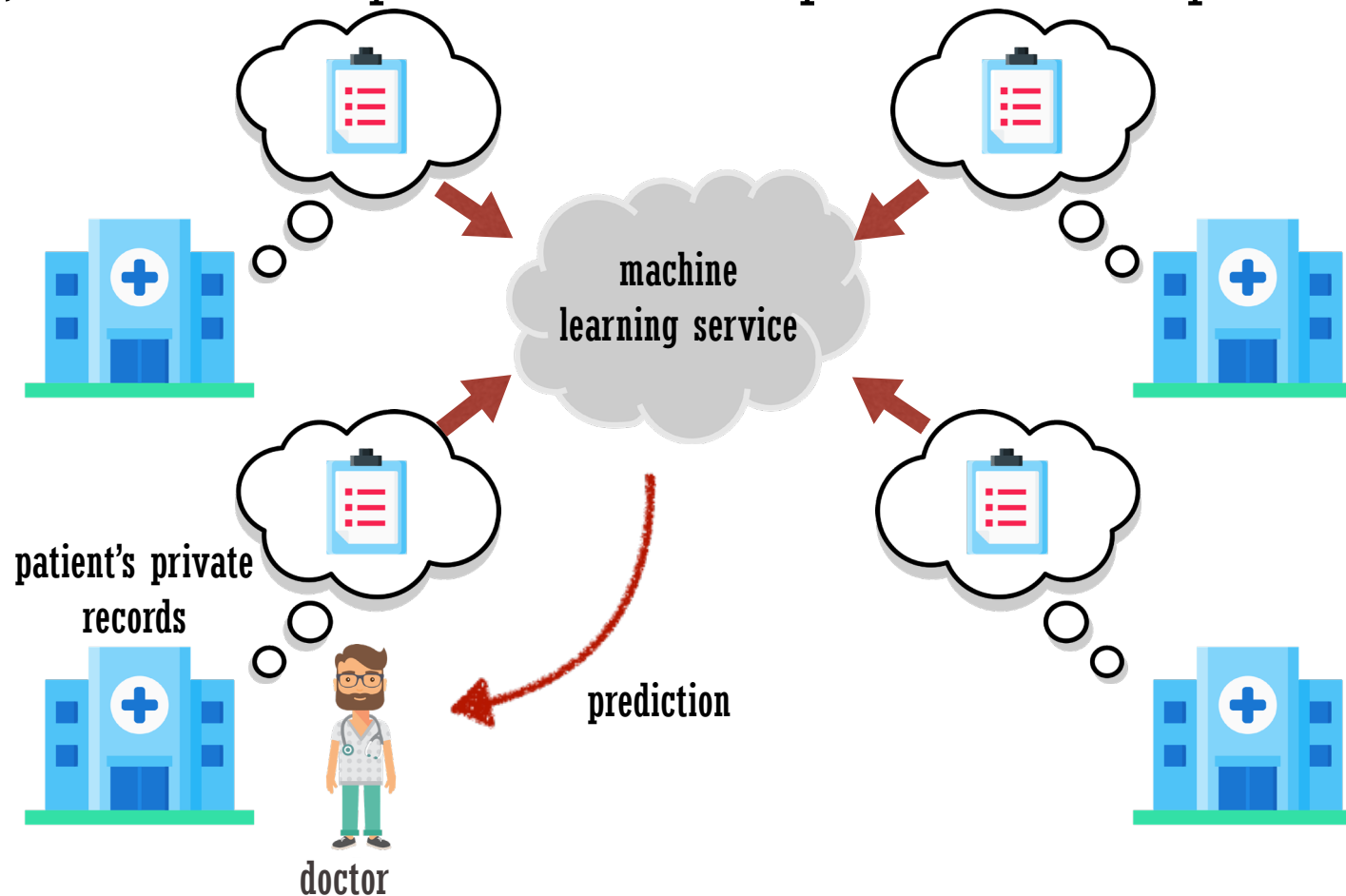


SECURE MULTI-PARTY COMPUTATION



SECURE MULTI-PARTY COMPUTATION (SMPC)

- **Goal:** For multiple parties to jointly compute a function over their inputs, while each party keeping their inputs private from other parties.
- **Example:** Suppose we wish to train a diagnostic model by data from multiple hospitals, while each hospital wishes to keep their own data private.



S MPC EXAMPLE

Suppose three people have monthly salary 5K, 100K, 22K, respectively. How can they compute the average salary while keeping their own salary in secret?



5,000



100,000



22,000



S MPC EXAMPLE

Suppose three people have monthly salary 5K, 100K, 22K, respectively. How can they compute the average salary while keeping their own salary in secret?

Step 1: Every one split their own salary into three parts.



$$5K = (1K, 1K, 3K)$$



$$100K = (22K, 5K, 73K)$$



$$22K = (5K, 11K, 6K)$$

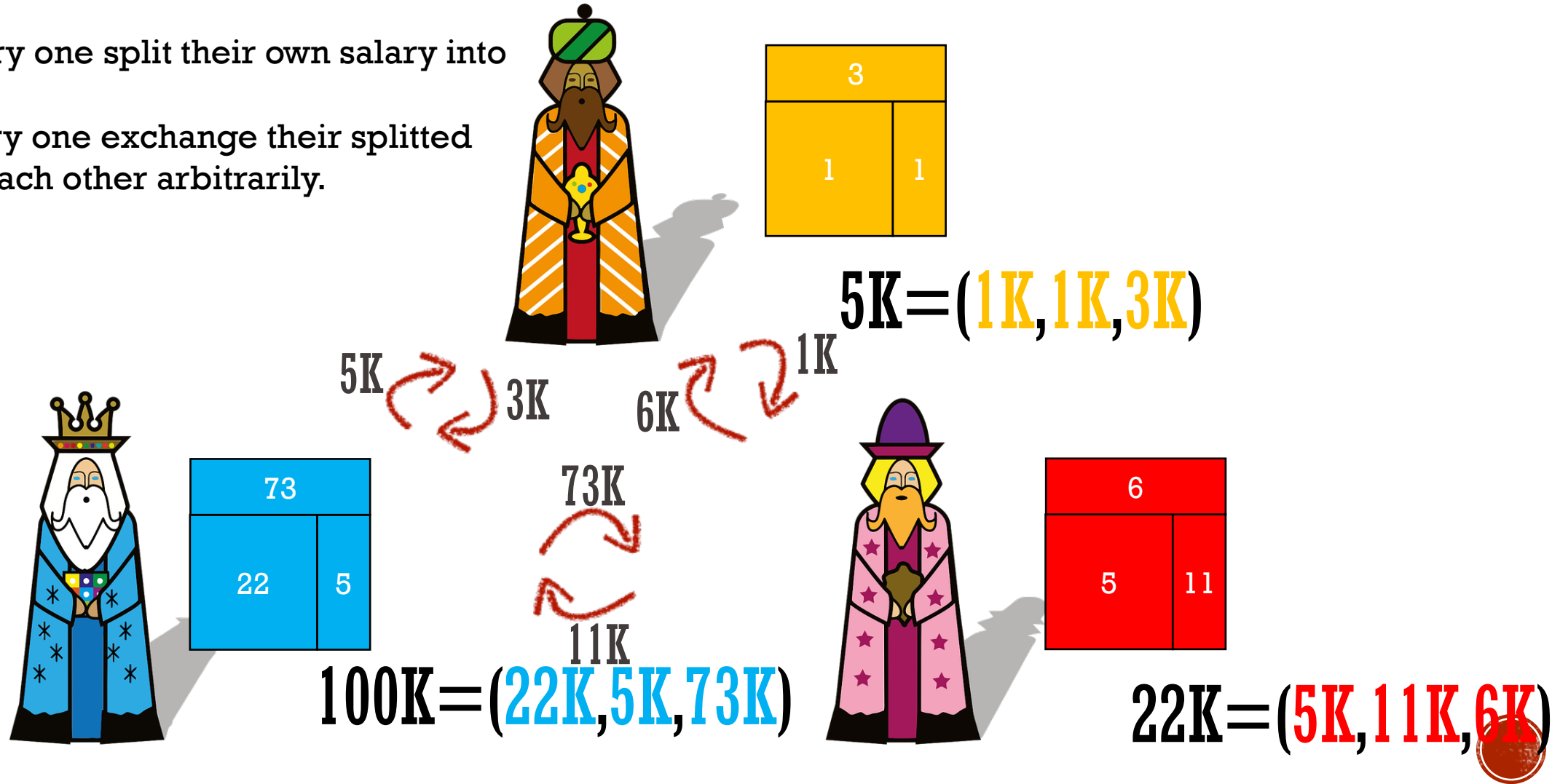


S MPC EXAMPLE

Suppose three people have monthly salary 5K, 100K, 22K, respectively. How can they compute the average salary while keeping their own salary in secret?

Step 1: Every one split their own salary into three parts.

Step 2: Every one exchange their splitted parts with each other arbitrarily.



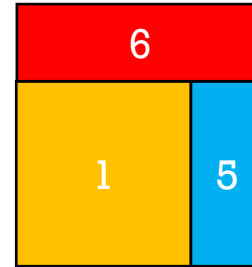
SMPc EXAMPLE

Suppose three people have monthly salary 5K, 100K, 22K, respectively. How can they compute the average salary while keeping their own salary in secret?

Step 1: Every one split their own salary into three parts.

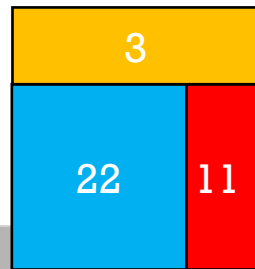
Step 2: Every one exchange their splitted parts with each other arbitrarily.

Step 3: Every one computes the sum of their received part, then compute total average.

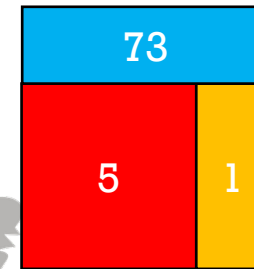


$$12K = 1K + 5K + 6K$$

$$42.3K = (36K + 12K + 79K) / 3$$



$$36K = 22K + 11K + 3K$$



$$79K = 5K + 1K + 73K$$



SMPC FOR MACHINE LEARNING

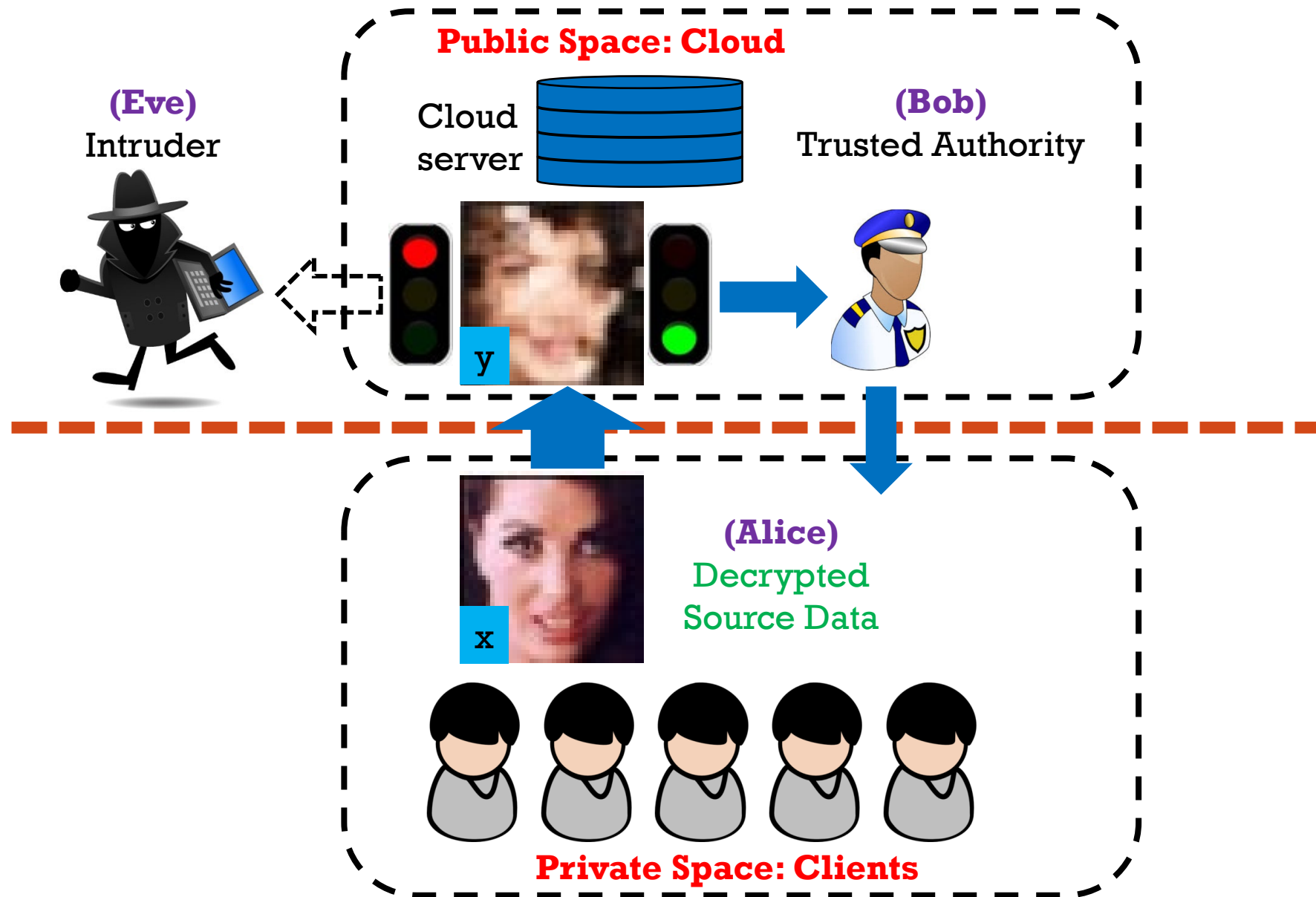
- **SecureML:** Apply SMPC to linear regression, logistic regression and neural network training using the stochastic gradient descent method.
- **CrypTen:** Open source framework built on PyTorch. Developed by Facebook.
- Computation cost is a BIG issue:
Upon training a 3-layer fully-connected neural network:
 - Plain model: 9 ms per epoch
 - SecureML: 4 mins per epoch
 - CrypTen: 15 mins per epoch



COMPRESSIVE PRIVACY EXAMPLE



Compressive Privacy Paradigm

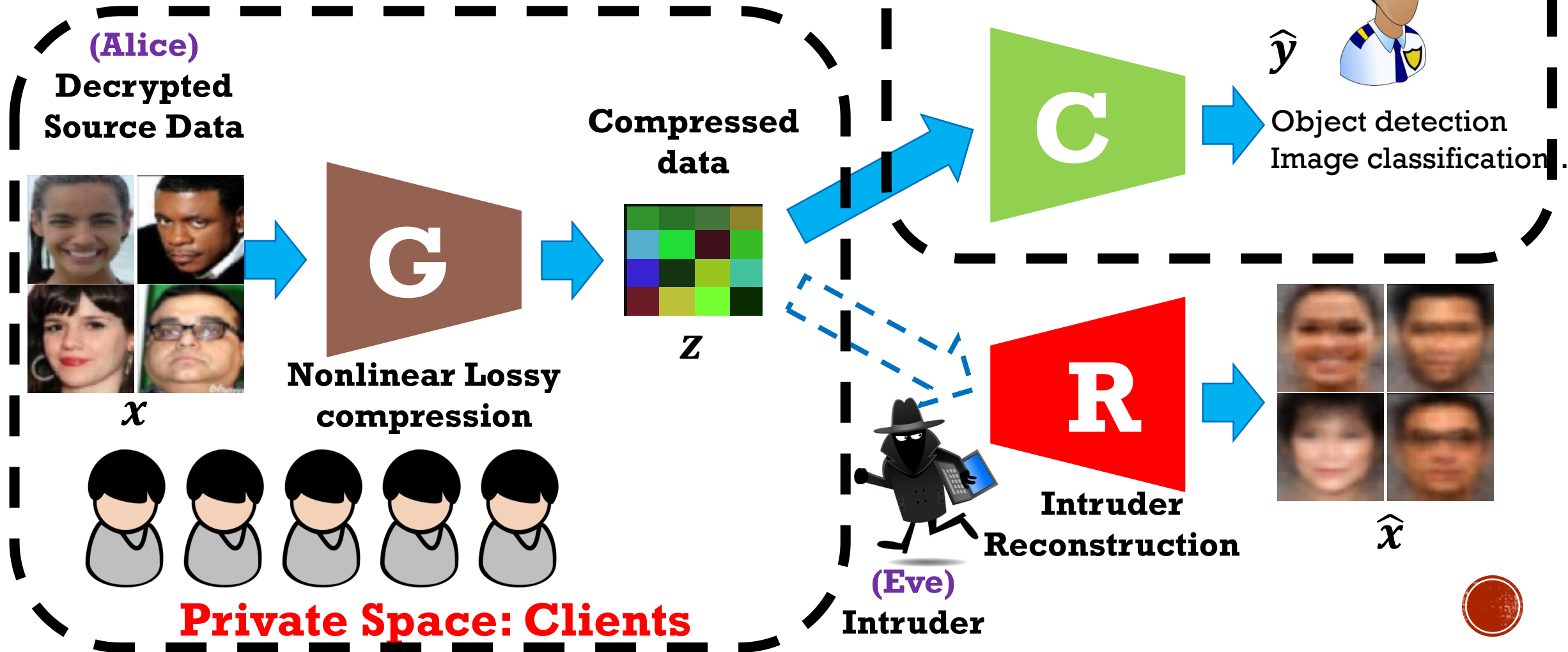


Kung, S. Y. (2018). A Compressive Privacy approach to Generalized Information Bottleneck and Privacy Funnel problems. *Journal of the Franklin Institute*, 355(4), 1846-1872.

Compressive Privacy Generative Adversarial Network

Objective: Bob does well but Eve does poorly

$$\max_{\mathbf{G}} \left(\min_{\mathbf{R}} \sum_i \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \lambda \max_{\mathbf{C}} \sum_i \log P(\hat{t}_i = t_i) \right)$$



CPGAN FOR BENCHMARK DATASET

- **Synthetic dataset:**

- Sampled from Gaussian mixture data model with binary class.
- Training/testing samples: 20K/2K

- **MNIST:**

- Training/testing samples: 55000/10000

- Examples:

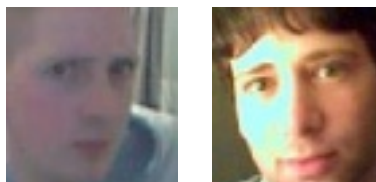


- **UCI Human activity recognition (HAR) dataset**

- Given the time-series sensor record from ten identities.
- Six activities: walking, sitting, standing e.t.c.

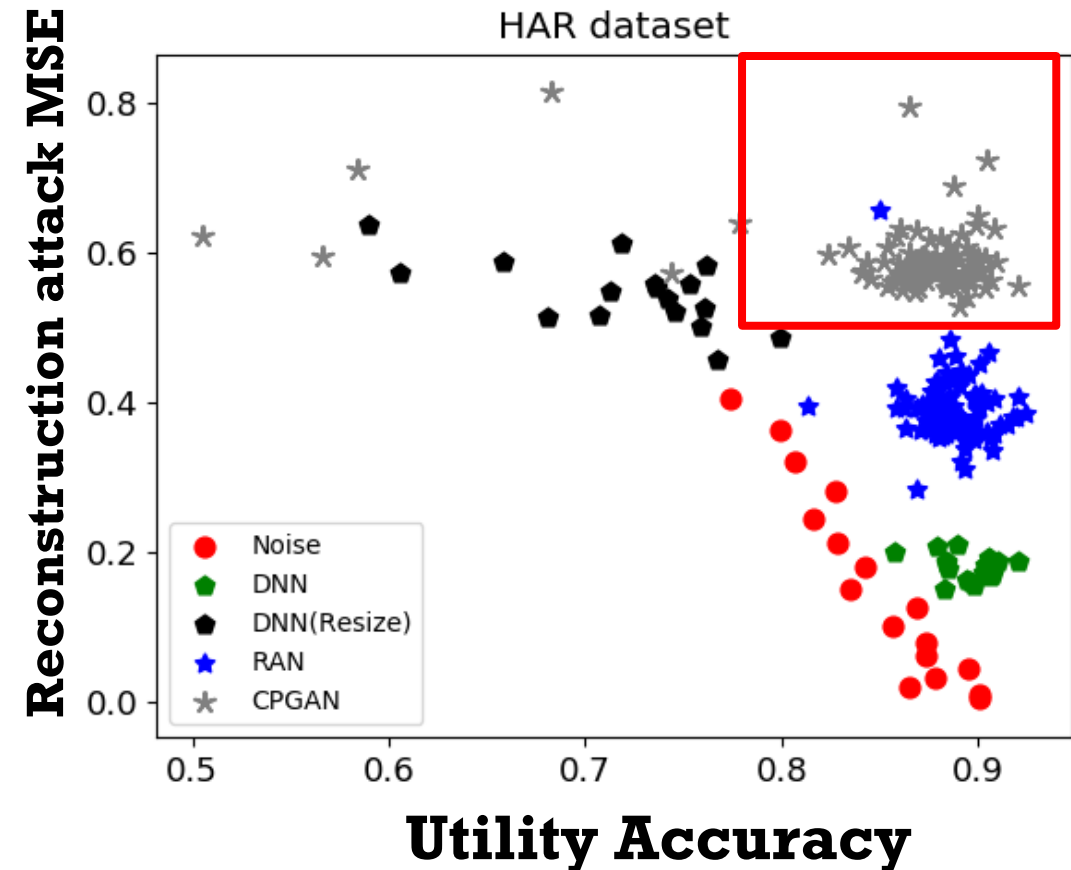
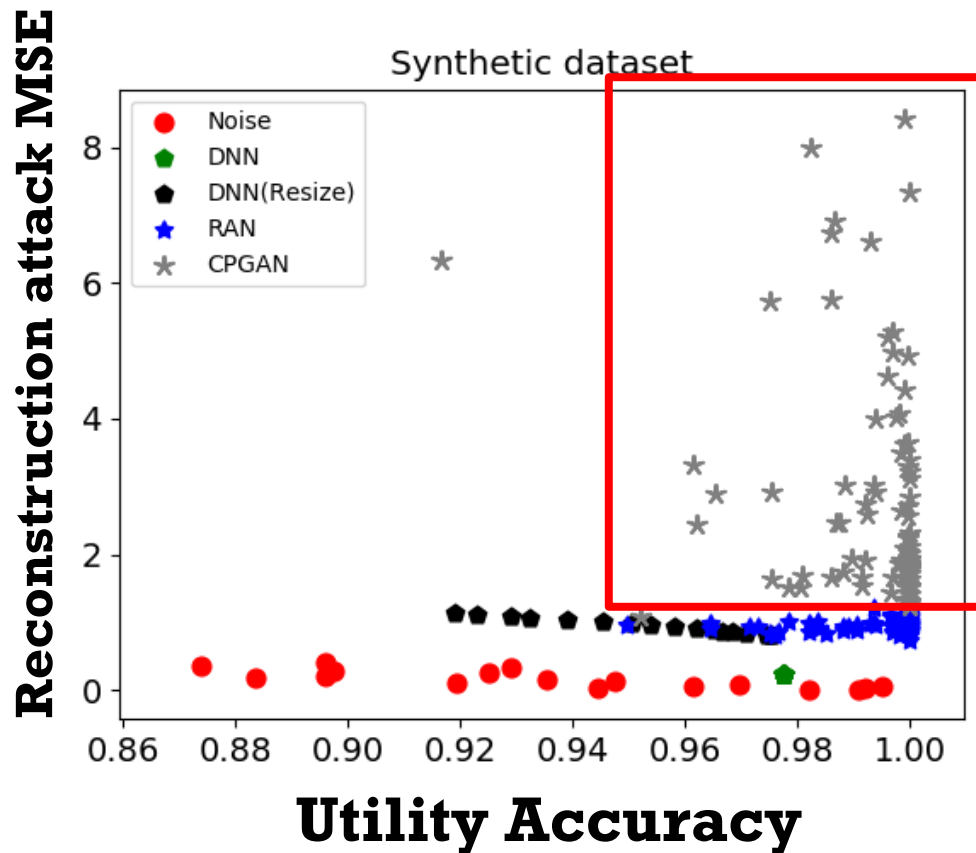
- **Genki-4K dataset:**

- Face images with 400 sample. Detect the expression of this image.
- Example:



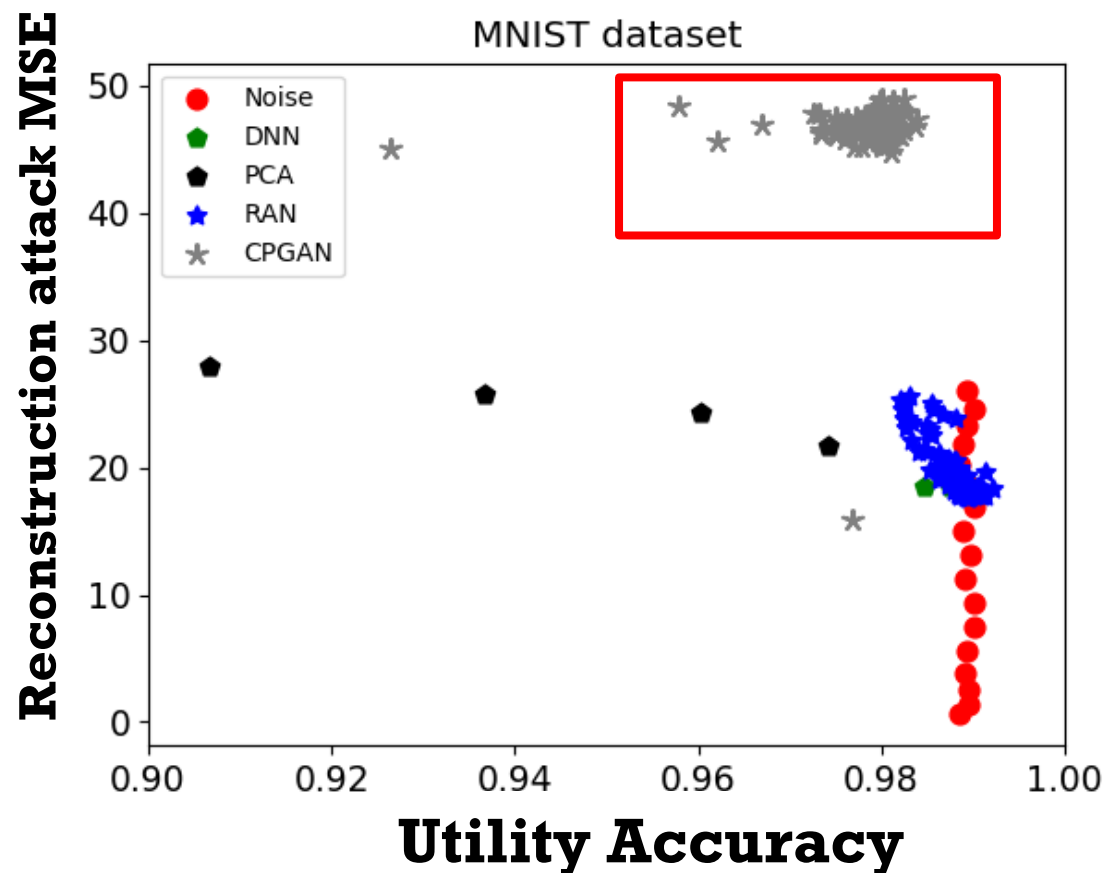
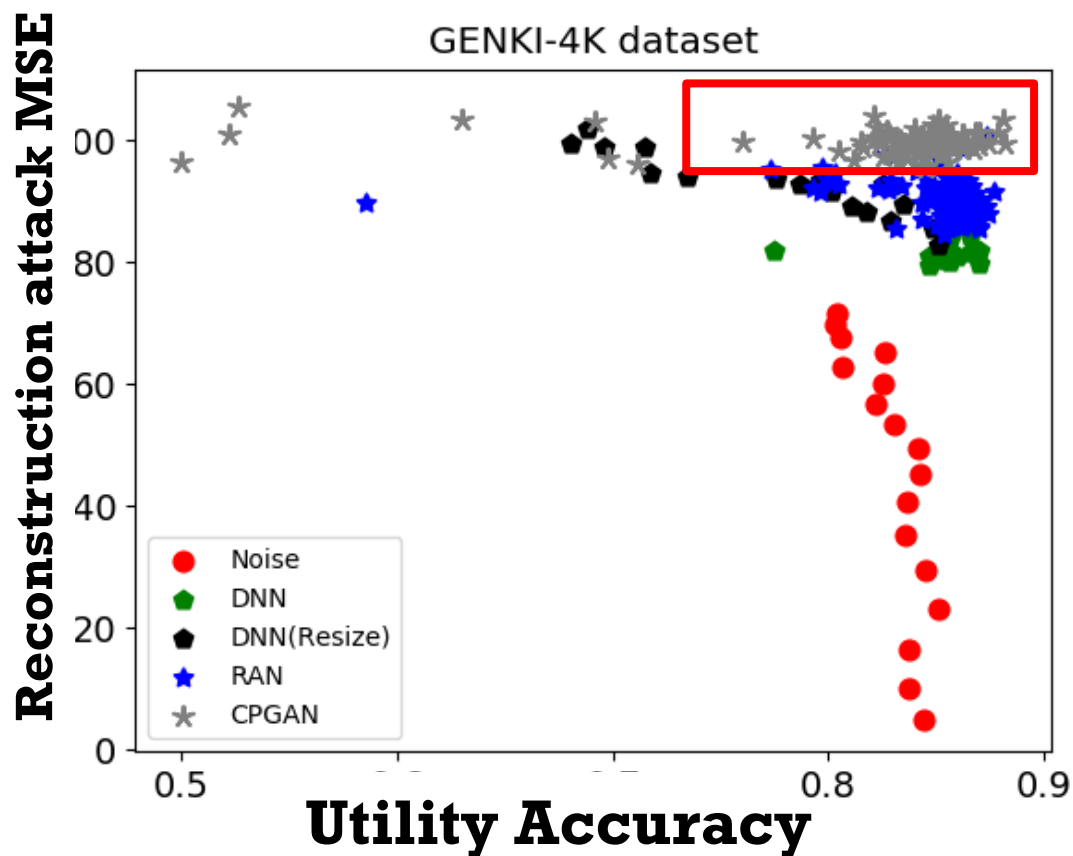
COMPARISON WITH PREVIOUS WORKS

CPGAN (Gray star) outperforms other methods on privacy perspective, but slightly drops (less than 1%) the utility accuracy (trade-off).



COMPARISON WITH PREVIOUS WORKS (CONT.)

CPGAN (Gray star) outperforms other methods on privacy perspective, but slightly drops (less than 1%) the utility accuracy (trade-off).



RESULTS ON CIFAR-10 AND SVHN

CIFAR-10:

Adversary reconstruction

Original images



Classify by original image: 96.5%

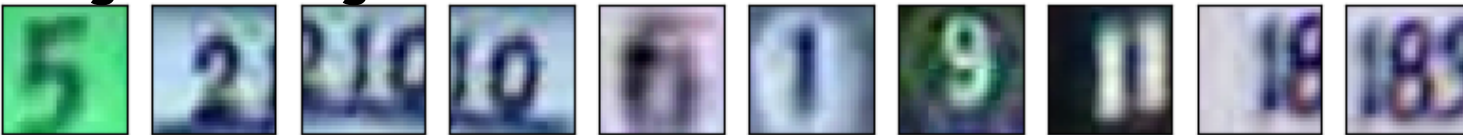
Reconstructed image by adversary



Classify by compressed data: 93.9%

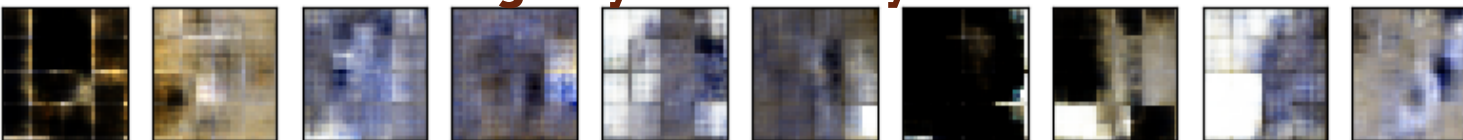
SVHN:

Original images



Classify by original image: 98.6%

Reconstructed image by adversary



Classify by compressed data: 97.7%

CPGAN defends the *reconstruction attack* while achieving satisfactory *utility performance*

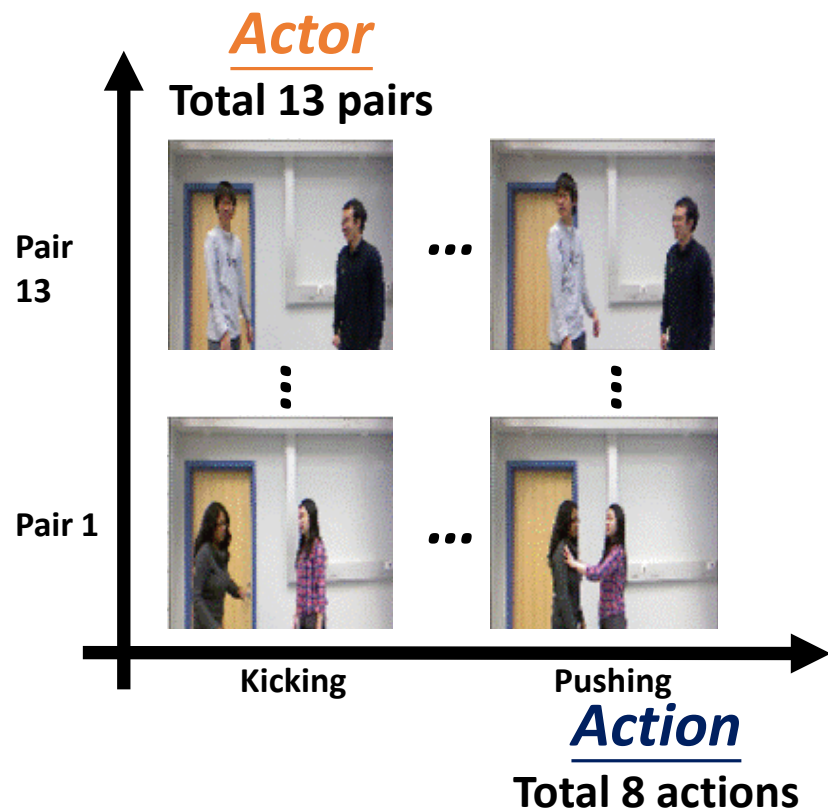


Compressive Privacy on Videos

Identity Privacy Preserving

SBU Kinect Interaction Dataset:

Training / Testing: 346 / 36

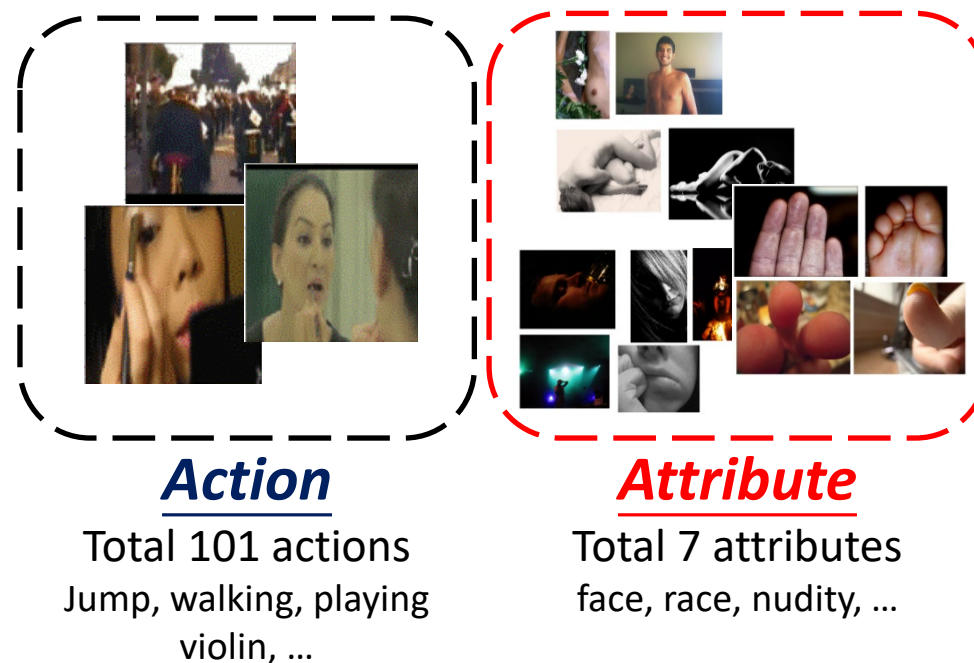


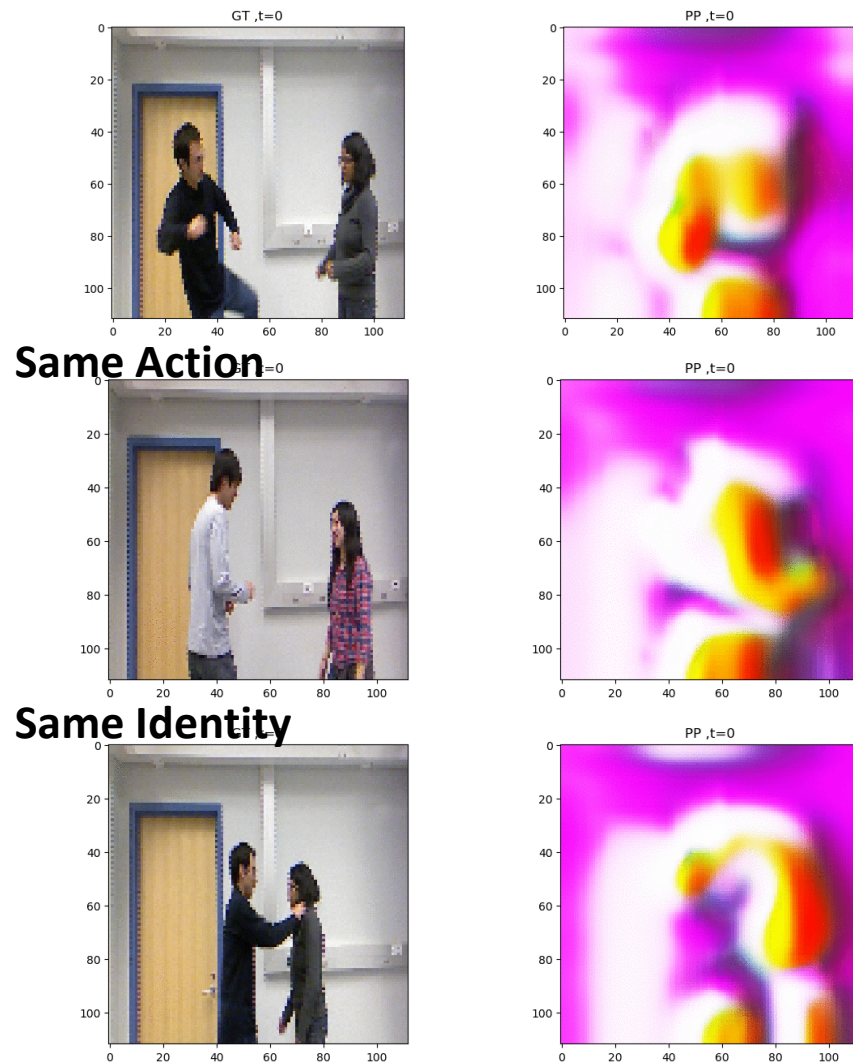
UCF101 Dataset

Training / Testing: 9537 /

3783

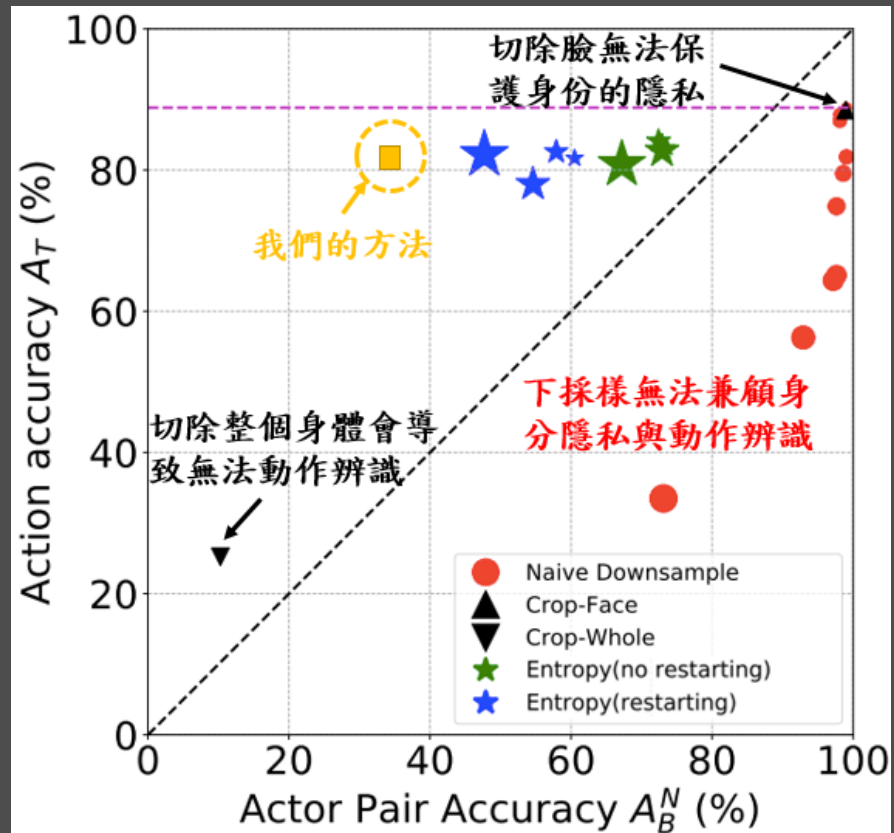
Source: YouTube

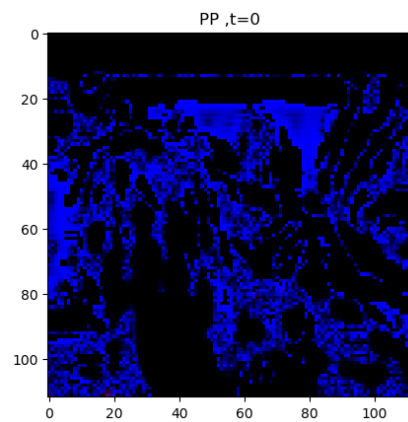
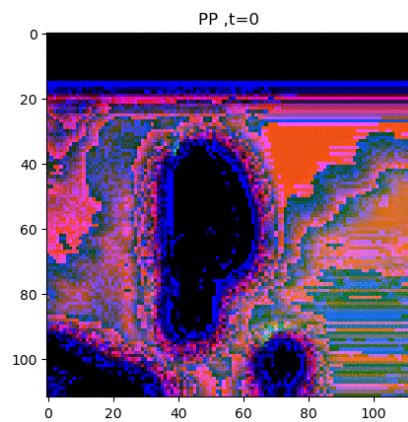
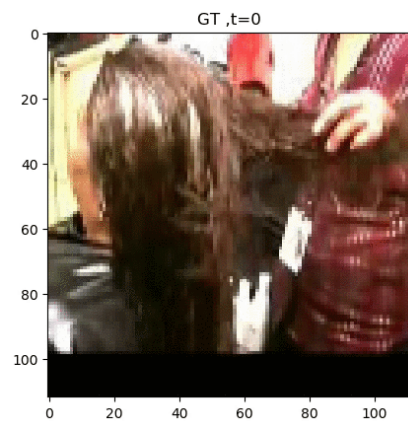
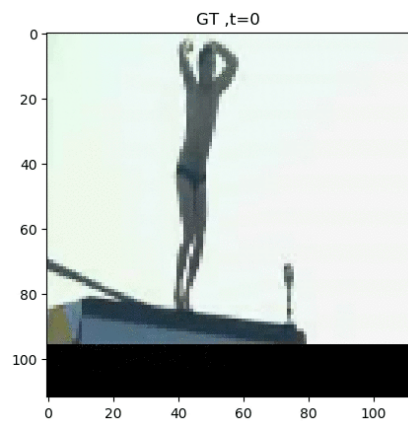




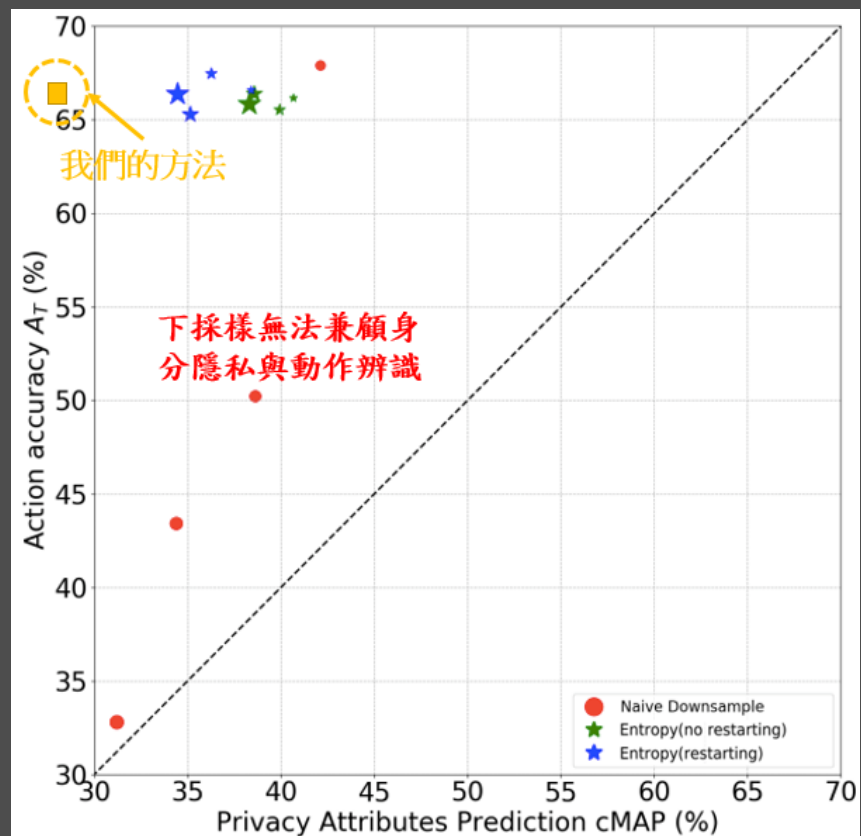
SBU Kinect Interaction Dataset

Result





UCF101 Dataset Result



41

ADVERSARIAL EXAMPLE ATTACK

ADVERSARIAL EXAMPLES

■ Adversary Target

- Perturbs a sample x to $x + \Delta x$ to fool model θ .

■ Possible Causes

- Curse of dimensionality

- ✓ Adversary generates unseen images off the manifold

- ReLU activation function

- ✓ ReLU network is piecewise linear, so carefully designed small noise can aggregate to alter decision.

■ Action Items

- Detection and removal of adversarial examples
- Robust design of neural network.

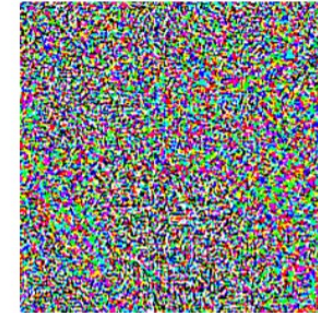


x

“panda”

57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

$=$



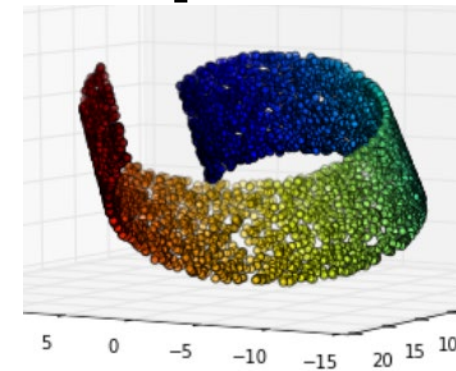
$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

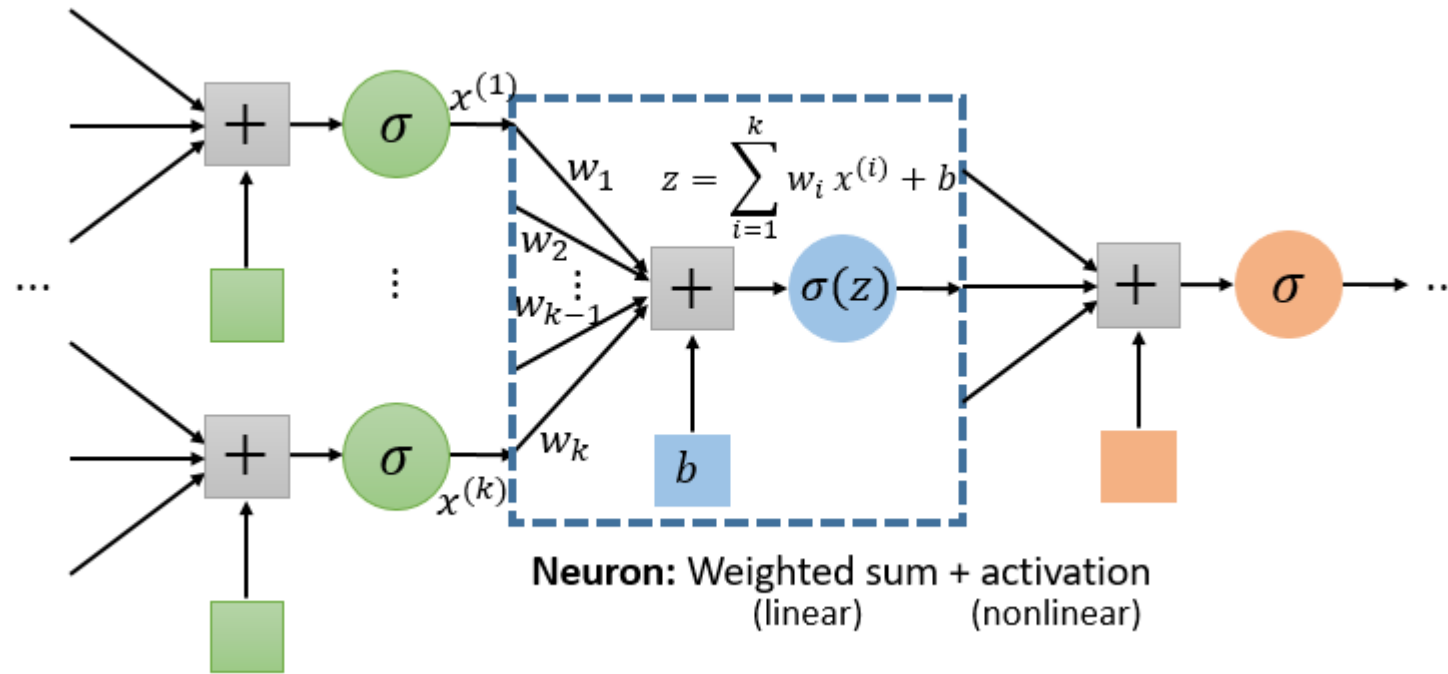
“gibbon”

99.3 % confidence

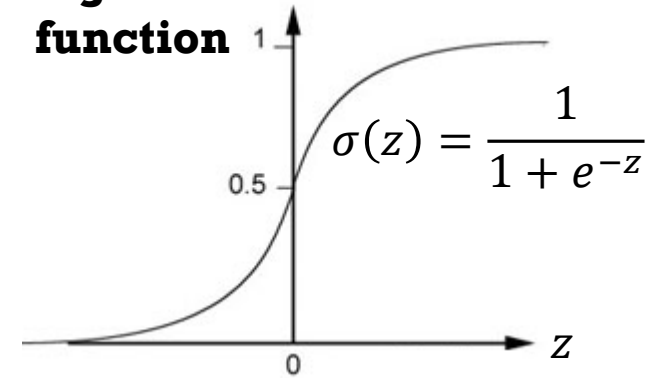
An example of manifold



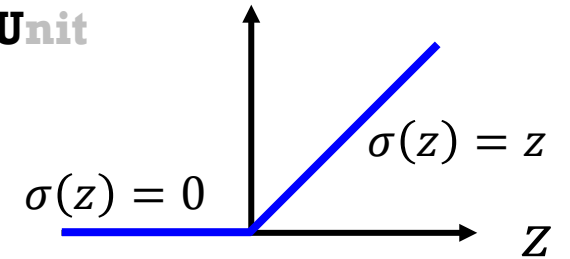
RECAP: NEURAL NETWORK



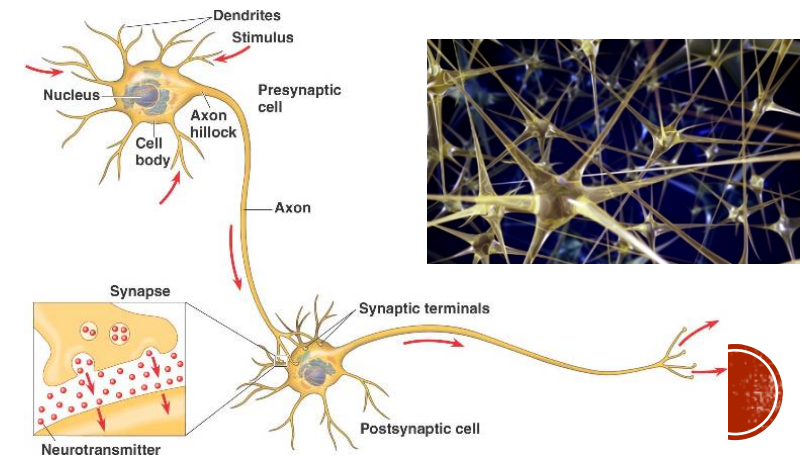
Sigmoid function



Rectified Linear Unit



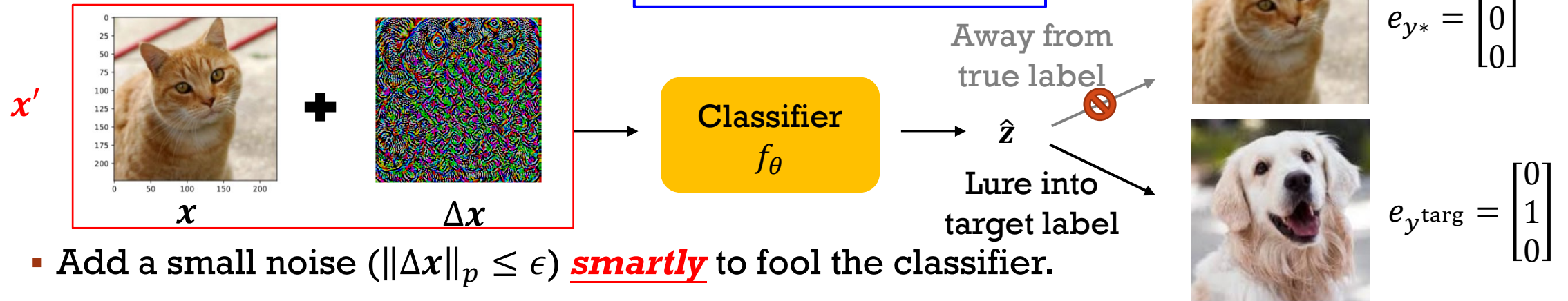
- Different connection leads to different network structures
- Network parameter θ : all the weights and biases in the “neurons”



ADVERSARIAL EXAMPLE ATTACK REMEDIES

▪ Attacker's goal

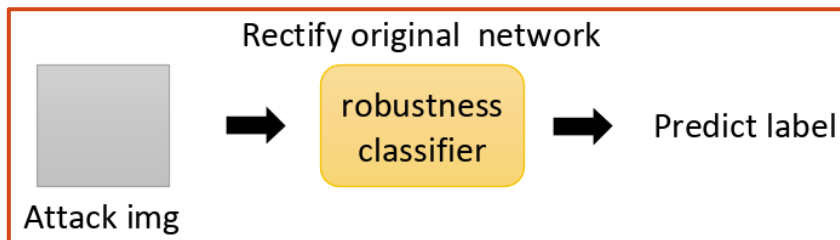
▪ Architecture



- Add a small noise ($\|\Delta x\|_p \leq \epsilon$) **smartly** to fool the classifier.

▪ Defender methods

- Enhance robustness against **specific** attack. (Cat-and-mouse game)
- Pre-processing the adversarial example (May also filter out the important feature)



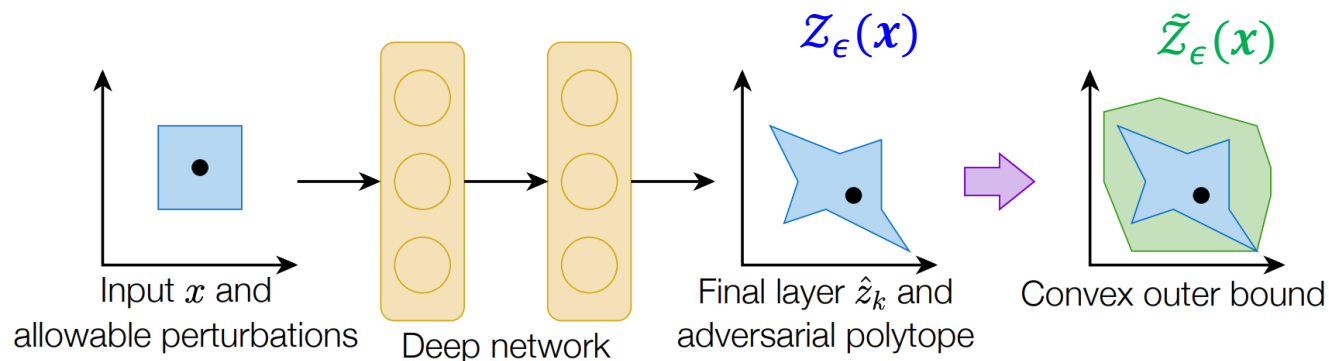
$$\min_{\theta} \sum_{i=1}^N L(f_{\theta}(x_i + \Delta x_i), y_i) \text{ where } \Delta x_i \text{ is by some specific attack.}$$

Problem: No theoretical guarantees on the worst case attacking scenario.



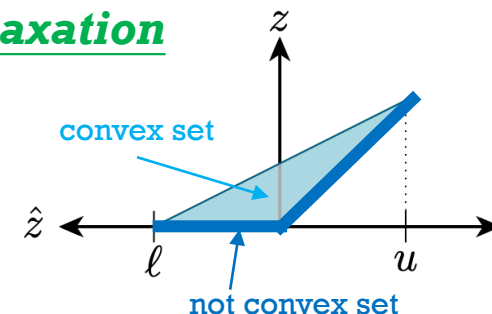
PROVABLE DEFENSE VIA CONVEX OUTER ADVERSARIAL POLYTOPE

- Base on deep fully-connected network (with ReLU activation)



Each network layer:

- Linear transform (convex constraint)
 - ReLU (NOT convex constraint)
- ➔ **Convex relaxation**



Key idea:

$$\mathcal{Z}_\epsilon(x) = \{f_\theta(x + \Delta x) : \|\Delta x\|_\infty \leq \epsilon\}$$

- Convex relaxation on cover $\mathcal{Z}_\epsilon(x)$ with convex polytope $\tilde{\mathcal{Z}}_\epsilon(x)$
- Any solution in dual problem is a lower bound to primal problem.

Difference in confidence value between true label and target label (by adversary)

$$\begin{aligned} \min & \mathbf{e}_{y^*}^T \hat{\mathbf{z}}_k - \mathbf{e}_{y^{\text{targ}}}^T \hat{\mathbf{z}}_k \\ \text{subject to } & \hat{\mathbf{z}}_k = f_\theta(x + \Delta x) \\ & \|\Delta x\|_\infty \leq \epsilon \end{aligned}$$

$\mathcal{Z}_\epsilon(x)$ (not convex)

Primal problem

$$\begin{aligned} \min & \mathbf{c}^T \hat{\mathbf{z}}_k \\ \text{subject to } & \hat{\mathbf{z}}_k \in \tilde{\mathcal{Z}}_\epsilon(x) \end{aligned}$$

(convex constraints)

Dual problem

$$\begin{aligned} \max & \theta(\xi, \nu, \mu, \tau, \lambda) \\ \text{subject to } & \text{convex constraints} \end{aligned}$$

(convex constraints)

tunable parameter $\mathcal{N}_\epsilon(x, y^{\text{targ}}, \alpha)$

In neural network form!

Eric Wong and J. Zico Kolter, "Provable defenses against adversarial examples via the convex outer adversarial polytope" ICML 2018

Tractable bound on worst case adversarial attack scenario.

PRIMAL PROBLEM AFTER CONVEX RELAXATION

minimize

$$\mathbf{c}^T \hat{\mathbf{z}}_k$$

$$\hat{\mathbf{z}}_{i+1} = \mathbf{W}_i \mathbf{z}_i + \mathbf{b}_i, \quad i = 1, \dots, k-1$$

weighted sum

$$\mathbf{z}_1 \leq \mathbf{x} + \boldsymbol{\epsilon}$$

$$\mathbf{z}_1 \geq \mathbf{x} - \boldsymbol{\epsilon}$$

$\|\cdot\|_\infty$ constraint on noise

subject to

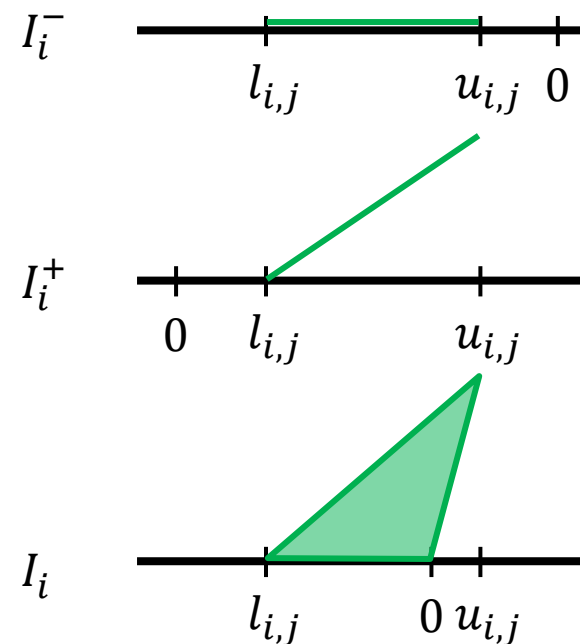
$$\left. \begin{array}{l} z_{i,j} = 0, \quad j \in \mathcal{I}_i^- \\ z_{i,j} = \hat{z}_{i,j}, \quad j \in \mathcal{I}_i^+ \\ z_{i,j} \geq 0 \\ z_{i,j} \geq \hat{z}_{i,j} \\ (u_{i,j} - l_{i,j})z_{i,j} - u_{i,j}\hat{z}_{i,j} \leq -u_{i,j}l_{i,j} \end{array} \right\} j \in \mathcal{I}_i$$

variables

$$\mathbf{z}_i \in \mathbb{R}^{d_i}, \quad i = 1, \dots, k-1$$

$$\hat{\mathbf{z}}_i \in \mathbb{R}^{d_i}, \quad i = 2, \dots, k$$

ReLU layer (with convex relaxation)



DERIVE THE DUAL PROBLEM

Introduce dual variables

$$\begin{aligned}\hat{z}_{i+1} = W_i z_i + b_i &\Rightarrow \nu_{i+1} \in \mathbb{R}^{|\hat{z}_{i+1}|} \\ z_1 \leq x + \epsilon &\Rightarrow \xi^+ \in \mathbb{R}^{|x|} \\ -z_1 \leq -x + \epsilon &\Rightarrow \xi^- \in \mathbb{R}^{|x|} \\ -z_{i,j} \leq 0 &\Rightarrow \mu_{i,j} \in \mathbb{R} \\ \hat{z}_{i,j} - z_{i,j} \leq 0 &\Rightarrow \tau_{i,j} \in \mathbb{R} \\ -u_{i,j} \hat{z}_{i,j} + (u_{i,j} - l_{i,j}) z_{i,j} \leq -u_{i,j} l_{i,j} &\Rightarrow \lambda_{i,j} \in \mathbb{R}\end{aligned}$$

Write down the Lagrangian

$$\begin{aligned}L(\mathbf{z}, \hat{\mathbf{z}}, \boldsymbol{\xi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\tau}, \boldsymbol{\lambda}) = & \mathbf{c}^T \hat{\mathbf{z}}_k + \sum_{i=1}^{k-1} \nu_{i+1}^T (\hat{\mathbf{z}}_{i+1} - (\mathbf{W}_i \mathbf{z}_i + \mathbf{b}_i)) + \xi_+^T (\mathbf{z}_1 - (\mathbf{x} + \epsilon)) + \xi_-^T (-\mathbf{z}_1 + (\mathbf{x} - \epsilon)) \\ & + \sum_{i=2}^{k-1} \left(\sum_{j \in \mathcal{I}_i^- \cup \mathcal{I}_i} \mu_{i,j} (-z_{i,j}) + \sum_{j \in \mathcal{I}_i^+ \cup \mathcal{I}_i} \tau_{i,j} (\hat{z}_{i,j} - z_{i,j}) + \sum_{j \in \mathcal{I}_i} \lambda_{i,j} ((u_{i,j} - l_{i,j}) z_{i,j} - u_{i,j} \hat{z}_{i,j} + u_{i,j} l_{i,j}) \right)\end{aligned}$$



$$\theta(\boldsymbol{\xi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\tau}, \boldsymbol{\lambda}) = \inf_{\mathbf{z}, \hat{\mathbf{z}}} L(\mathbf{z}, \hat{\mathbf{z}}, \boldsymbol{\xi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\tau}, \boldsymbol{\lambda})$$

Dual Problem

$$\begin{aligned}\text{maximize} \quad & \theta(\boldsymbol{\xi}, \boldsymbol{\nu}, \boldsymbol{\mu}, \boldsymbol{\tau}, \boldsymbol{\lambda}) = -\xi_+^T (\mathbf{x} + \epsilon) + \xi_-^T (\mathbf{x} - \epsilon) - \sum_{i=1}^{k-1} \nu_{i+1}^T \mathbf{b}_i + \sum_{i=2}^{k-1} \lambda_i^T (\mathbf{u}_i \odot \mathbf{l}_i) \\ \text{subject to} \quad & \nu_k = -\mathbf{c} \\ & \mathbf{W}_1^T \nu_2 = \xi_+ - \xi_- \\ & \left. \begin{aligned} \tau_i + \nu_i &= \mathbf{u}_i \odot \lambda_i \\ \mathbf{W}_i^T \nu_{i+1} + \tau_i + \mu_i &= (\mathbf{u}_i - \mathbf{l}_i) \odot \lambda_i \end{aligned} \right\} i = 2, \dots, k-1 \\ & \xi_+, \xi_- \in \mathbb{R}_{\geq 0}^{d_0} \\ & \nu_i \in \mathbb{R}^{d_i}, \quad i = 2, \dots, k \\ \text{variables} \quad & \left. \begin{aligned} \mu_{i,j} \in \mathbb{R}, \quad \tau_{i,j} = 0, \quad \lambda_{i,j} = 0, \quad j \in \mathcal{I}_i^- \\ \mu_{i,j} \geq 0, \quad \tau_{i,j} \geq 0, \quad \lambda_{i,j} \geq 0, \quad j \in \mathcal{I}_i \\ \mu_{i,j} = 0, \quad \tau_{i,j} \in \mathbb{R}, \quad \lambda_{i,j} = 0, \quad j \in \mathcal{I}_i^+ \end{aligned} \right\} i = 2, \dots, k-1\end{aligned}$$



DUAL PROBLEM IN NEURAL NETWORK FORM

maximize $J_\epsilon(\mathbf{x}, \boldsymbol{\nu}) = -\mathbf{x}^T \hat{\boldsymbol{\nu}}_1 - \epsilon |\hat{\boldsymbol{\nu}}_1| - \sum_{i=1}^{k-1} \mathbf{b}_i^T \boldsymbol{\nu}_{i+1} + \sum_{i=2}^{k-1} \sum_{j \in \mathcal{I}_i} \frac{u_{i,j} l_{i,j}}{u_{i,j} - l_{i,j}} [\nu_{i,j}]_+$

$\boldsymbol{\nu}_k = -\mathbf{c}$ input

$\hat{\nu}_{i,j} = (\mathbf{W}_i^T \boldsymbol{\nu}_{i+1})_j, i = 1, \dots, k-1$ linear transform

subject to $\nu_{i,j} = \begin{cases} 0 & , \text{ if } j \in \mathcal{I}_i^- \\ \frac{u_{i,j}}{u_{i,j} - l_{i,j}} [\hat{\nu}_{i,j}]_+ - \alpha_{i,j} [\hat{\nu}_{i,j}]_- & , \text{ if } j \in \mathcal{I}_i \\ \hat{\nu}_{i,j} & , \text{ if } j \in \mathcal{I}_i^+ \end{cases} \quad i = 2, \dots, k-1$

$0 \leq \alpha_{i,j} \leq 1, j \in \mathcal{I}_i$ Leaky ReLU

variables $\boldsymbol{\nu}, \hat{\boldsymbol{\nu}}$

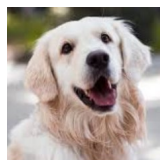


Can be written as neural network $\mathcal{N}_\epsilon(\mathbf{x}, y^{\text{targ}}, \alpha)$

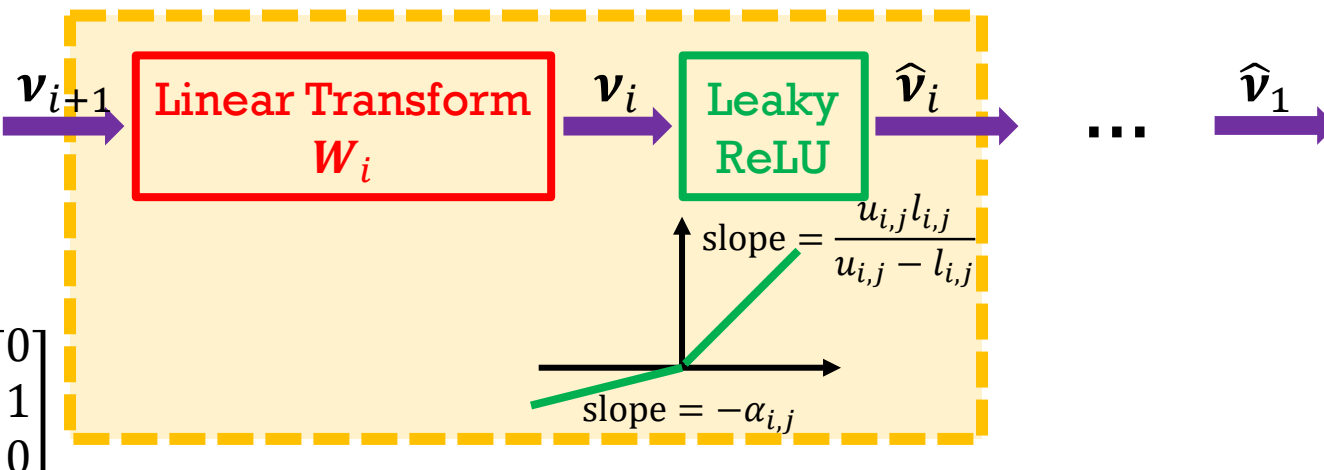
$\mathbf{v}_k = \mathbf{c} = \mathbf{e}_{y_*} - \mathbf{e}_{y^{\text{targ}}}$
(depends on y^{targ})



$\mathbf{e}_{y_*} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$



$\mathbf{e}_{y^{\text{targ}}} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$



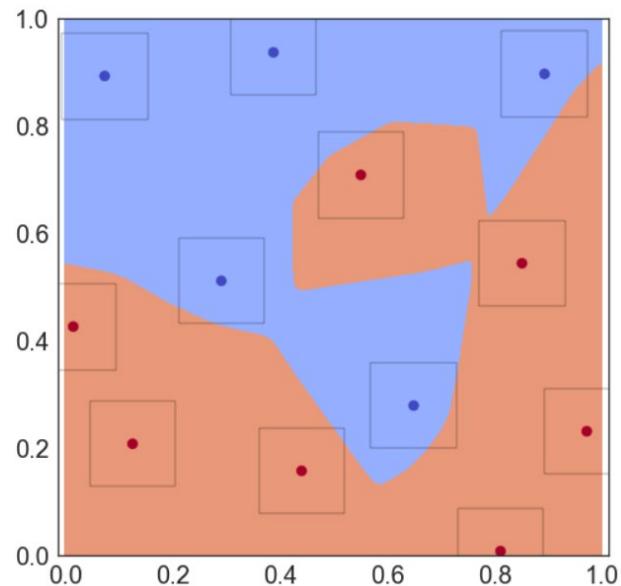
ROBUST LEARNING

$$\begin{aligned} \min \mathbf{e}_{y_*}^T \hat{\mathbf{z}}_k - \mathbf{e}_{y^{\text{targ}}}^T \hat{\mathbf{z}}_k \\ \text{subject to } \|\Delta \mathbf{x}\|_\infty \leq \epsilon \end{aligned} \geq \mathcal{N}_\epsilon(\mathbf{x}, y^{\text{targ}}, \alpha)$$

If $\mathcal{N}_\epsilon(\mathbf{x}, \alpha)$ is positive for all y^{targ} , then adversary cannot fool the classifier

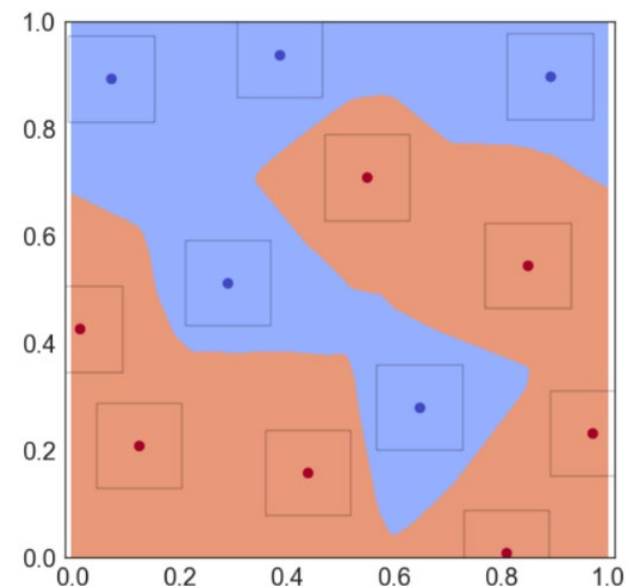
Standard training

$$\min_{\theta} \sum_{i=1}^N L(f_{\theta}(x_i), y_i)$$



Robust training

$$\max_{\theta} \min_{y^{\text{targ}}} \sum_{i=1}^N \mathcal{N}_\epsilon(\mathbf{x}, y^{\text{targ}}, \alpha)$$



TAKE AWAY MESSAGES

- MLaaS raises security issues such as
 - Model inversion attack: Infer training data from model.
 - Membership inference attack: Infer membership from model.
 - Adversarial example attack: Fool the classifier with unperceivable noises.
and many others...
- Methods for secure machine learning:
 - Differential Privacy: Adding noise to machine learning models.
 - Homomorphic Encryption: Cryptographic approach. Secure but costly.
 - Compressive Privacy: Nonlinear lossy compression
 - ✓ Preserve sufficient information for machine learning service.
 - ✓ Difficult for intruder to reconstruct from compressed data.
 - ✓ CPGAN defends the reconstruction attack while achieving satisfactory utility performance.

