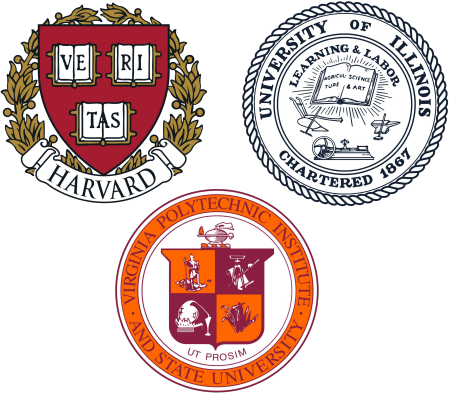


# Improving Robustness to Model Inversion Attacks via Mutual Information Regularization

Tianhao Wang<sup>1</sup>, Yuheng Zhang<sup>2</sup>, Ruoxi Jia<sup>3</sup>

<sup>1</sup>Harvard University <sup>2</sup>University of Illinois Urbana-Champaign <sup>3</sup>Virginia Tech

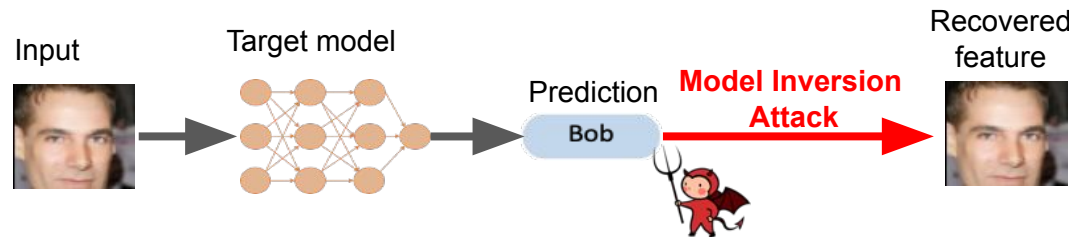
TL; DR: We formally analyze model inversion attack and propose the Mutual Information Regularization based Defense (MID).



## Motivation

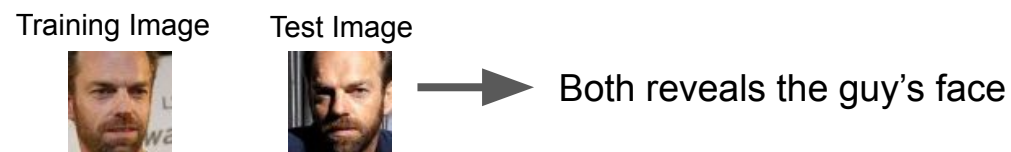
- Existing defense mechanisms against model inversion attack rely on model-specific heuristics or noise injection.
- Existing defense mechanisms significantly hinder model performance.
- We need to design a defense mechanism that is **applicable to a variety of models** and achieves **better utility-privacy tradeoff**.

## Model Inversion Attack



## Defense Goal

Both the recovery of training images and test images would incur privacy loss to the target identity. We need to design an algorithm to **protect the training data distribution**, instead of just training data set.



## MID: Mutual Information Regularization based Defense

**Intuition:** If the output distribution is independent from input distribution, the attacker cannot learn anything about X's distribution.

**Method:** Regularize the loss function by the mutual information between model's input and output distribution.

$$\min_{f \in \mathcal{H}} E_{(x,y) \sim p_{X,Y}(x,y)} [\mathcal{L}(y, f(x))] + \lambda \mathcal{I}(X, \hat{Y})$$

Original Loss Function

Regularizer Coefficient

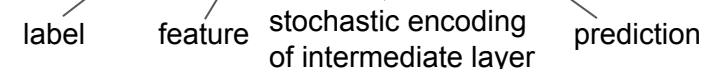
mutual information between input and prediction

$$\mathcal{I}(X, \hat{Y}) = \int_{x,y} p_{X,Y}(x,y) \log \left( \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} \right) dx dy$$

## Instantiation of MID

- Linear regression: Taylor-expansion based approximation
- Decision tree: modify ID3
- Deep Neural Networks: information bottleneck technique

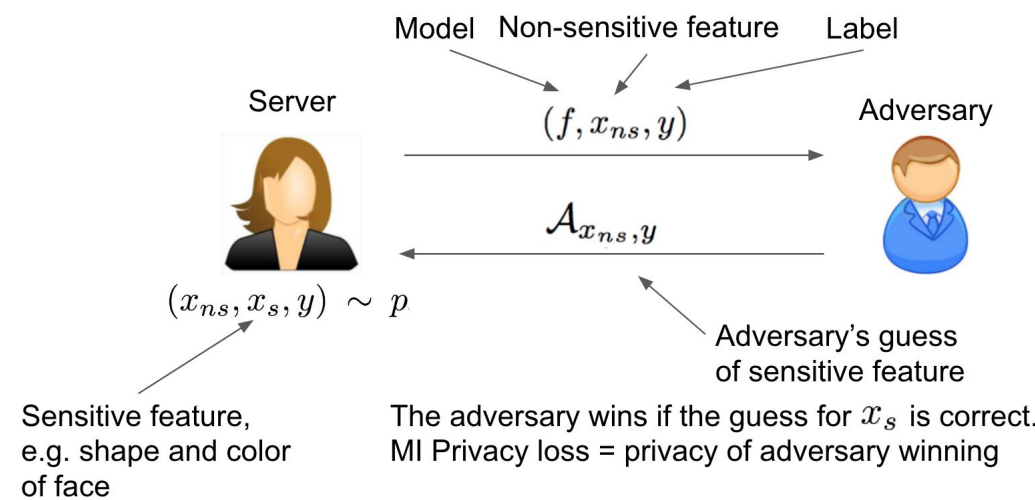
For example, we can regard the neural network as a Markov chain:  $Y - X - Z - \hat{Y}$



By Data Processing Inequality, we have  $\mathcal{I}(X, \hat{Y}) \leq \mathcal{I}(X, Z)$  and we can obtain a new training loss  $\min_{\theta} -\mathcal{I}(Z; Y) + \lambda \mathcal{I}(Z; X)$  which boils down to the **classic information bottleneck**.

## Formalizing Model Inversion Attacks

We present a methodology for formalizing model inversion attacks. Unlike previous works that only capture the privacy loss of members in the training set, this is the first attempt of modeling privacy loss of members **in the population**.



## Charactering MI Privacy Loss for Differentially Private Models

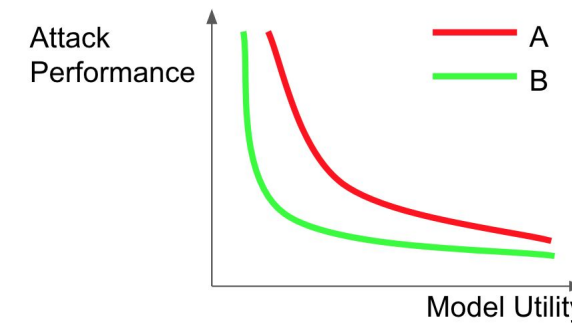
Main Result: when the learning algorithm is  $(\epsilon, \delta)$ -differentially private, the MI privacy loss is **tightly** upper bounded by

$$\frac{e^{n\epsilon} - 1}{e^{n\epsilon} + 1} + \frac{2(e^{n\epsilon} - 1)}{(e^{n\epsilon} + 1)(e^{\epsilon} - 1)} \delta$$

where n is the size of training set. To make bound small, the privacy budget  $\epsilon$  needs to be set as  $\mathbf{o(1 / \#training\ data)}$ !

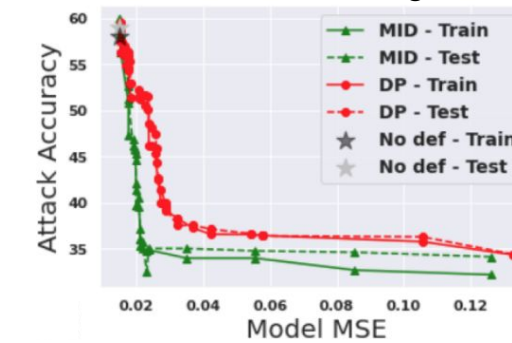
## Evaluation Metric

- Defense mechanisms are evaluated in terms of **privacy-utility tradeoff**.
- In the illustration below, the green line is more preferable as it is **more robust against the attack at any fixed model utility**.

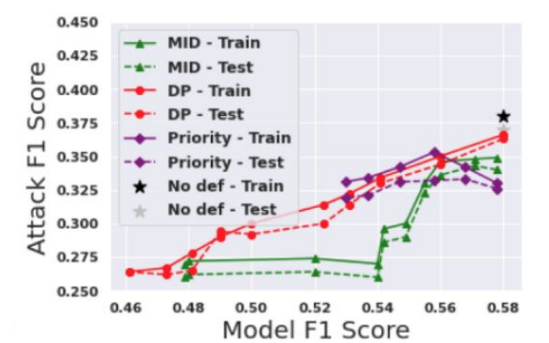


## Evaluation on Defending against Various MI Attacks

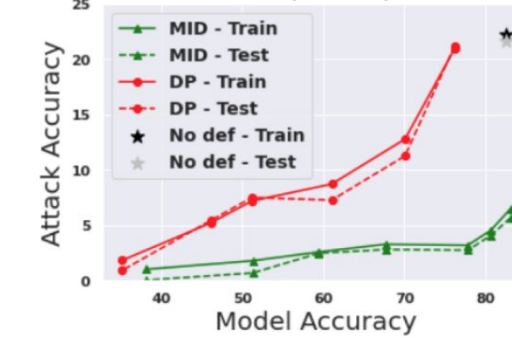
### MAP on linear regression



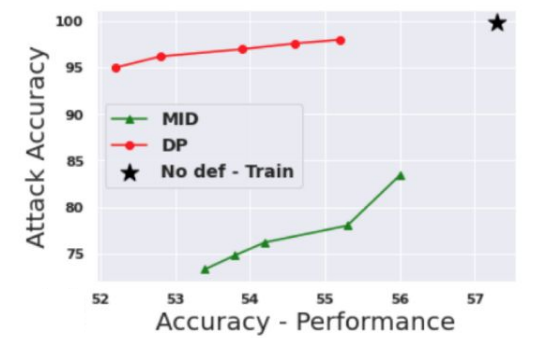
### MAP on decision tree



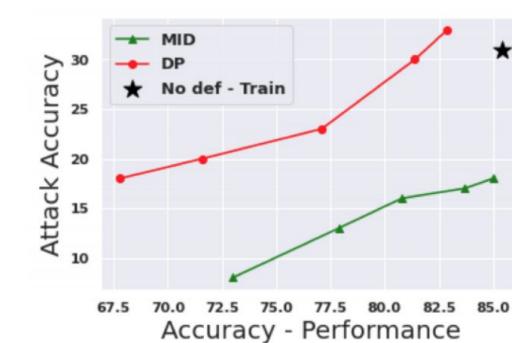
### Knowledge Alignment



### Update-Leaks



### GMI



### White-box MAP

