

# Machine Learning in Cybersecurity

- Adversarial Perturbations & Game Theory
- Membership Inference Attacks

Prof. Dr. Mario Fritz | 17.12.2020

- Act Responsible
- Differentiate between “what can be done” and “what should be done”

- Security & Privacy
  - S&P (Oakland), CCS, Usenix Security, NDSS
- ML
  - Neurips (Dec 8<sup>th</sup>), ICML, KDD
    - Workshop on Security in Machine Learning
      - <https://secml2018.github.io>
    - Workshop on Privacy Preserving Machine Learning
      - <https://ppml-workshop.github.io/ppml/>
  - AI
    - AAAI, IJCAI
  - CSRankings.org

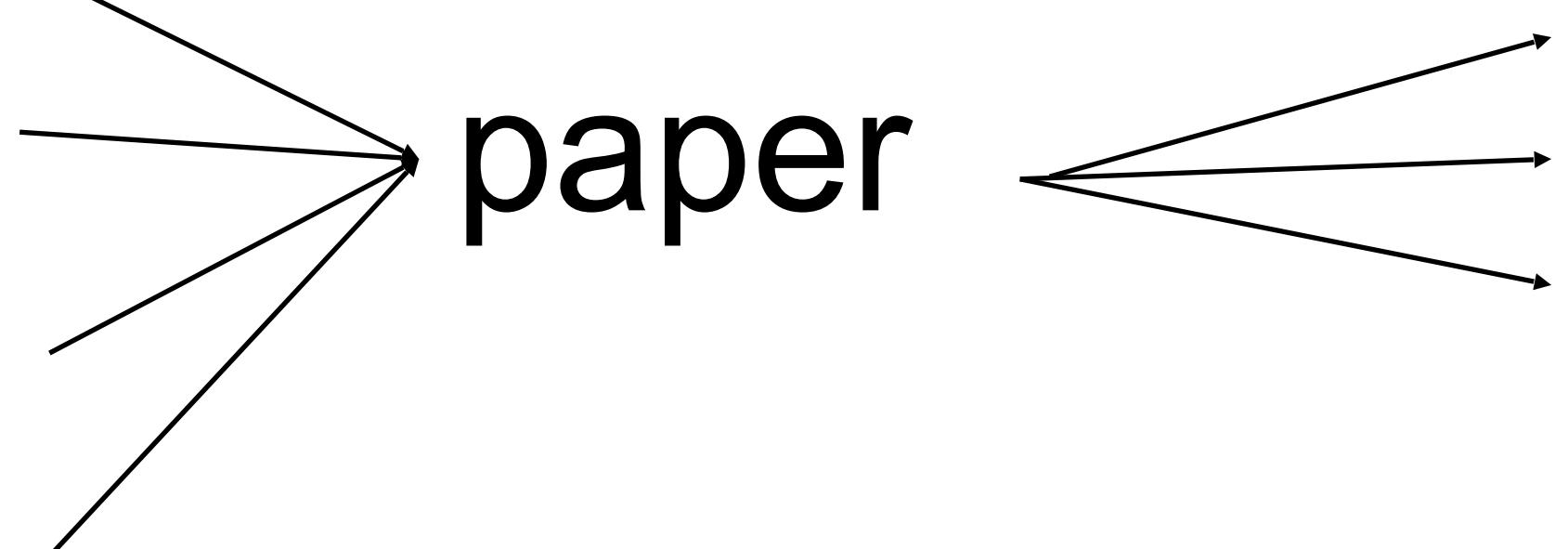
## Past

Paper Bibliography

- REFERENCES**
- [1] "Cifar," <https://www.cs.toronto.edu/~kriz/cifar.html>. [2](#) [3](#)
  - [2] "Labeled Faces in the Wild: A Database for Studying Face Recognition in Uncontrolled Environments," <http://vis-www.cs.umass.edu/lfw/>. [3](#)
  - [3] "Lasagne," <http://lasagne.readthedocs.io>.
  - [4] "Mnist," [3](#)
  - [5] "UCI Machine Learning Repository: Data Sets," <http://archive.ics.uci.edu/ml/>.
  - [6] M. Balazs, P. Berrang, M. Biagi, R. Ellis, C. Hermann, M. Humbert, and J. Lehmann, "Identifying Personal DNA Methylation Profiles by Genome Inference," in *Proceedings of the 3rd IEEE Symposium on Security and Privacy (S&P)*, IEEE, 2017, pp. 387-398. [12](#)
  - [7] M. Balazs, P. Berrang, M. Humbert, and P. Munozman, "Membership Privacy in MicroRNA-based Studies," in *Proceedings of the 24th ACM SIGSAC Conference on Computer and Communications Security (CCS)*, ACM, 2016, pp. 313-330. [1](#) [12](#)
  - [8] M. Balazs, M. Humbert, J. Peng, and Y. Zhang, "Walk2friends: Inferring Social Links from Mobility Profiles," in *Proceedings of the 24th ACM SIGSAC Conference on Computer and Communications Security (CCS)*, ACM, 2017, pp. 1943-1957. [13](#)
  - [9] M. Barai, F. Fallu, R. Lazzeretti, A.-R. Salioglu, and T. Schneider, "Privacy-Preserving ECG Classification With Branching Programs and Neural Networks," *IEEE Transactions on Information Forensics and Security*, vol. 6, no. 2, pp. 452-466, 2011. [12](#)
  - [10] K. Bennett, I. Ivanov, B. Krauter, A. Marrodot, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Smith, "Practical Secure Aggregation for Privacy-Preserving Machine Learning," in *Proceedings of the 24th ACM SIGSAC Conference on Computer and Communications Security (CCS)*, ACM, 2017, pp. 1175-1191. [13](#)
  - [11] R. Boni, R. A. Popa, S. Tu, and S. Goldwasser, "Machine Learning Classification over Encrypted Data," in *Proceedings of the 22th Network and Distributed System Security Symposium (NDSS)*, 2012. [13](#)
  - [12] N. Carlini and D. Wagner, "Towards Evaluating the Robustness of Neural Networks," in *Proceedings of the 35th IEEE Symposium on Security and Privacy*, 2014. [13](#)
  - [13] B. Chen, Y. Liu, and X. Wang, "A Practical and Efficient Privacy-Preserving Machine Learning," *IEEE Trans. Dependable Secur. Comput.*, vol. 13, no. 1, pp. 10-21, 2016. [13](#)
  - [14] K. Lang, "Newsreader: Learning to filter news," in *Proceedings of the 12th International Conference on Machine Learning (ICML)*, [2](#) [3](#)
  - [15] B. Li and Y. Vorobeychik, "Scalable Optimization of Randomized Operational Decisions in Adversarial Classification Settings," in *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS)*, PMLR, 2015, pp. 599-607. [13](#)
  - [16] J. Lin, M. Jaiswal, Y. Lu, and N. Asokan, "Oblivious Neural Network Predictions via Min-ONN Transformations," in *Proceedings of the 34th ACM SIGSAC Conference on Computer and Communications Security (CCS)*, ACM, 2017, pp. 619-631. [13](#)
  - [17] Y. Long, V. Bindnagel, and C. A. Gunter, "Towards Measuring Membership Privacy," *arXiv:1712.09446*, 2017. [12](#)
  - [18] Y. Long, V. Bindnagel, L. Wang, D. Bai, X. Wang, H. Tang, C. A. Gunter, and K. Chen, "Understanding Membership Inference on Well-Generalized Learning Models," *arXiv:1902.04989*, 2019. [13](#)
  - [19] L. Malis, C. Song, E. D. Cristofaro, and V. Shmatkov, "Inference Attacks Against Collaborative Learning," *arXiv:1805.09049*, 2018. [13](#)
  - [20] P. Mohassel and Y. Zhang, "SecureML: A System for Scalable Privacy-Preserving Machine Learning," in *Proceedings of the 35th Symposium on Security and Privacy (S&P)*, IEEE, 2017, pp. 19-38. [13](#)
  - [21] S. J. Oh, M. Augustin, B. Schiele, and M. Fritz, "Towards Reverse-Engineering Black-Box Neural Networks," in *Proceedings of the 2018 International Conference on Learning Representations (ICLR)*, 2018. [1](#)
  - [22] S. J. Oh, M. Fein, and B. Schiele, "Adversarial Image Perturbation for Privacy Protection - A Game Theory Perspective," in *Proceedings of the 2017 IEEE International Conference on Computer Vision (ICCV)*, IEEE, 2017, pp. 1491-1500. [13](#)
  - [23] N. Papernot, P. D. McDaniel, T. J. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical Black-Box Attacks Against Machine Learning," in *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security (ASIA'CS)*, ACM, 2017, pp. 503-519. [13](#)
  - [24] N. Papernot, P. D. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The Limitations of Deep Learning in Adversarial Settings," in *Proceedings of the 1st IEEE European Symposium on Security and Privacy*, 2016. [13](#)

## Present

paper



## ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models

Ahmed Salem\*, Yang Zhang\*, Mathias Humbert†,  
Mario Fritz†, Michael Backes†  
\*CISPA, Saarland Informatics Campus  
†Swiss Data Science Center, ETH Zurich and EPFL  
†CISPA Helmholtz Center I.G., Saarland Informatics Campus

CS.CR] 4 Jun 2018

**Abstract**—Machine learning (ML) has become a core component of many real-world applications and training data is a key factor that drives current progress. This huge success has led Internet companies to deploy machine learning as a service (MLaaS). Recently, the first membership inference attack has shown that extraction of information on the training set is possible in such MLaaS settings, which has severe security and privacy implications.

However, the early demonstrations of the feasibility of such attacks have many assumptions on the adversary such as using multiple so-called shadow models, knowledge of the target model structure and having a dataset from the same distribution as the target model's training data. We relax all 3 key assumptions, disease. Previously, membership inference has been successfully conducted in many other domains, such as biomedical data [7] and mobility data [35].

Shokri et al. [36] present the first membership inference attack against machine learning models. The general idea behind this attack is to use multiple machine learning models (one for each prediction class), referred to as *attack models*, to make membership inference over the *target model*'s output, i.e., posterior probabilities. Given that the target model is a black-box API, Shokri et al. propose to construct multiple *shadow models* to mimic the target model's behavior and

## Future

e.g. google scholar

### MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models

[A Salem, Y Zhang, M Humbert, P Berrang... - arXiv preprint arXiv ..., 2018 - arxiv.org](#)

Machine learning (ML) has become a core component of many real-world applications and training data is a key factor that drives current progress. This huge success has led Internet companies to deploy machine learning as a service (MLaaS). Recently, the first membership ...

☆ 59 Cited by 43 Related articles All 8 versions

Privacy-preserving machine learning through data obfuscation

[T Zhang, Z He, B Li - arXiv preprint arXiv:1807.01960, 2018 - arxiv.org](#)

As machine learning becomes a practice and commodity, numerous cloud-based services and frameworks are provided to help customers develop and deploy machine learning applications. While it is prevalent to outsource model training and serving tasks in the cloud ...

☆ 59 Cited by 10 Related articles All 5 versions

Micapsule: Guarded offline deployment of machine learning as a service

[L Horiot, Y Zhang, A Bawali... - arXiv preprint arXiv ..., 2018 - arxiv.org](#)

With the widespread use of machine learning (ML) techniques, ML as a service has become increasingly popular. In this setting, an ML model resides on a server and users can query the model with their data via an API. However, if the user's input is sensitive, sending it to the ...

☆ 59 Cited by 15 Related articles All 5 versions

Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning

[M Naor, R Shokri, A Hohenberger - 2019 IEEE Symposium on ..., 2019 - ieeexplore.ieee.org](#)

Deep neural networks are susceptible to various inference attacks as they remember information about their training data. We design white-box inference attacks to perform a comprehensive privacy analysis of deep learning models. We measure the privacy leakage ...

☆ 59 Cited by 6 Related articles

Knockoff nets: Stealing functionality of black-box models

[T Orenstein, B Schiele, M Fritz - Proceedings of the IEEE ..., 2019 - openaccess.thecvf.com](#)

Abstract Machine Learning (ML) models are increasingly deployed in the wild to perform a wide range of tasks. In this work, we ask to what extent can an adversary steal functionality of such "victim" models based solely on blackbox interactions: image in, predictions out. In ...

☆ 59 Cited by 13 Related articles All 7 versions

Attributing fake images to GANs: Analyzing fingerprints in generated images

[N Yu, L Davis, M Fritz - arXiv preprint arXiv:1811.08180, 2018 - arxiv.org](#)

Research in computer graphics has been in pursuit of realistic image generation for a long time. Recent advances in machine learning with deep generative models have shown increasing success of closing the realism gap by using data-driven and learned ...

☆ 59 Cited by 8 Related articles All 2 versions

High-fidelity extraction of neural network models

[M Jagade, N Canini, D Berthier, A Kurakin... - arXiv preprint arXiv ..., 2019 - arxiv.org](#)

Model extraction allows an adversary to steal a copy of a remotely deployed machine learning model given access to its predictions. Adversaries are motivated to mount such attacks for a variety of reasons, ranging from reducing their computational costs, to ...

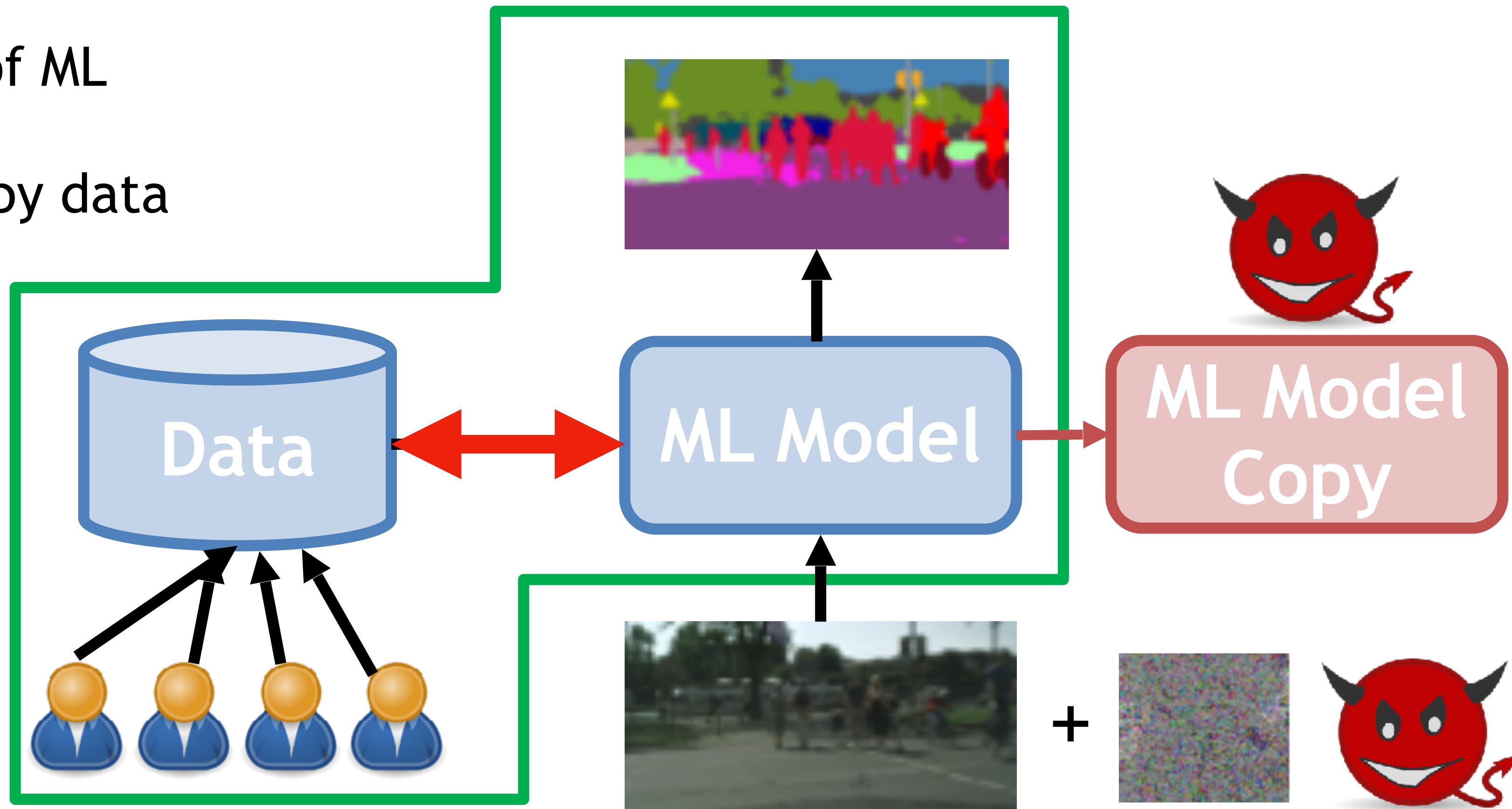
☆ 59 Cited by 3 All 2 versions

# Attacks on ML and Privacy

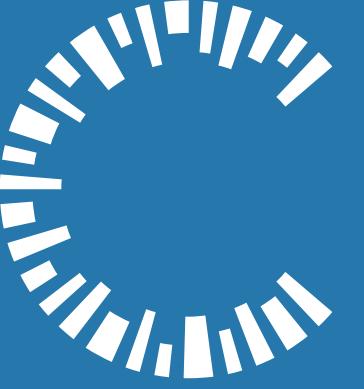
- Widespread deployment of ML
  - Future industry is fueled by data
  - How to make  
Machine Learning  
privacy compliant and  
secure?



- Membership Inference
  - Linkability Attack



# Adversarial Perturbations



# Adversarial Perturbations: Game Theoretic Analysis of Attacks and Defenses

- adversarial perturbation
- zero sum games

# Attacking for the good

- Adversarial perturbations as privacy protection
- Pros
  - Imperceptible
  - Effective (traditional obfuscation techniques are not)
- Cons
  - Model specific
- Game Theoretic analysis can help in attack/defense/situation



Fully visible

Gaussian blur

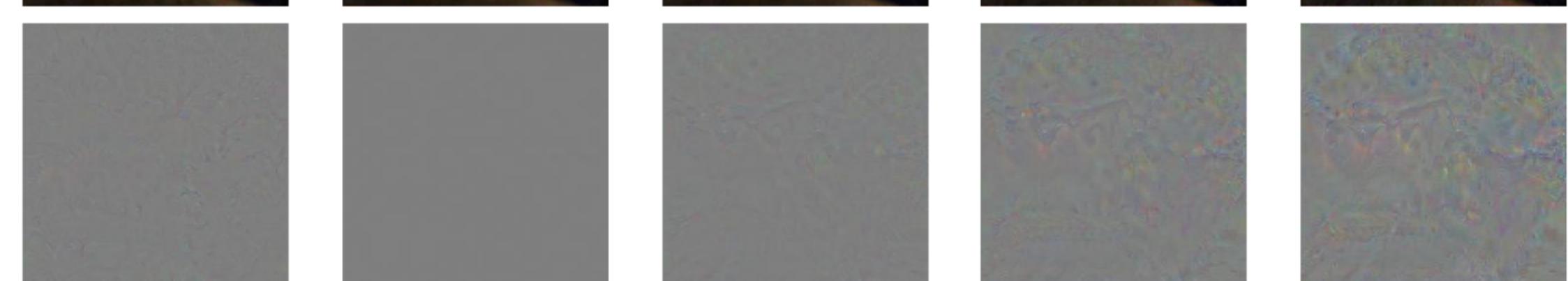
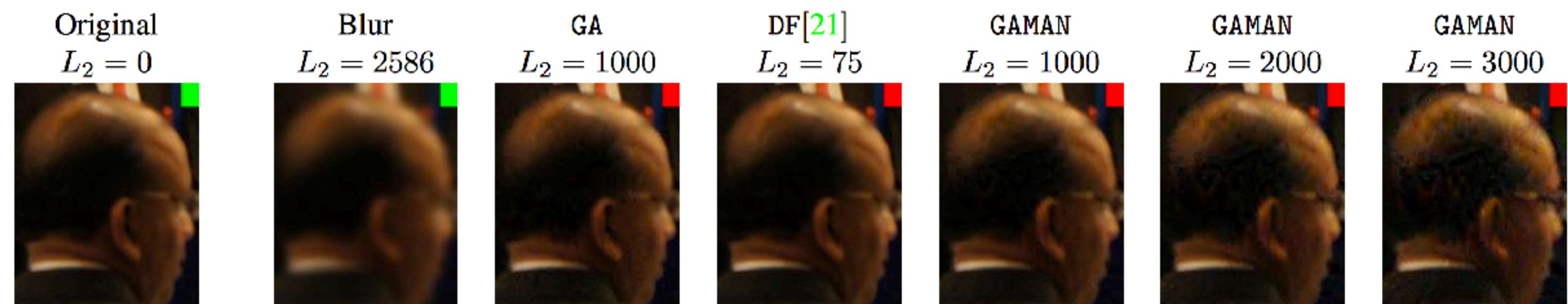
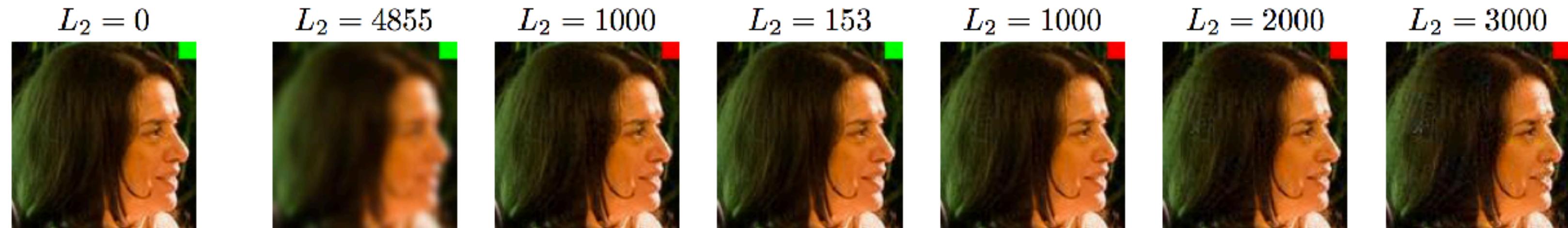


Black fill-in

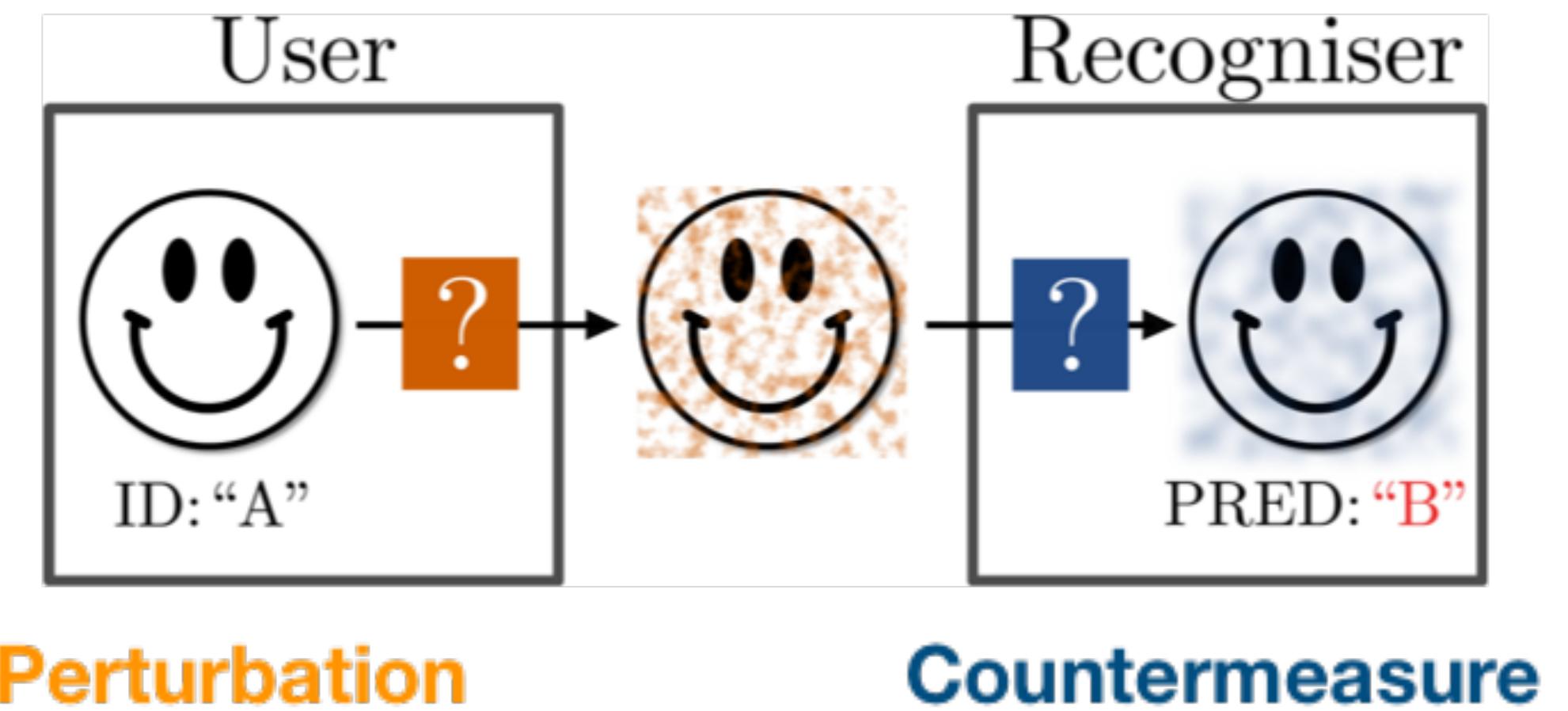


White fill-in

# Adversarial Image Perturbations

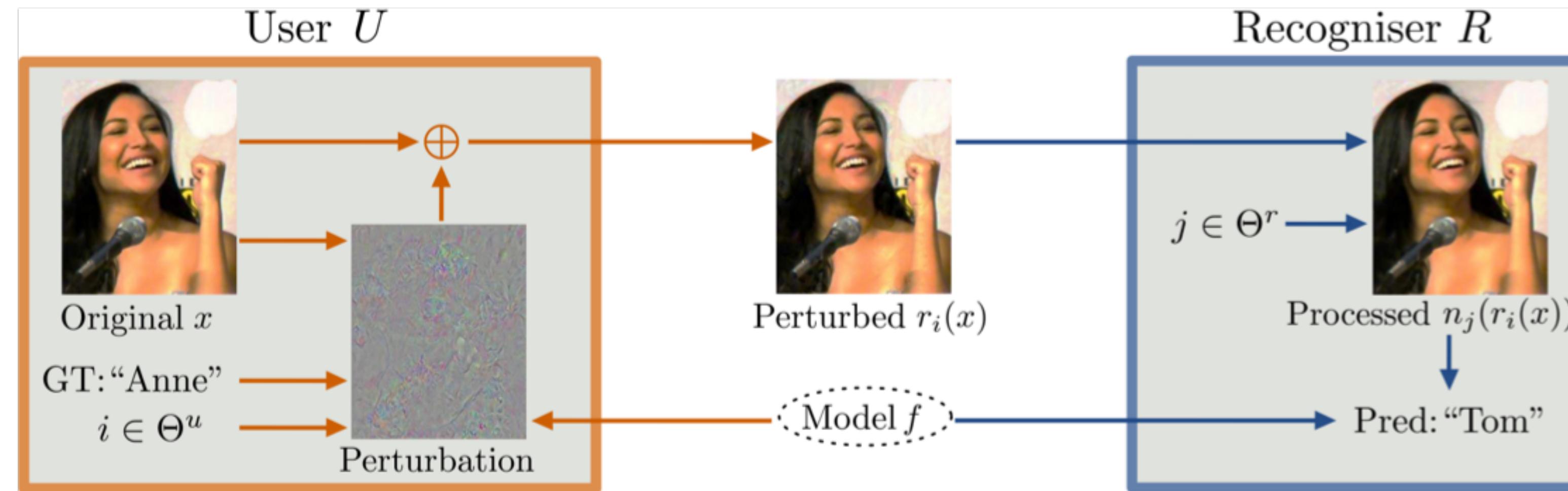


# Motivation: Game Theoretic Analysis



- Perturbation techniques and counter measures rapidly evolving
- Counter measures: blur, crop, jpg encoding
- Goal:
  - Analysis of action space
  - Upper bound on recognition rate

# Model: Zero-Sum Game



- Two players (User, Recogniser), antagonistic goals.
- Finite set of actions
- Strategy is sampled following an unknown distribution.
- Principled formulation of opponent-independent guarantee.
- Assumption: Known, fixed recognition model (white box)

# Model: Zero-Sum Game

- User (U) strategy

Defense-specific AIPs, Diverse set of AIP algorithms, AIP generation module weights (infinite space), ...

$$i \in \Theta^u$$

AIP 1

AIP 2

AIP 3

- Recogniser (R) strategy

Defense mechanism, Neural net architectures, Neural net weights (infinite space), ...

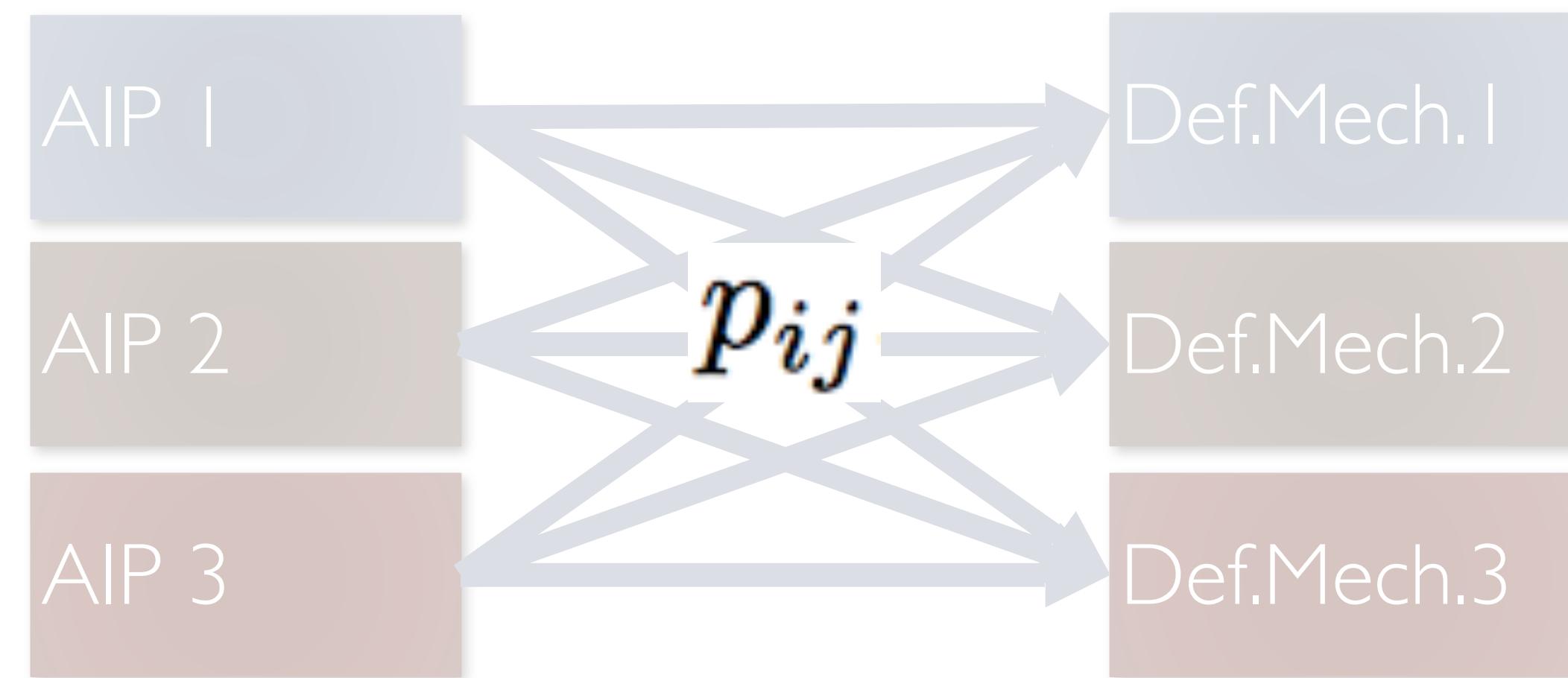
$$j \in \Theta^r$$

Def.Mech.1

Def.Mech.2

Def.Mech.3

- Reward for R,  $p_{ij}$ , is the recognition rate when U plays i and R plays j:

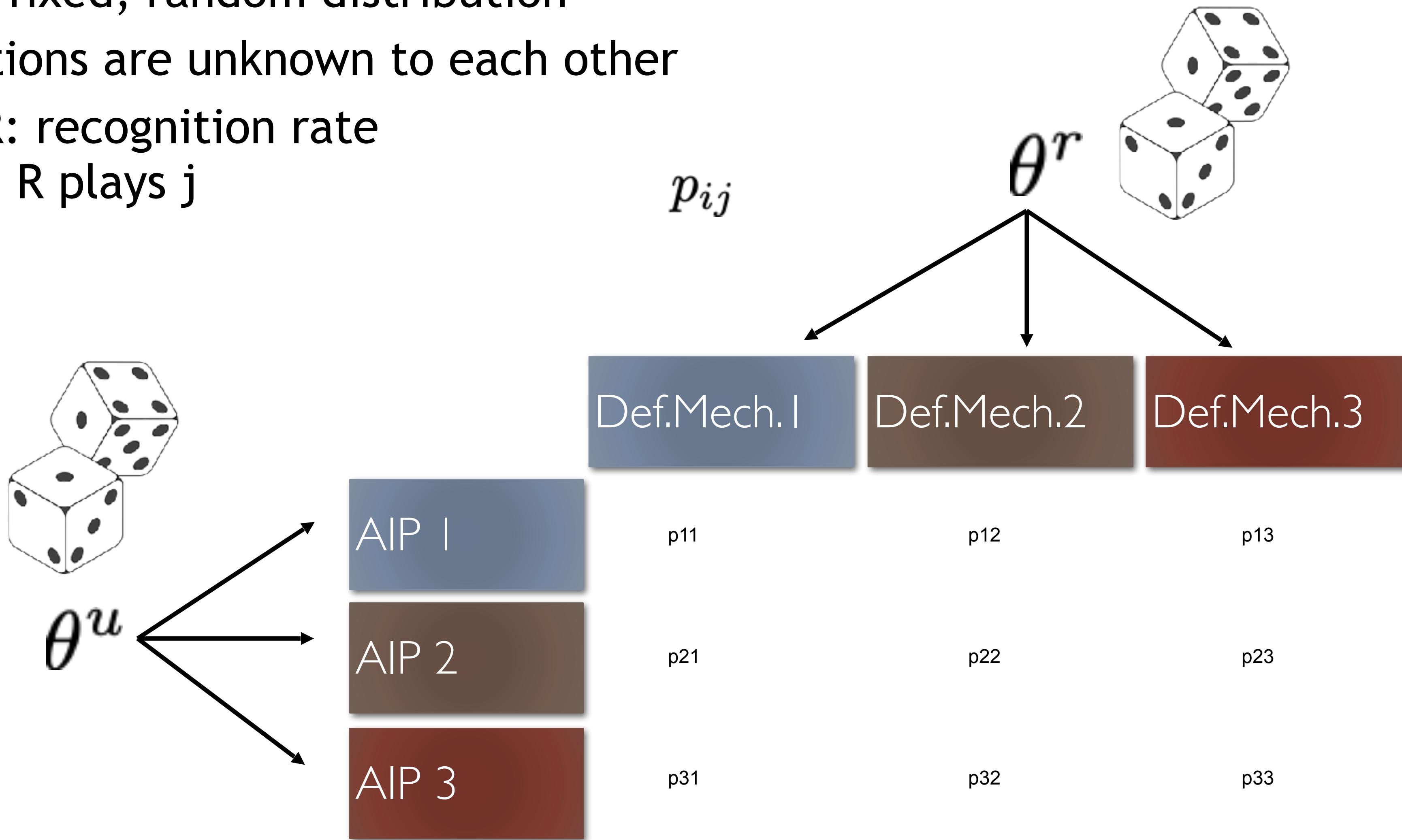


- Reward for U is the mis-recognition rate (ZSG):

$$1 - p_{ij}$$

# Model: Zero-Sum Game

- Strategy is a fixed, random distribution
- The distributions are unknown to each other
- Reward for R: recognition rate  
U plays i and R plays j



- R's expected reward:

$$\sum_{i,j} \theta_i^u \theta_j^r p_{ij}$$

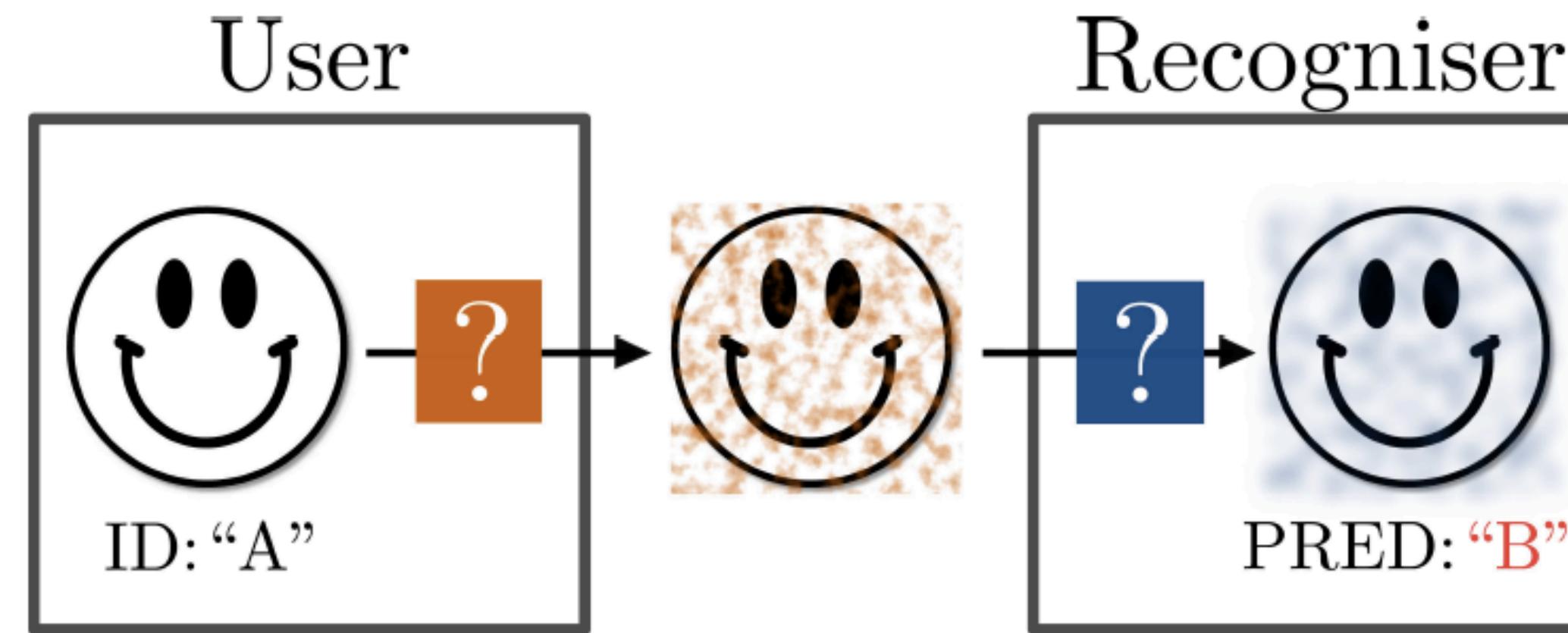
- U's “optimal” strategy (also value of the game v):

$$\arg \min_{\theta^u} \max_{\theta^r} \sum_{i,j} \theta_i^u \theta_j^r p_{ij} \text{ s.t. } \theta^u, \theta^r \text{ are distributions.}$$

- If U plays optimal strategy, recognition rates is upper bounded by v
- Bound independent of R's strategy

- Attack:
  - Adversarial perturbations make the recognizer fail
- Defense against attack:
  - Adversarial training: train with perturbed examples
- Attack against defense:
  - Adversarial perturbations of the adversarially trained classifier

# Case study: Action space



- **User actions:**  
AIP variants:  
 $\{\text{GAMAN}, \text{GAMAN}_T, \text{GAMAN}_N, \text{GAMAN}_B, \text{GAMAN}_C, \text{GAMAN}_{TNBC}\}$ .
- **Recogniser actions:**  
Defense measures:  
Translation, Noise, Blur, Crop.  
[Graese *et al.* ICMLA'16]  
 $\{f, f_T, f_N, f_B, f_C, f_{TNBC}\}$ .

# Results

Recogniser  $\Theta^r$

User $\Theta^u$	Proc	T	N	B	C	TNBC
GAMAN	4.0	6.6	15.0	22.2	16.7	9.9
/T	2.5	2.3	11.6	18.5	7.2	4.9
/N	5.8	7.6	4.6	23.6	16.6	9.1
/B	0.4	0.8	8.6	5.8	3.1	1.4
/C	2.6	2.2	11.8	18.1	3.4	4.3
/TNBC	0.7	0.9	5.2	9.5	3.2	2.0

recogniser  
pay-off  
table

- Reference: Clean image: 91.1%, No-image: 0.8%)
- Best deterministic: 8.6% (GAMAN/B: 100%) (fN: 100%)
- Best randomised: guarantees < 7.4% expected recognition rate
  - Optimal U strategy: (GAMAN/B:61%, GAMAN/TNBC:39%).
  - Optimal R strategy: (fN:52%, fN:48%).

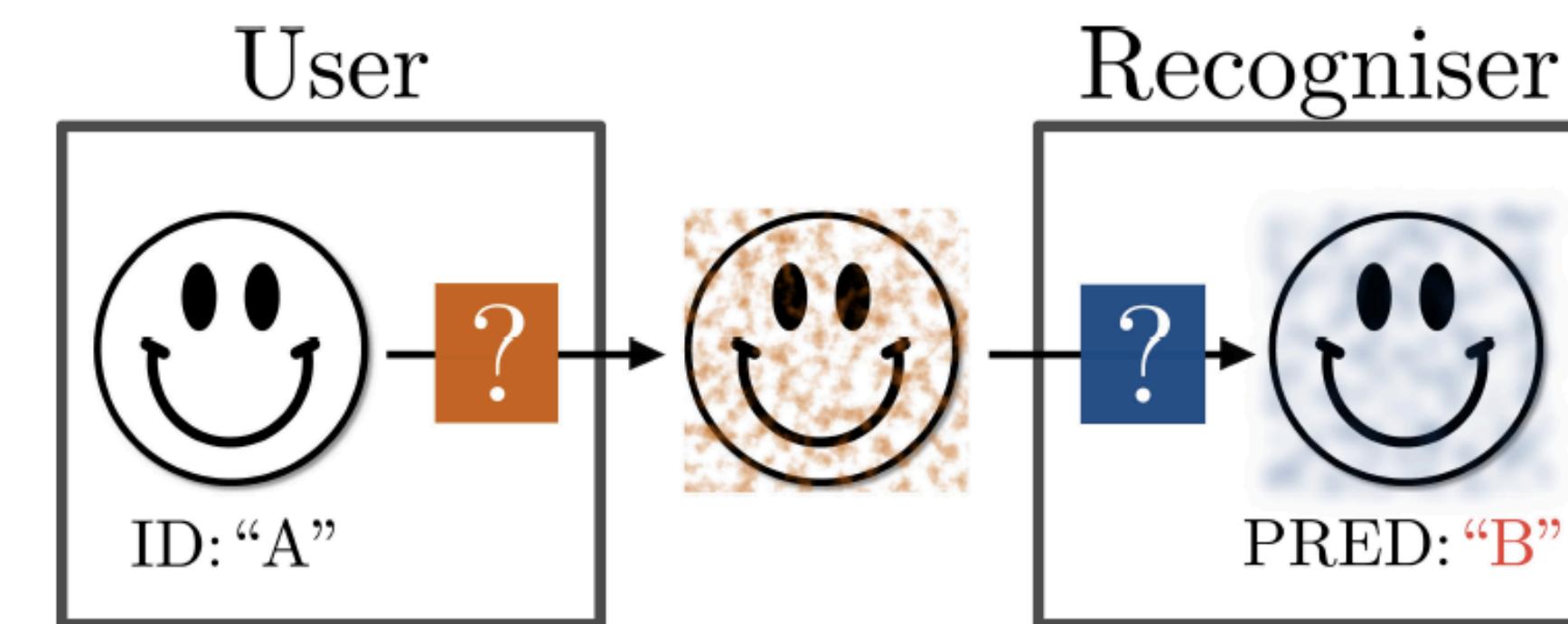
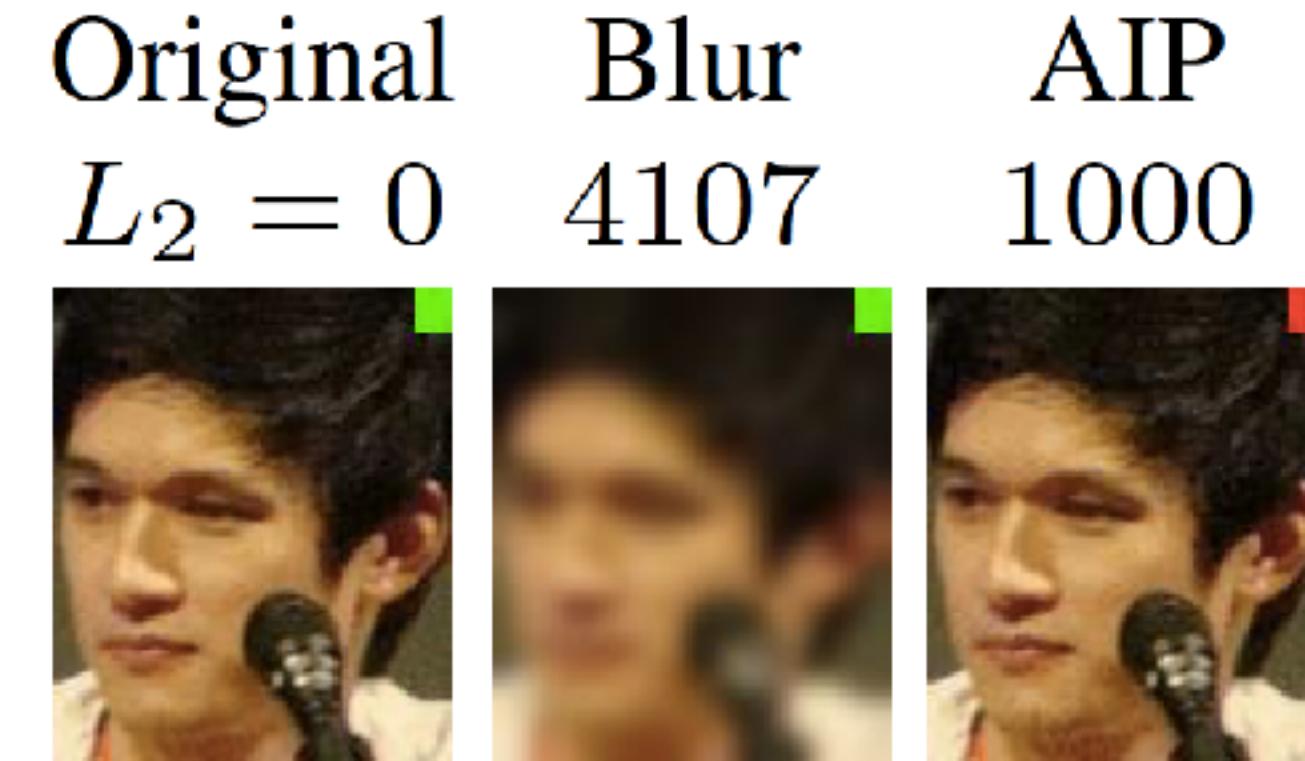
- AIP as an effective anonymisation tool for person recognition

- Robust AIP variants.

- Game theoretical framework:

- makes assumptions explicit
  - accounts for the lack of knowledge on the deployed recogniser.
  - analysis tool for a design space

- “Selective Confusion” possible





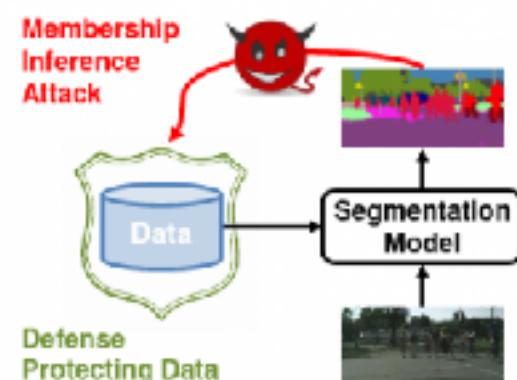
# Membership Inference Attacks

# Privacy Attacks and Privacy Preserving Machine Learning



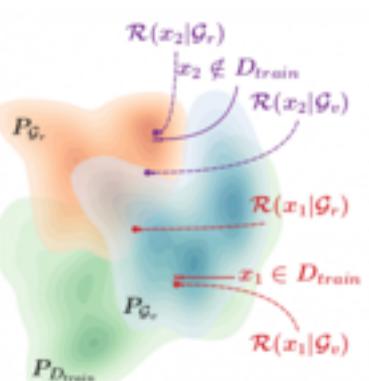
Salem; Zhang; Humbert; **Fritz**; Backes

**ML-Leaks**: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models  
**NDSS'19**



He; Rahimian; Schiele; **Fritz**

**Segmentations-Leak**: Membership Inference Attacks and Defenses in Semantic Image Segmentation  
**ECCV'20**



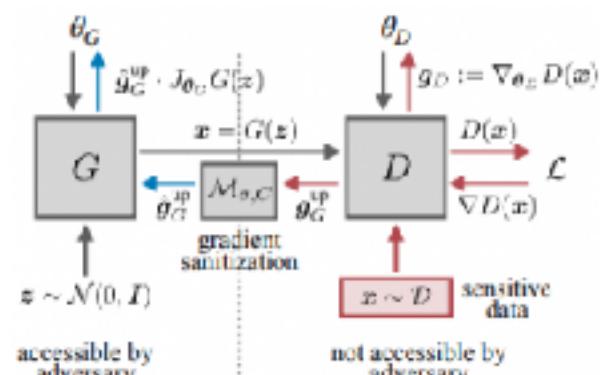
Chen; Yu; Zhang; **Fritz**

**GAN-Leaks**: A Taxonomy of Membership Inference Attacks against GANs  
**CCS'20**

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

Salem; Bhattacharyya; Backes; **Fritz**; Zhang

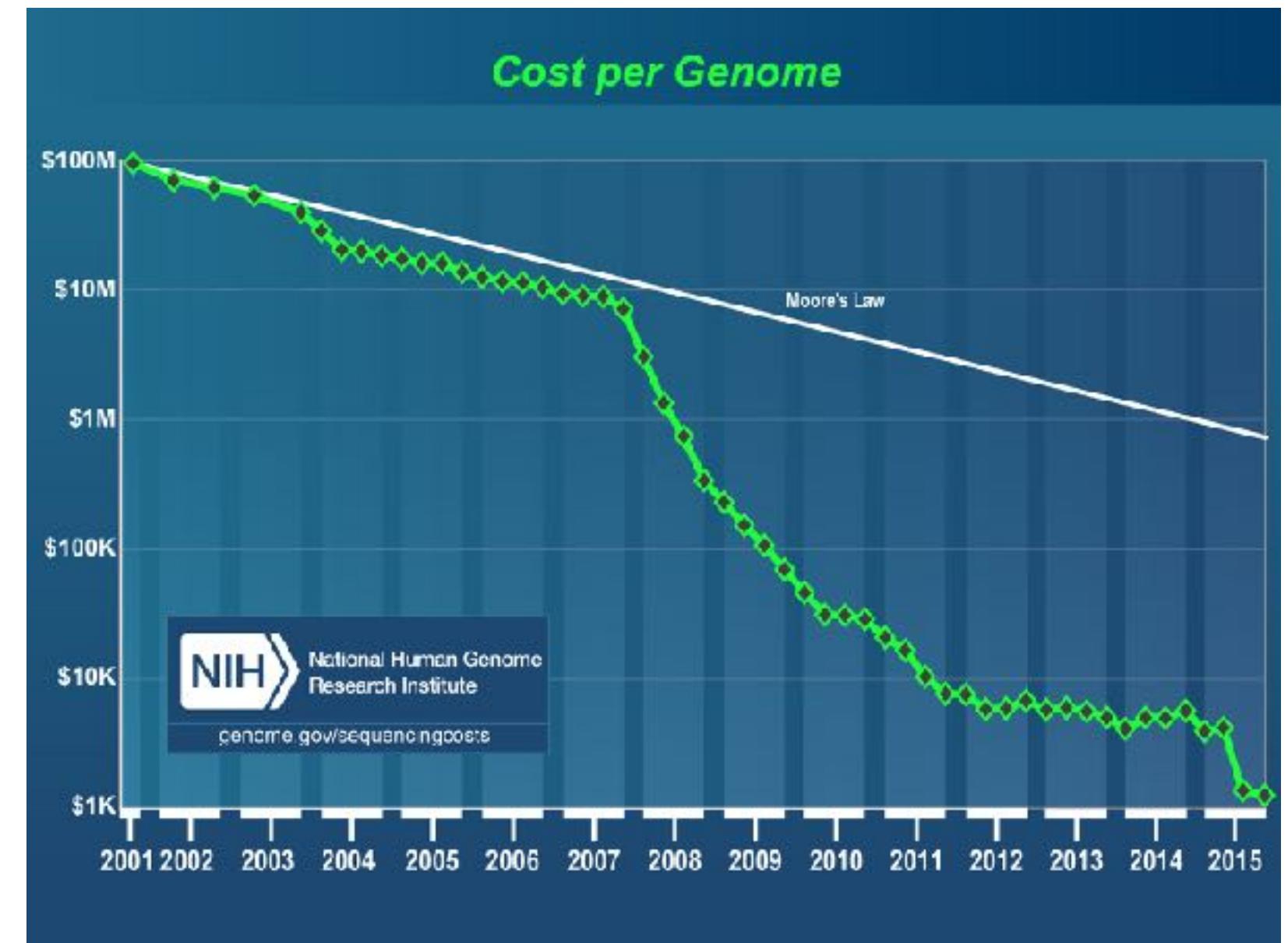
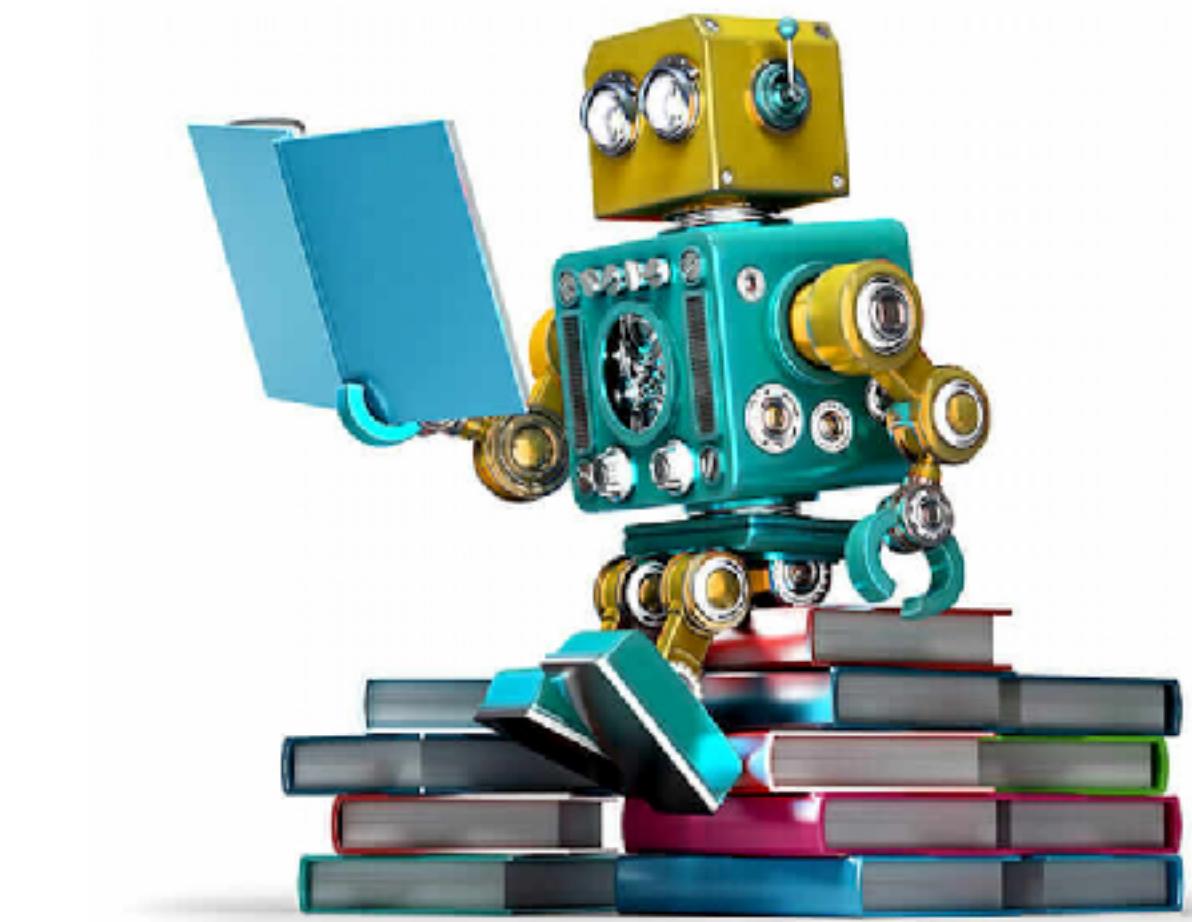
**Updates-Leak**: Data Set Inference and Reconstruction Attacks in Online Learning  
**USENIX Security'20**



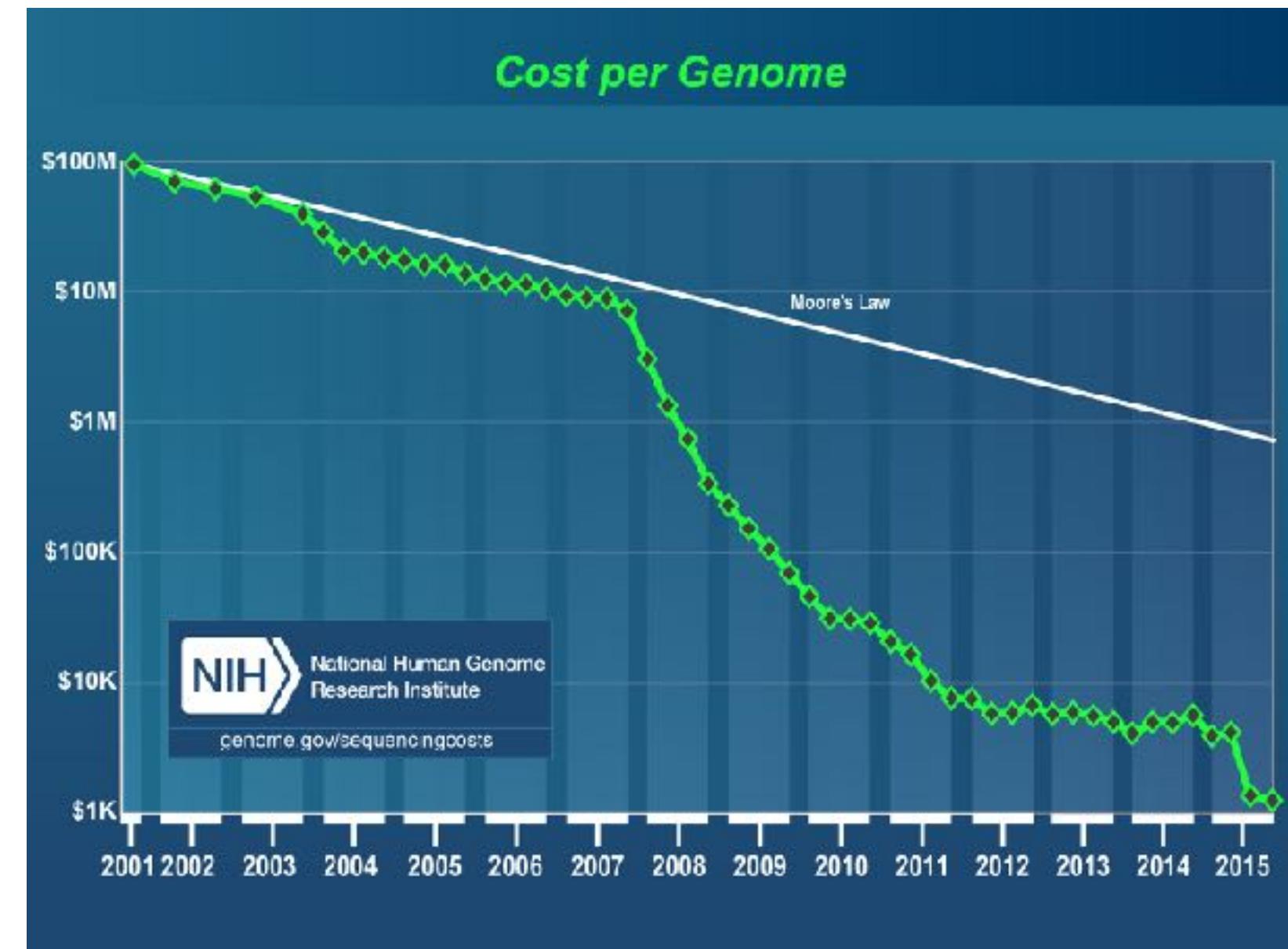
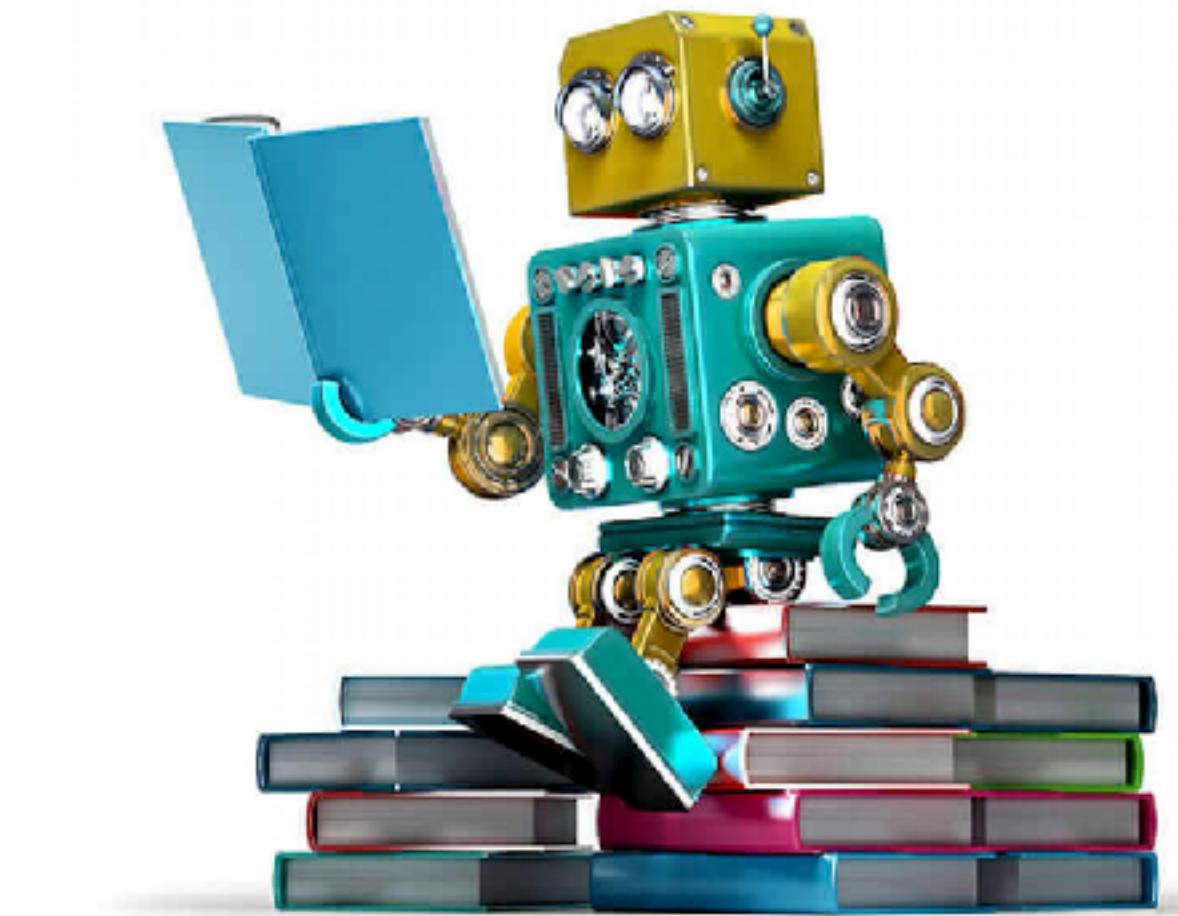
Chen; Orekondy; **Fritz**

**GS-WGAN**: A Gradient-Sanitized Approach for Learning Differentially Private Generators  
**NeurIPS'20**

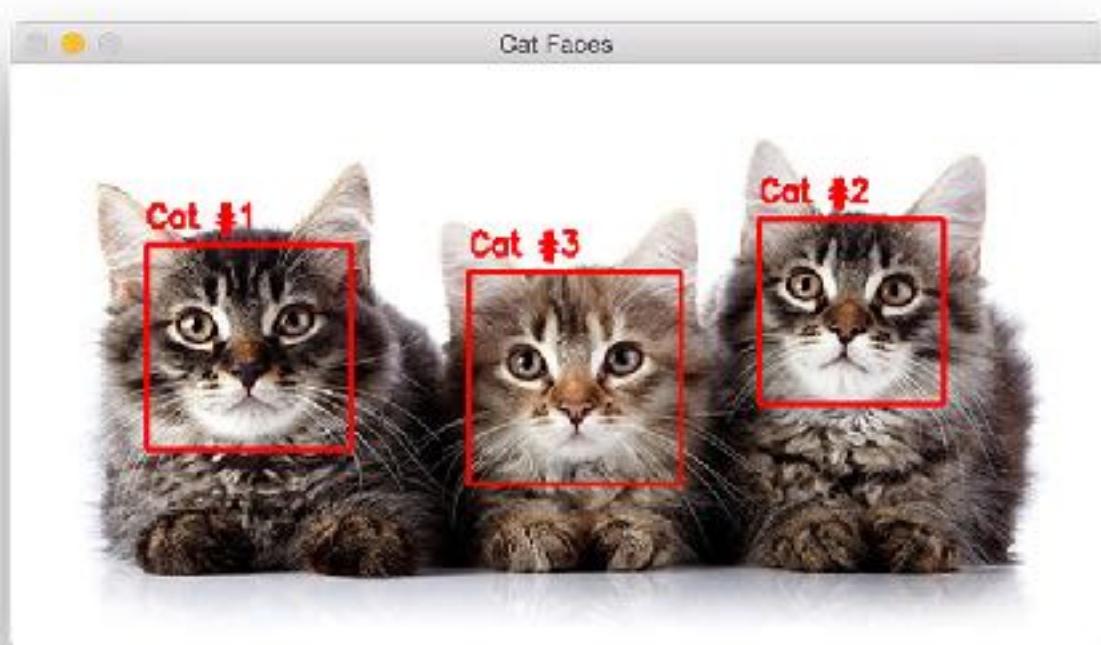
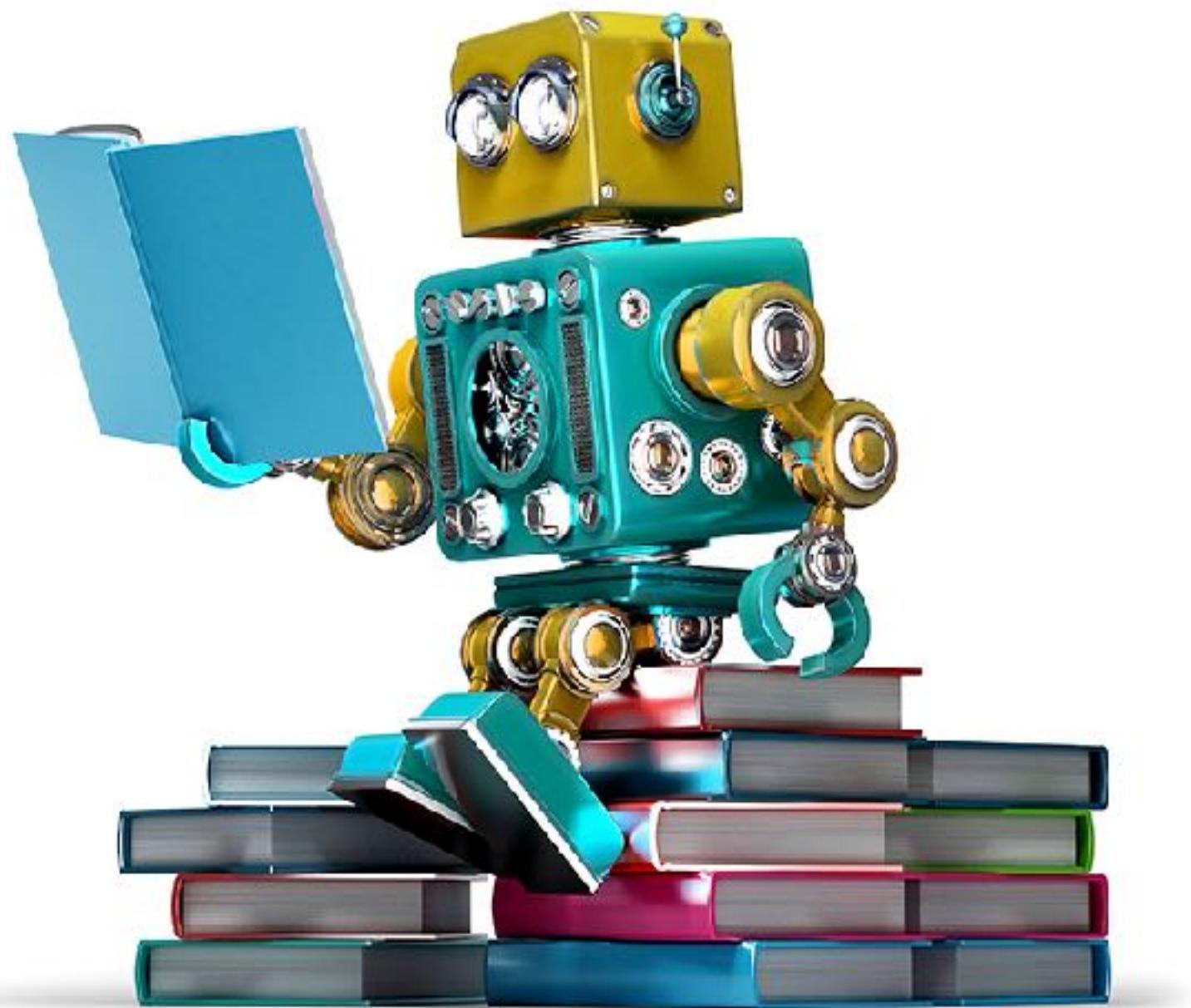
# New Technologies



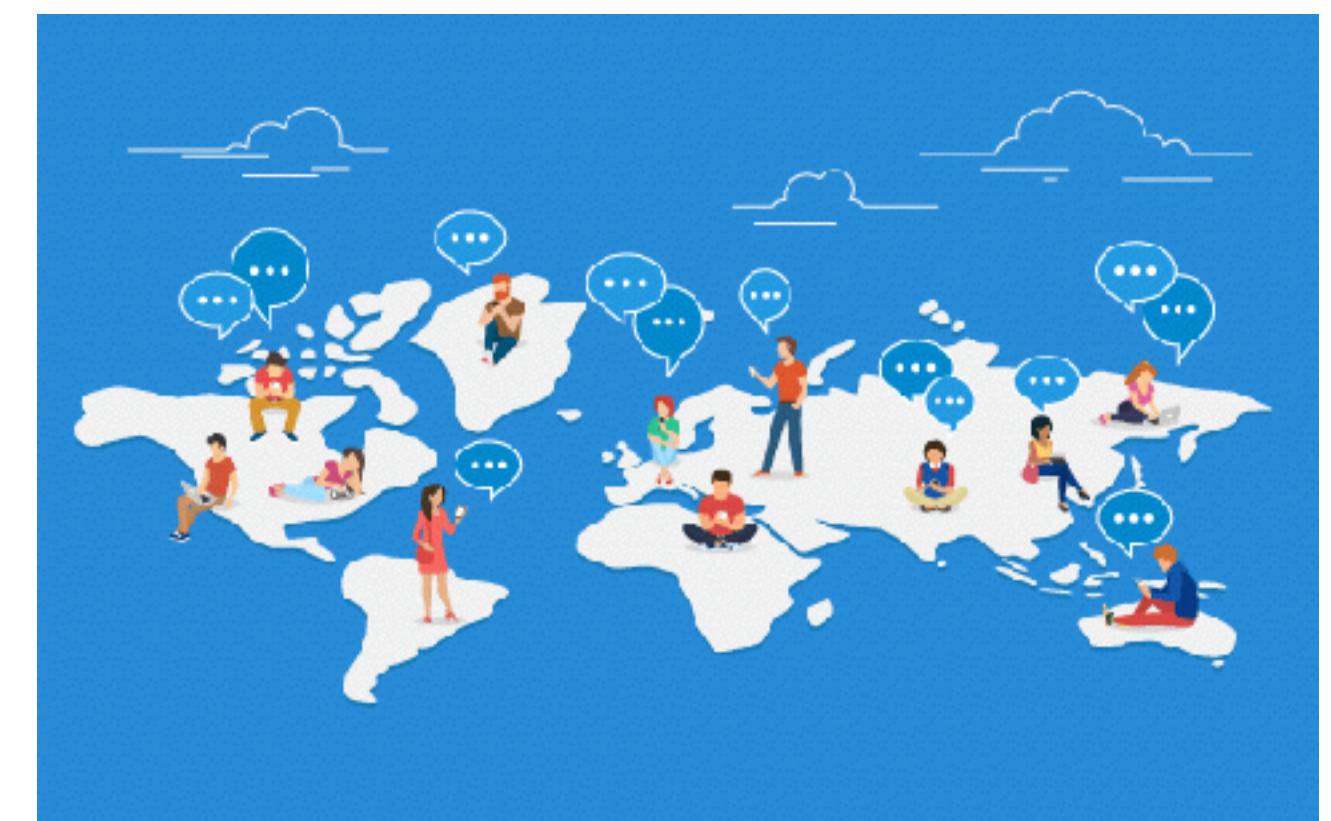
# New Privacy Challenges!



# Era of Machine Learning

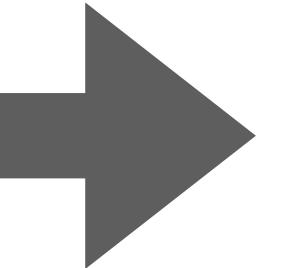
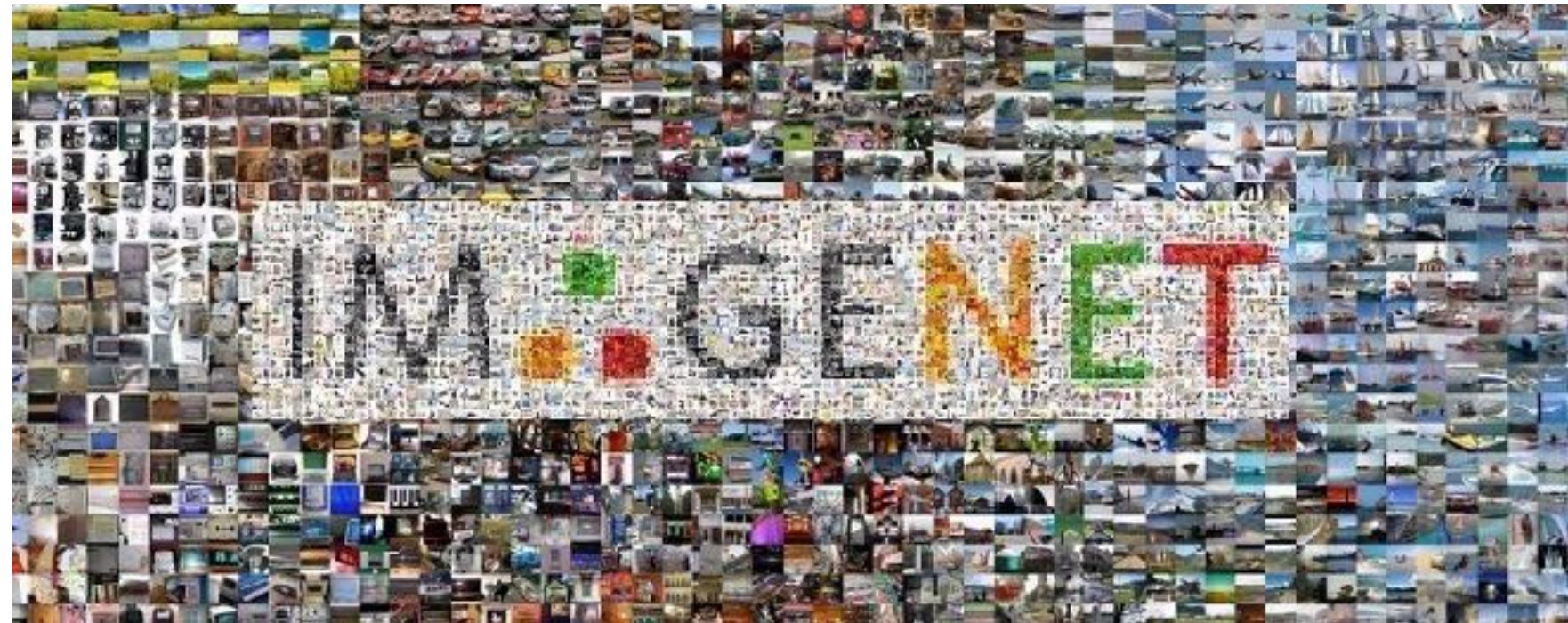


- ML models are often trained on sensitive data
  - Biomedical data
  - Keyboard input
  - Image / Language

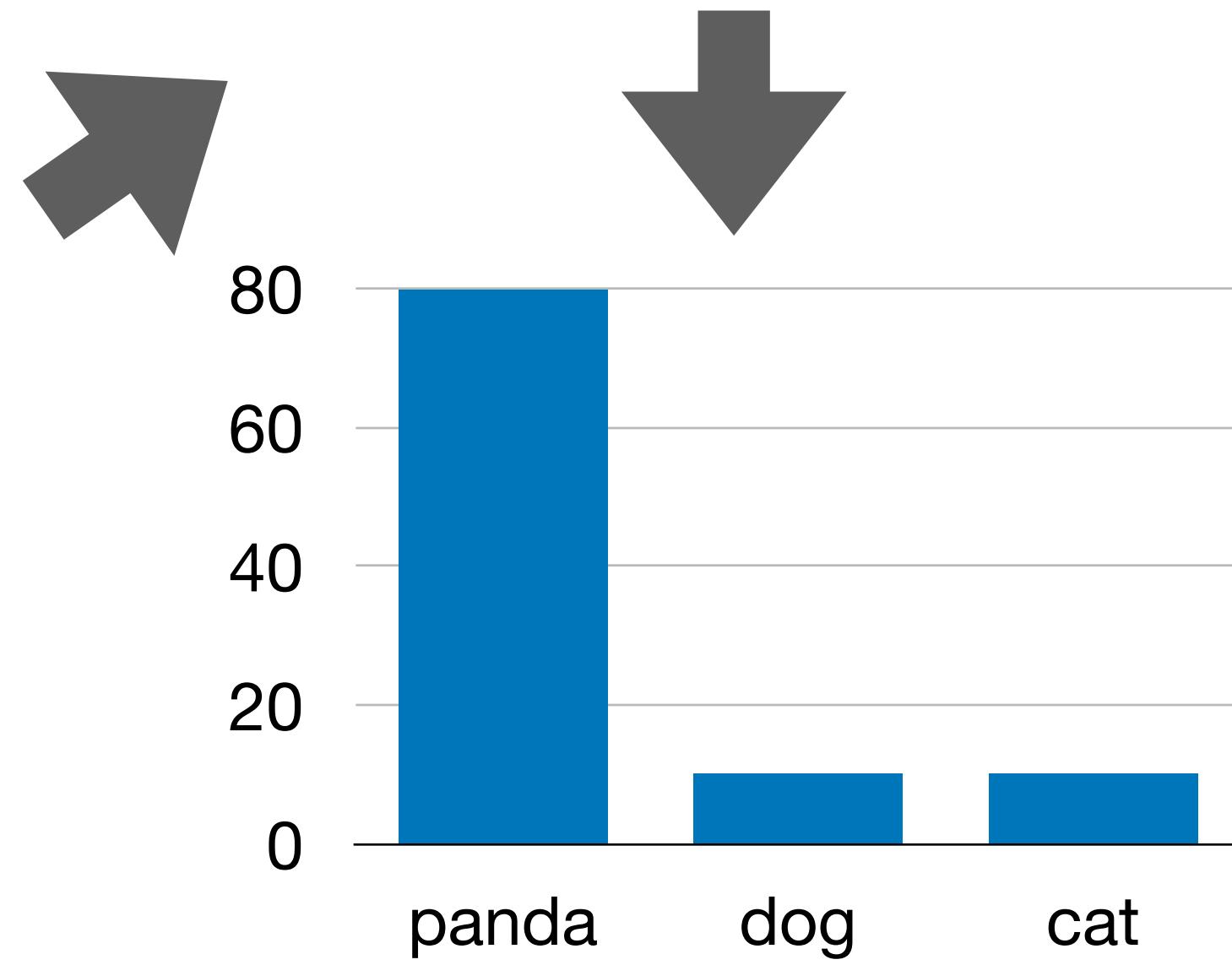


How About Privacy?

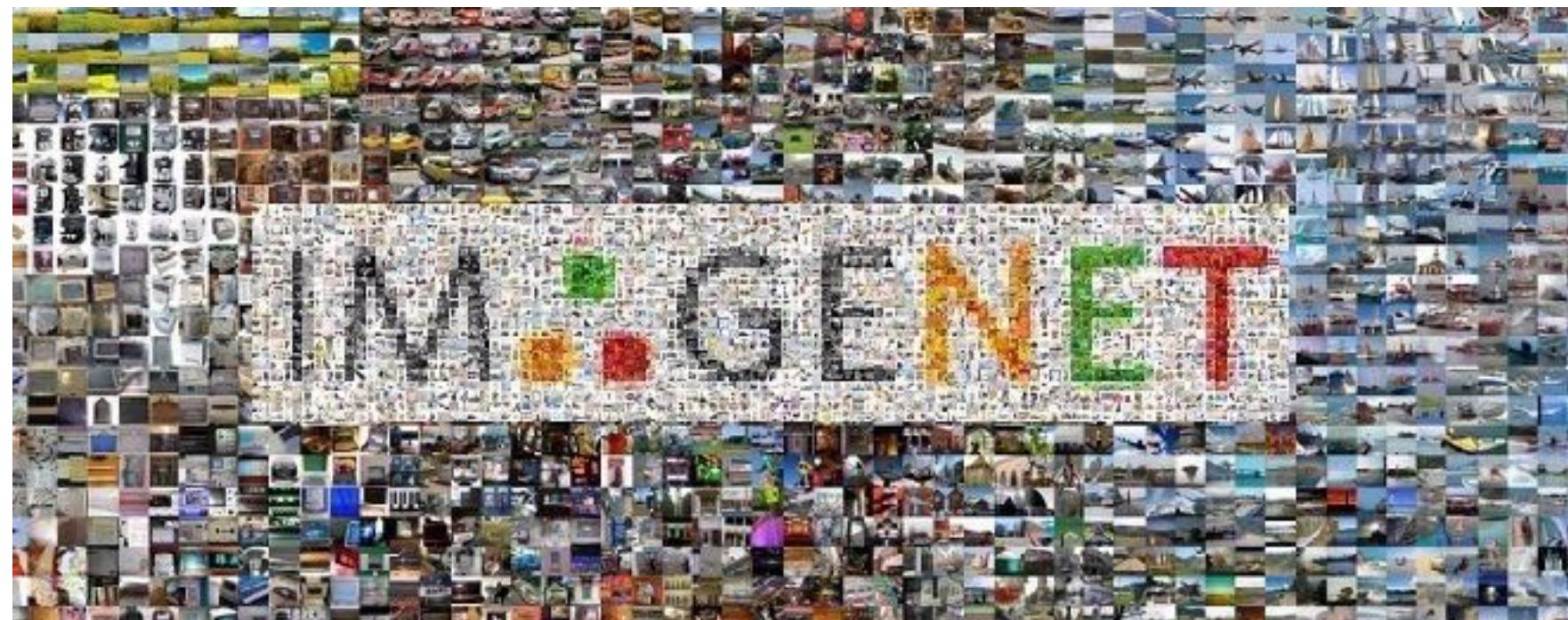
# Machine Learning Pipeline



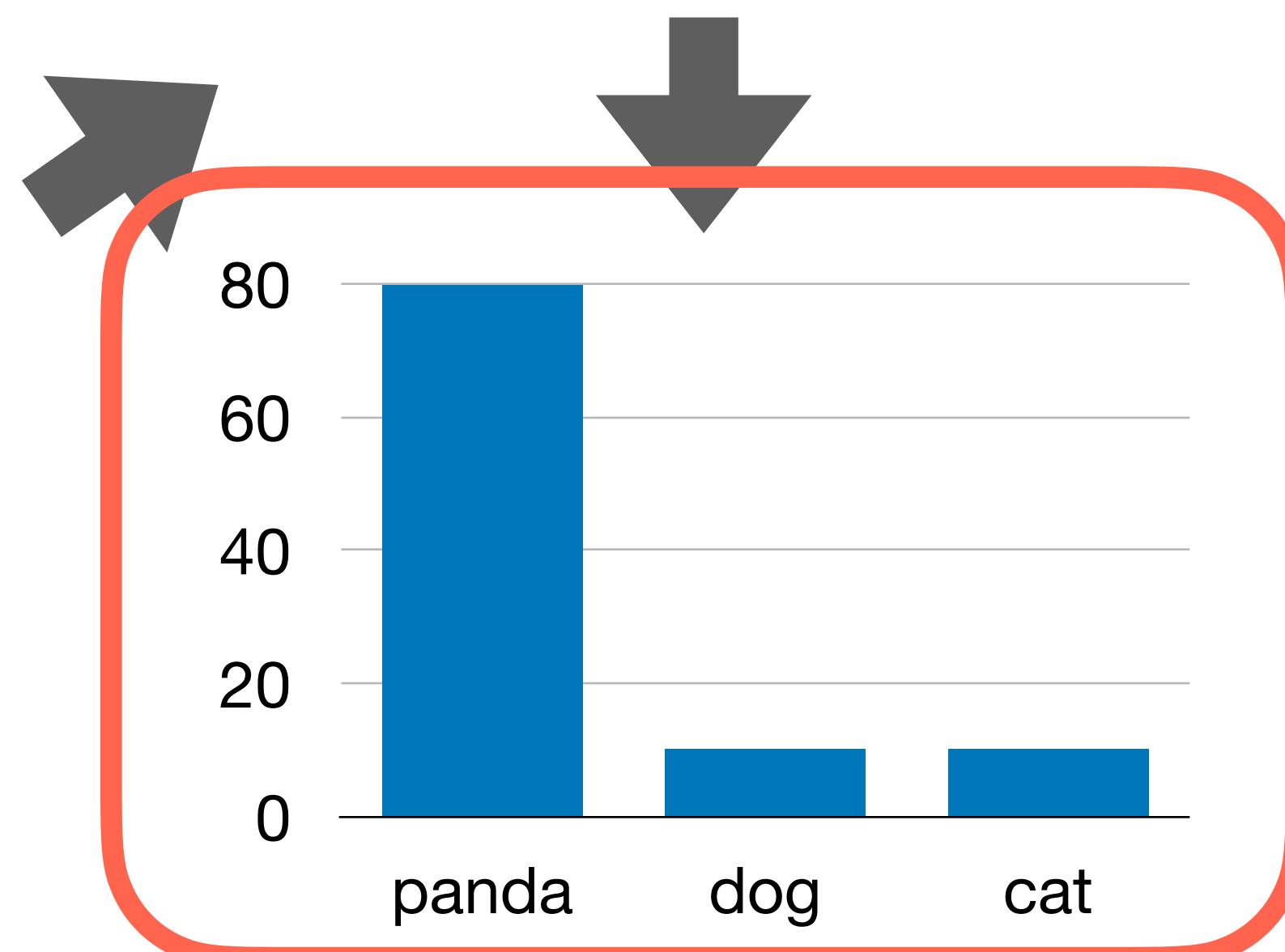
ML Model



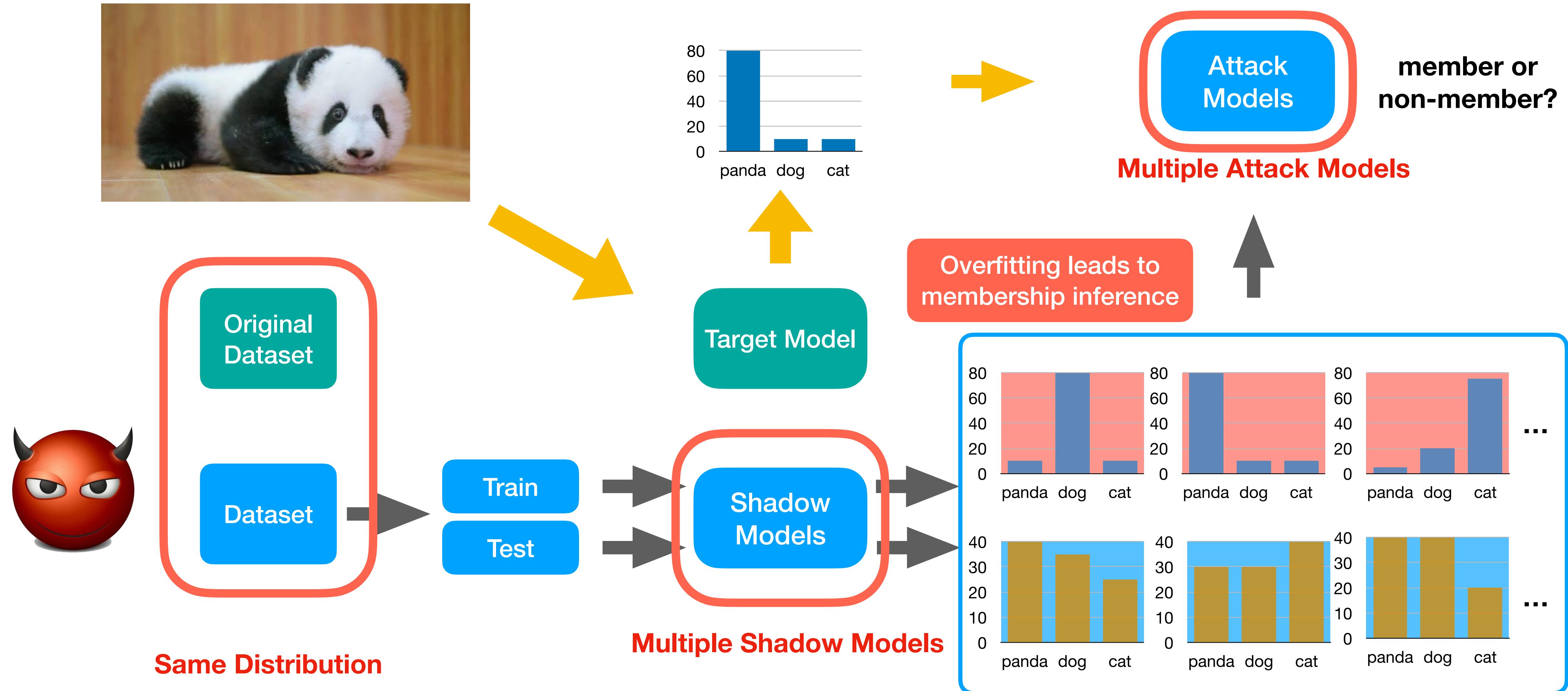
# Membership Inference against Machine Learning



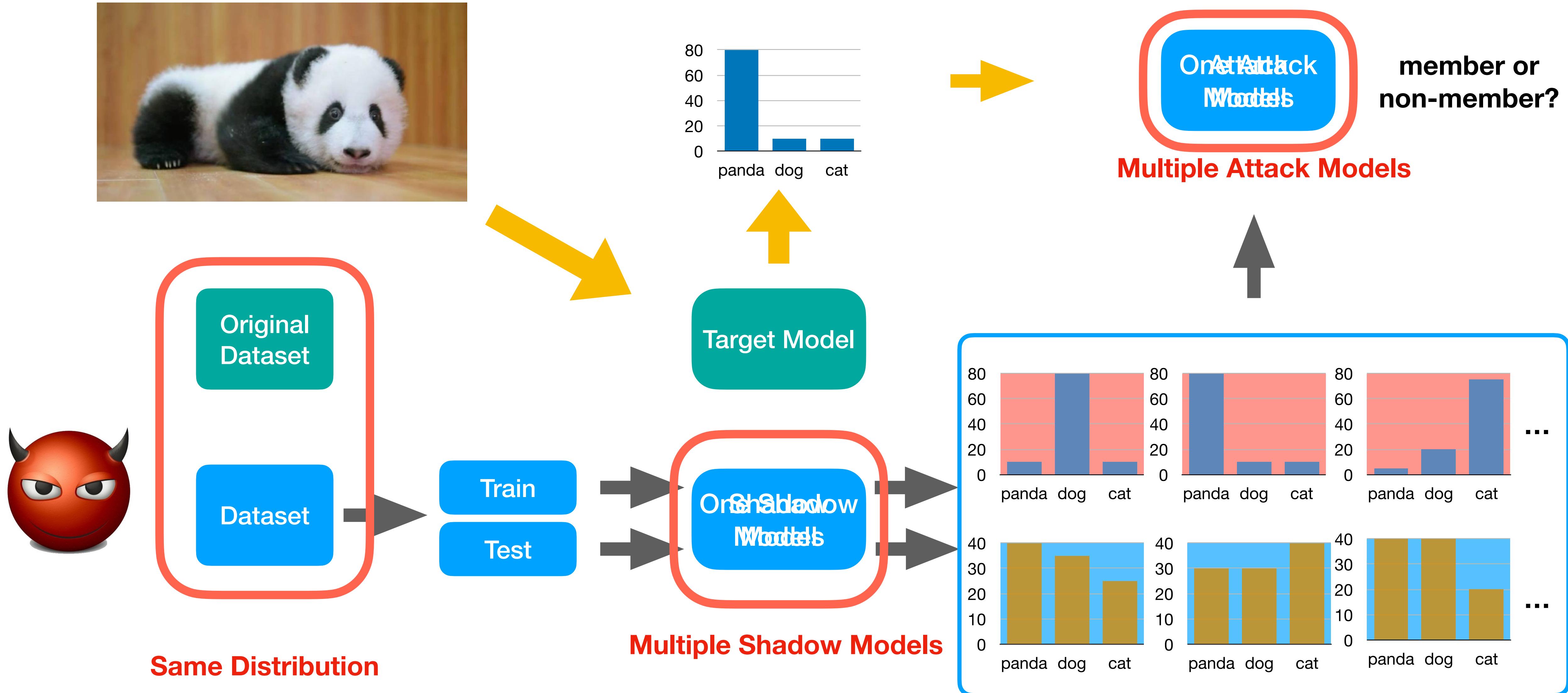
ML Model



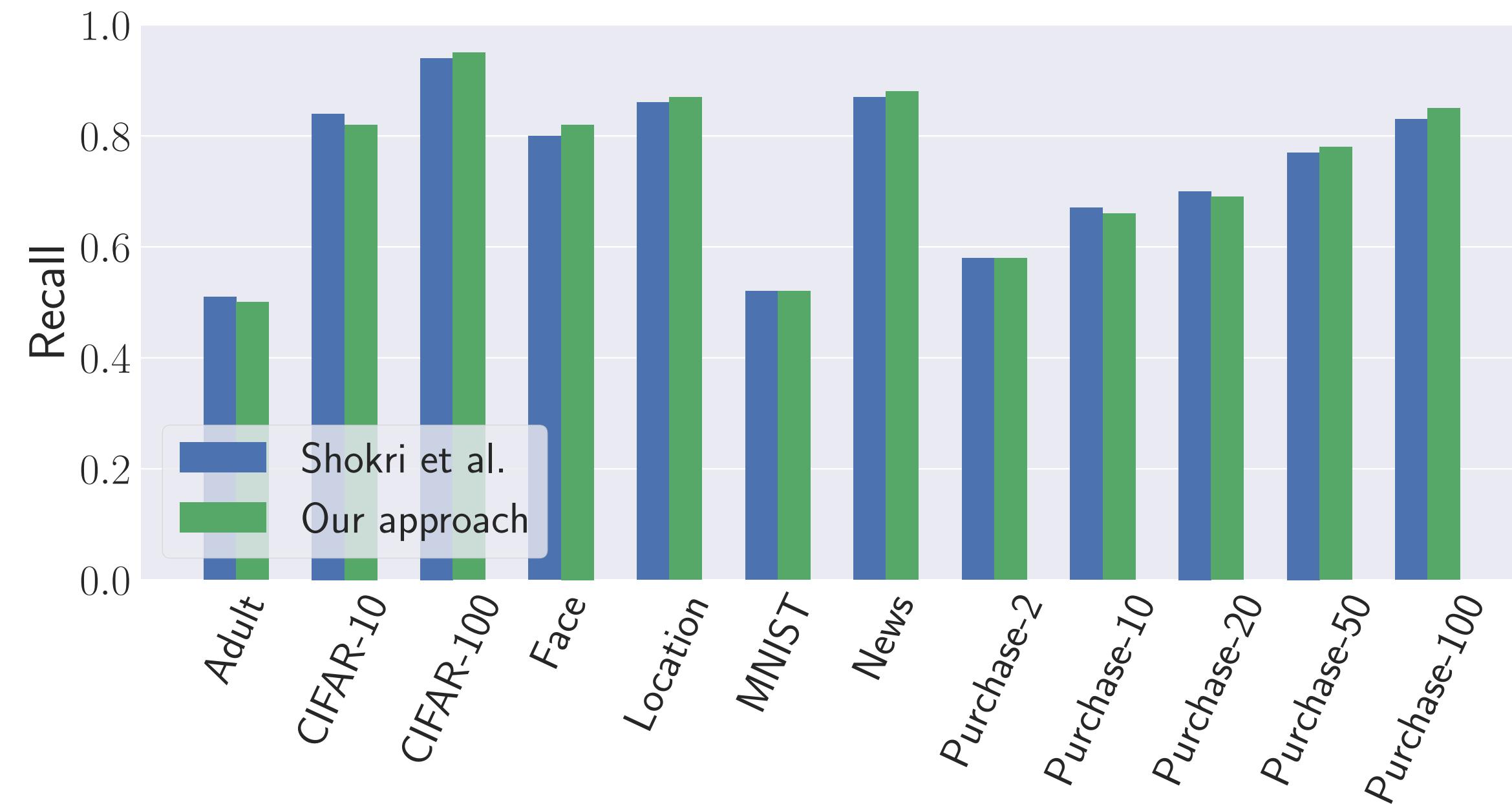
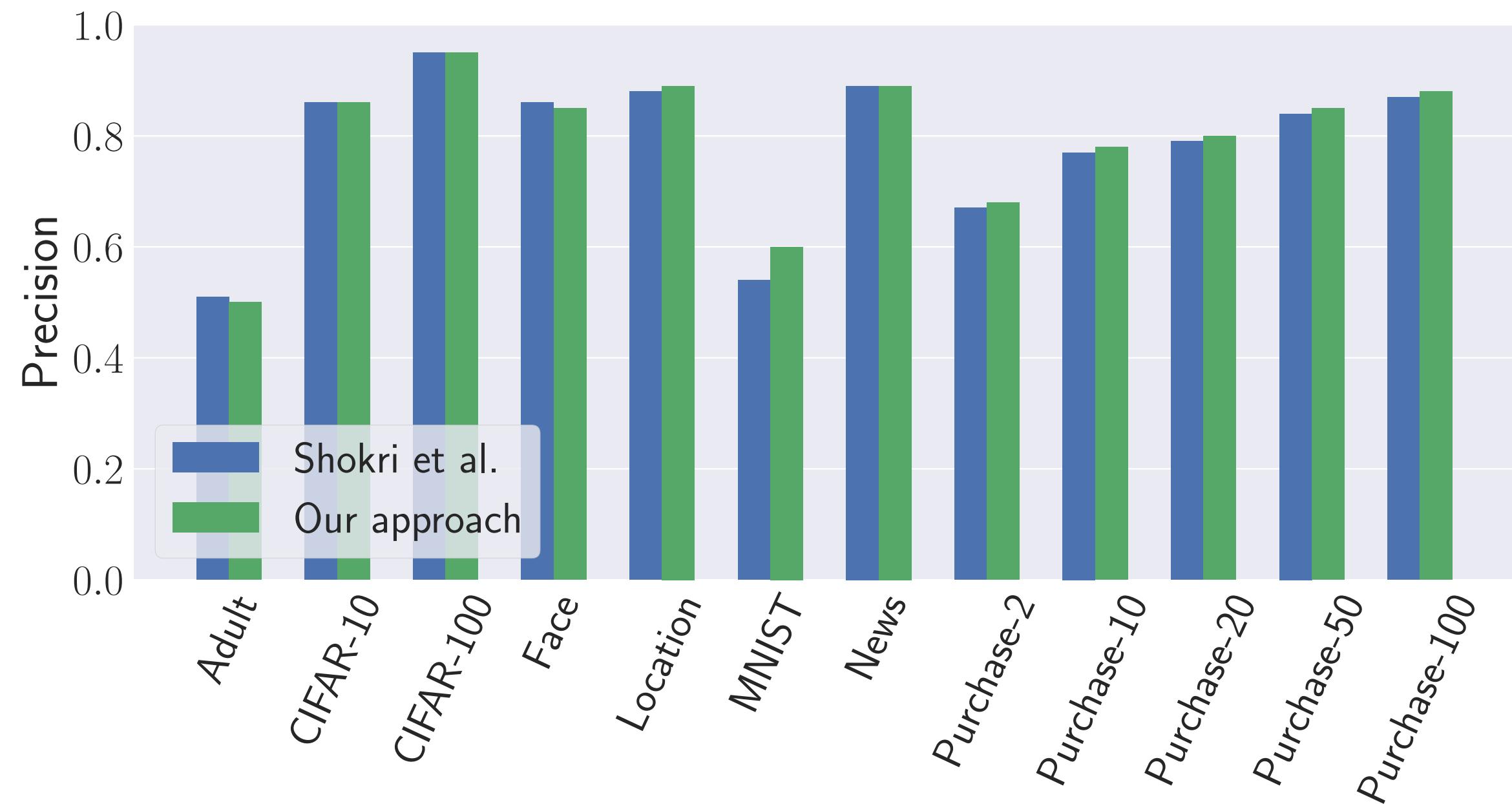
# Attack by Shokri et al.



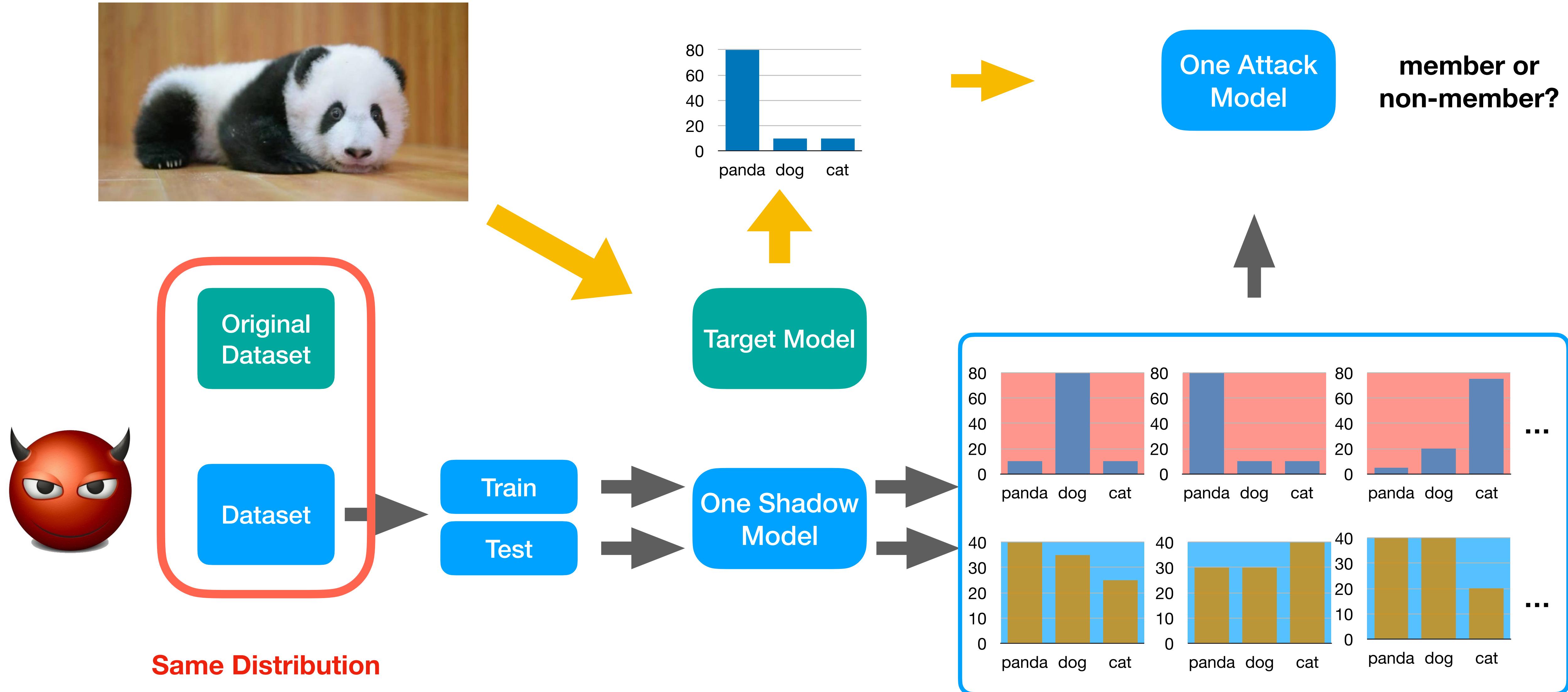
# Attack 1



# Attack 1



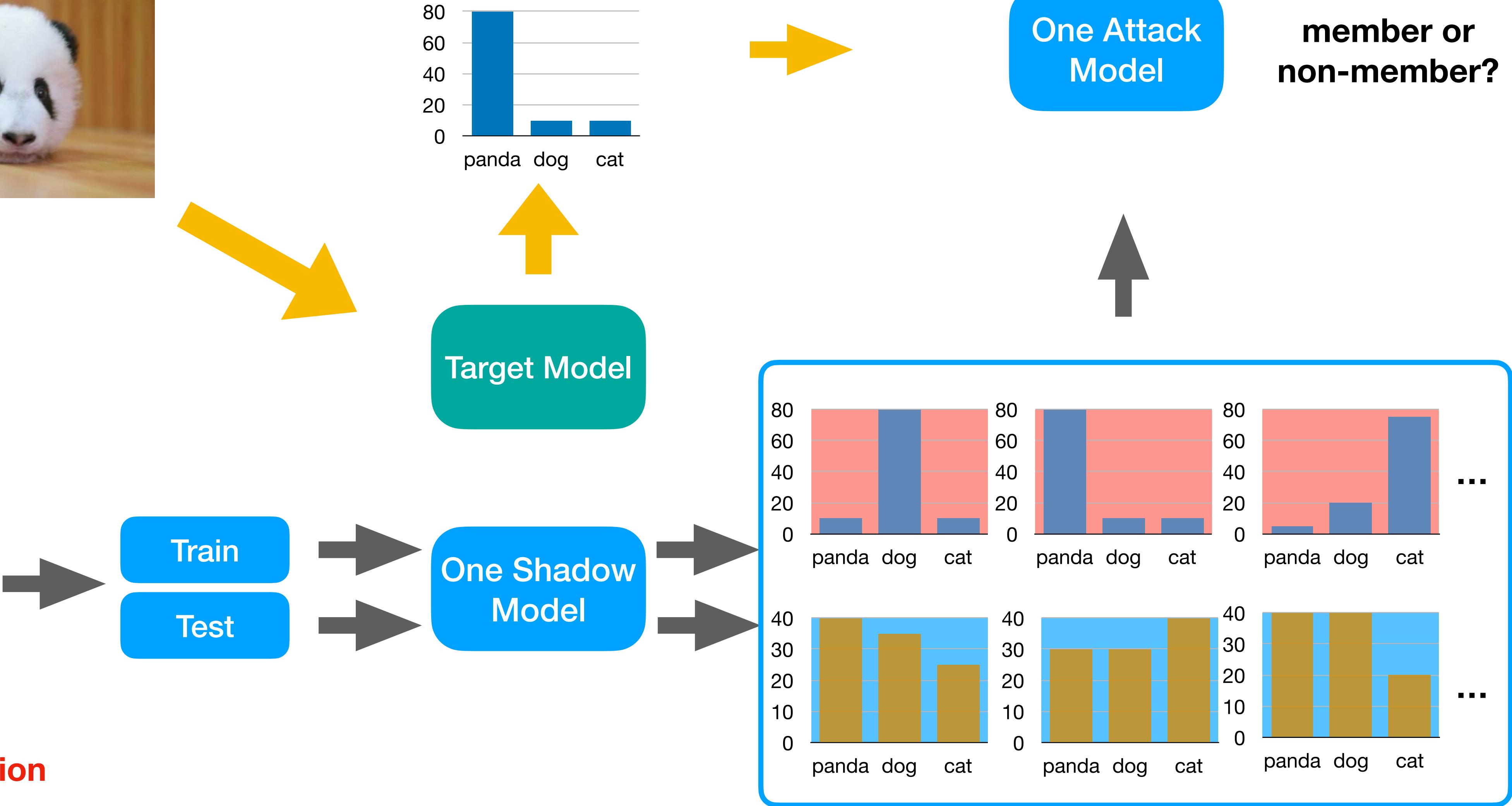
# Attack 1



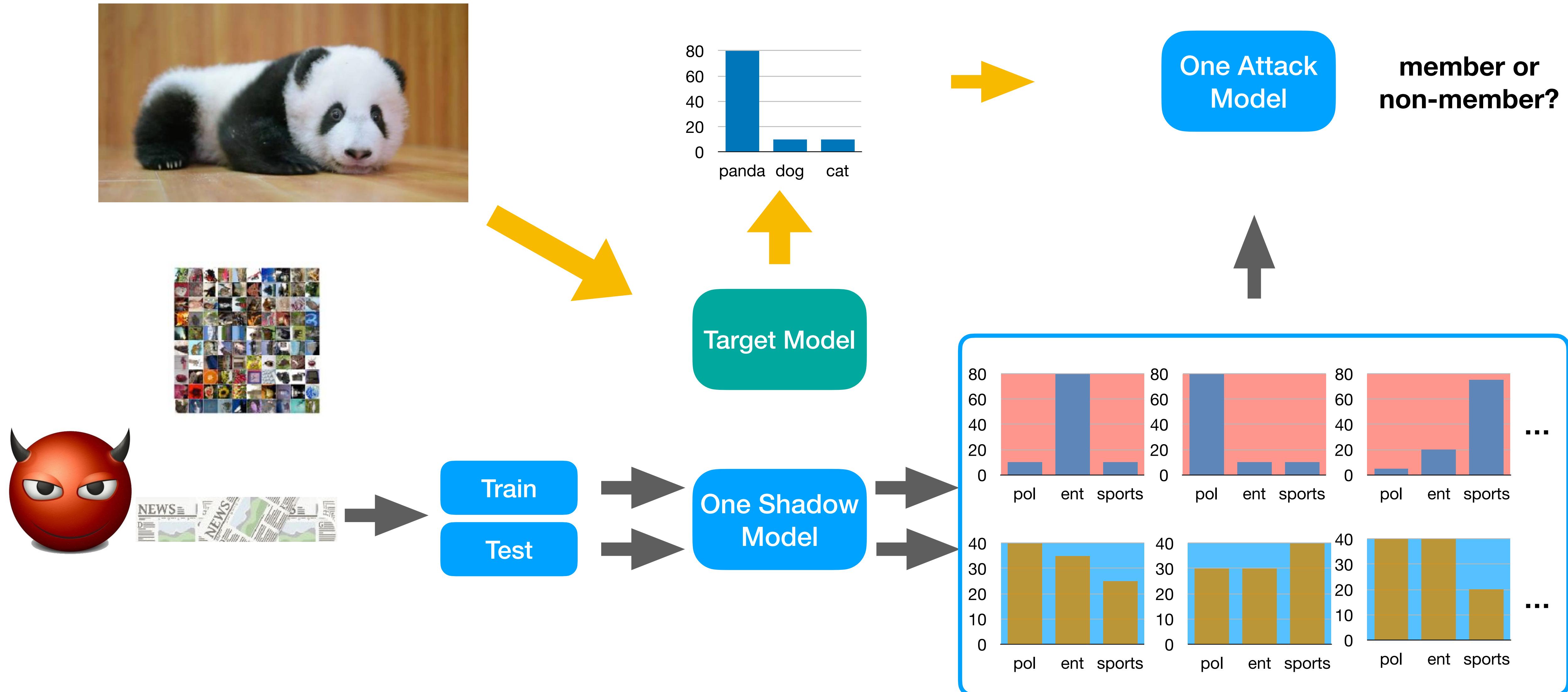
# Attack 1



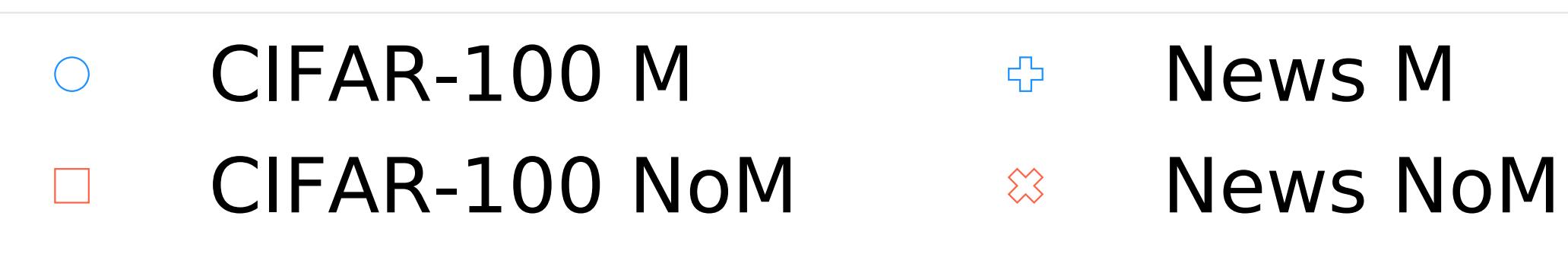
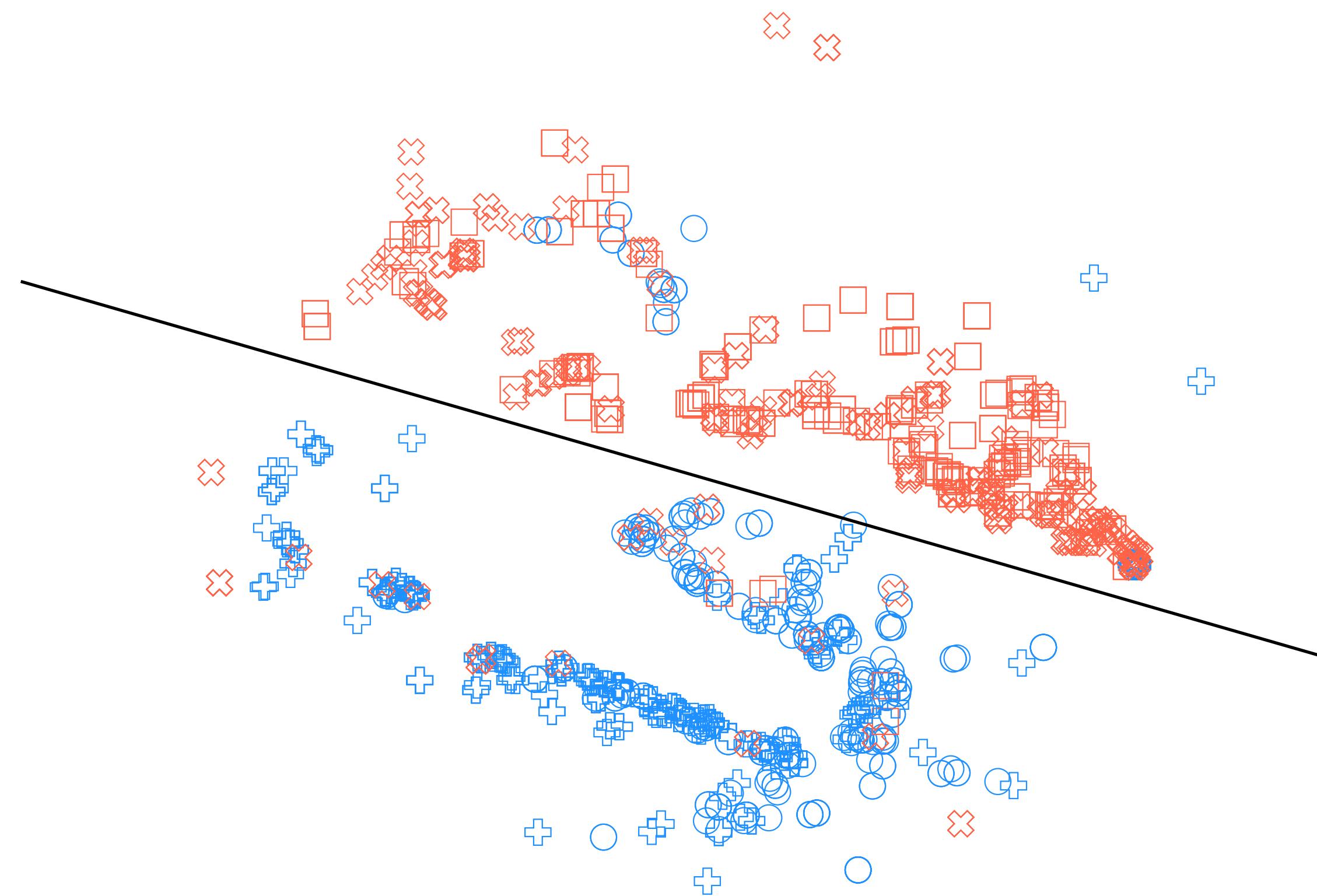
Same Distribution



# Attack 2



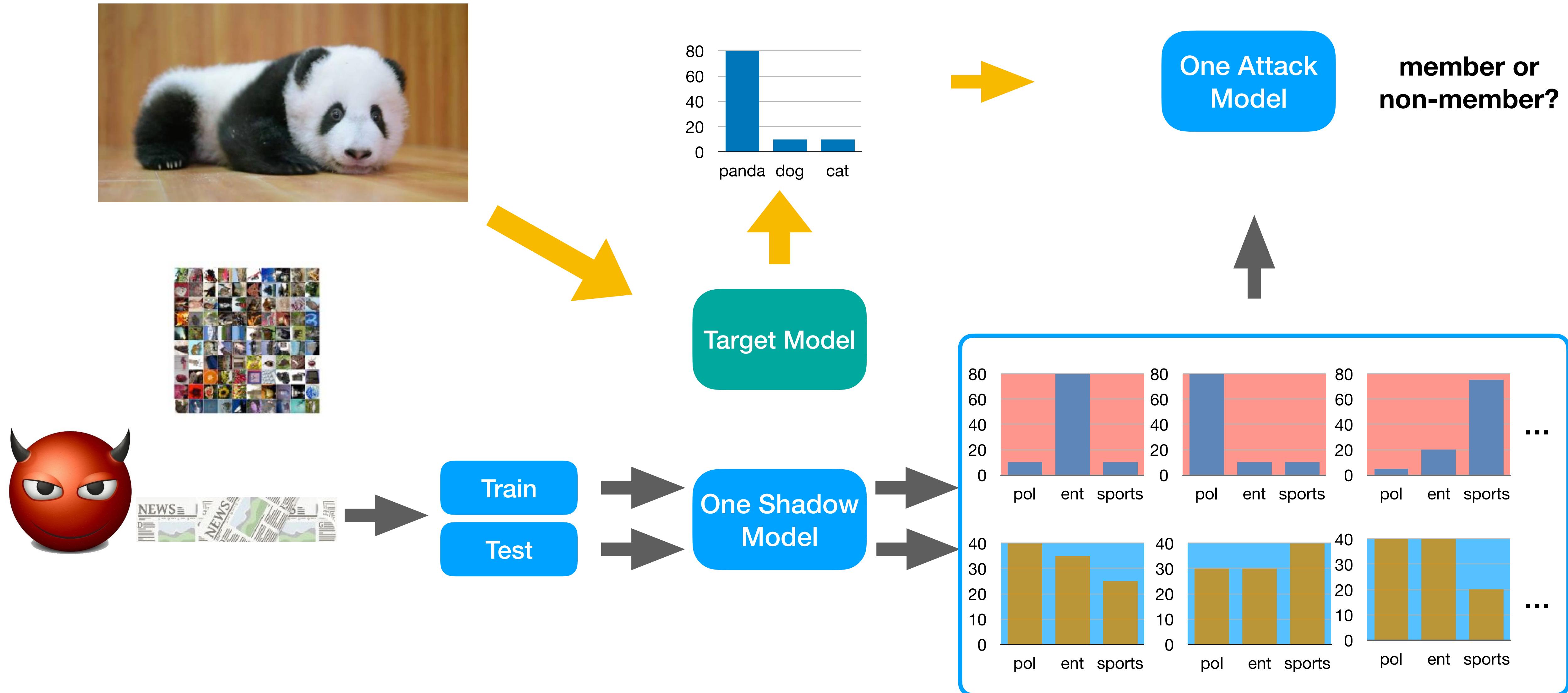
# Attack 2



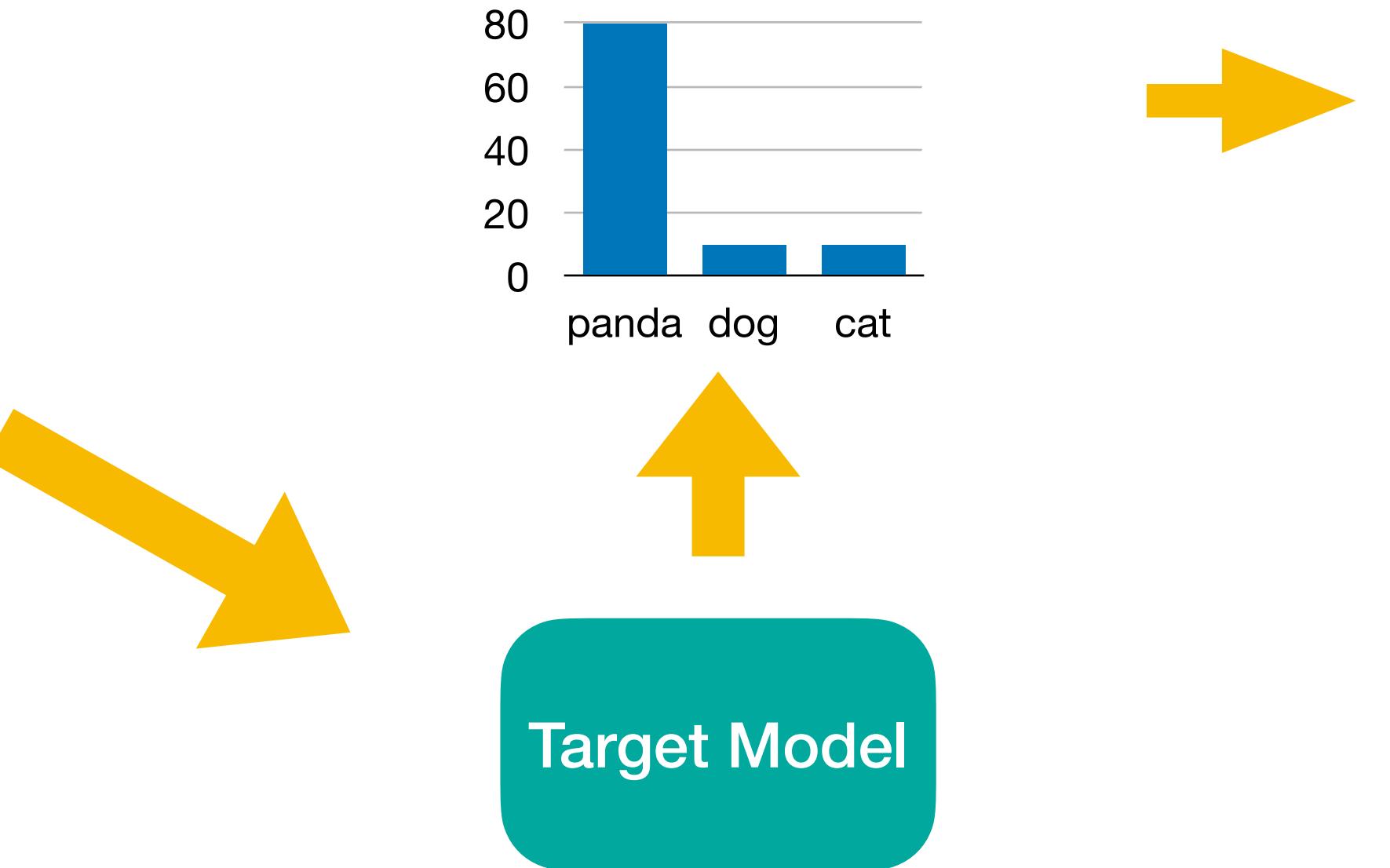
We didn't transfer the dataset!

ML models got overfitted in a similar way!

# Attack 2



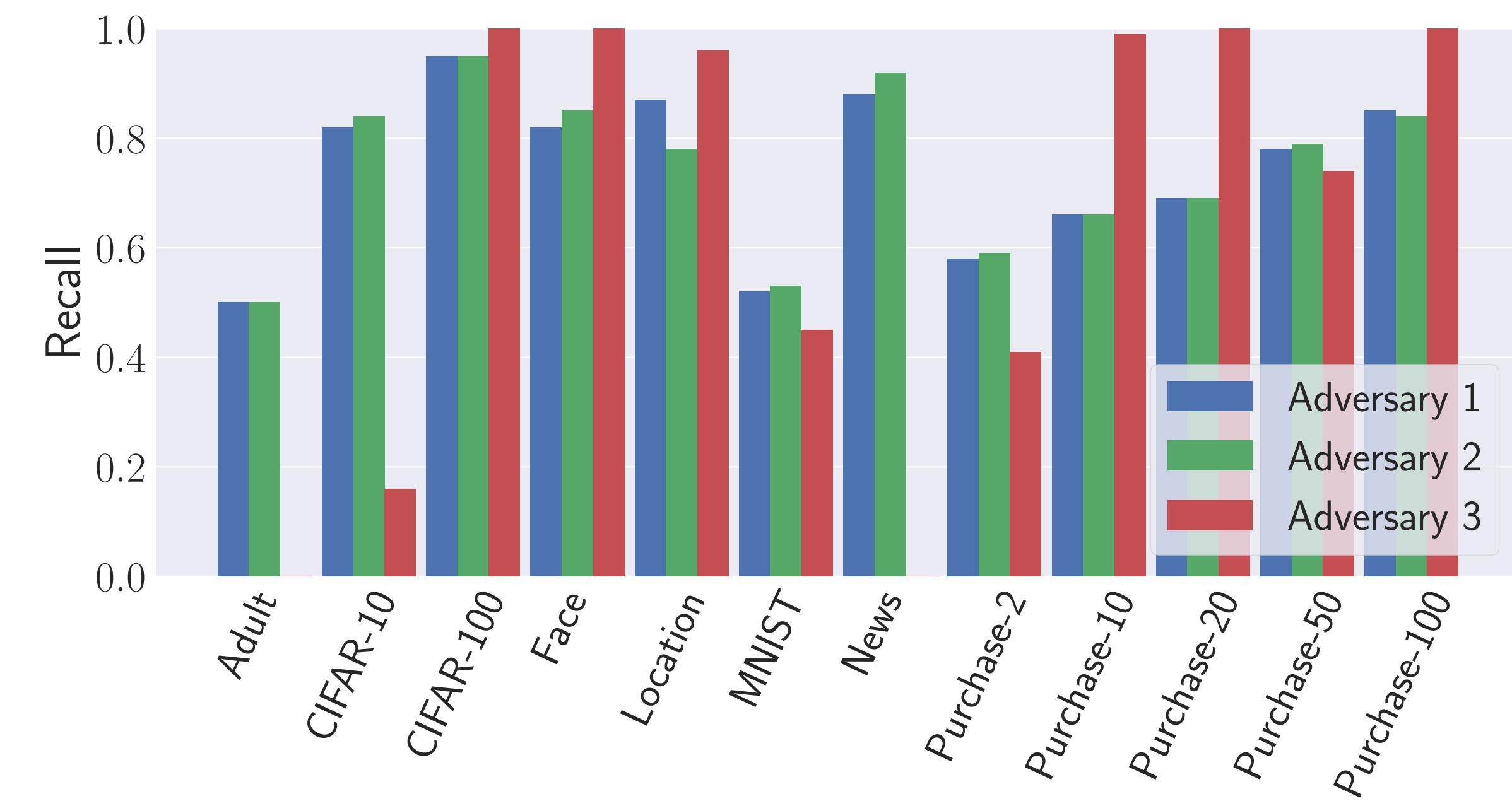
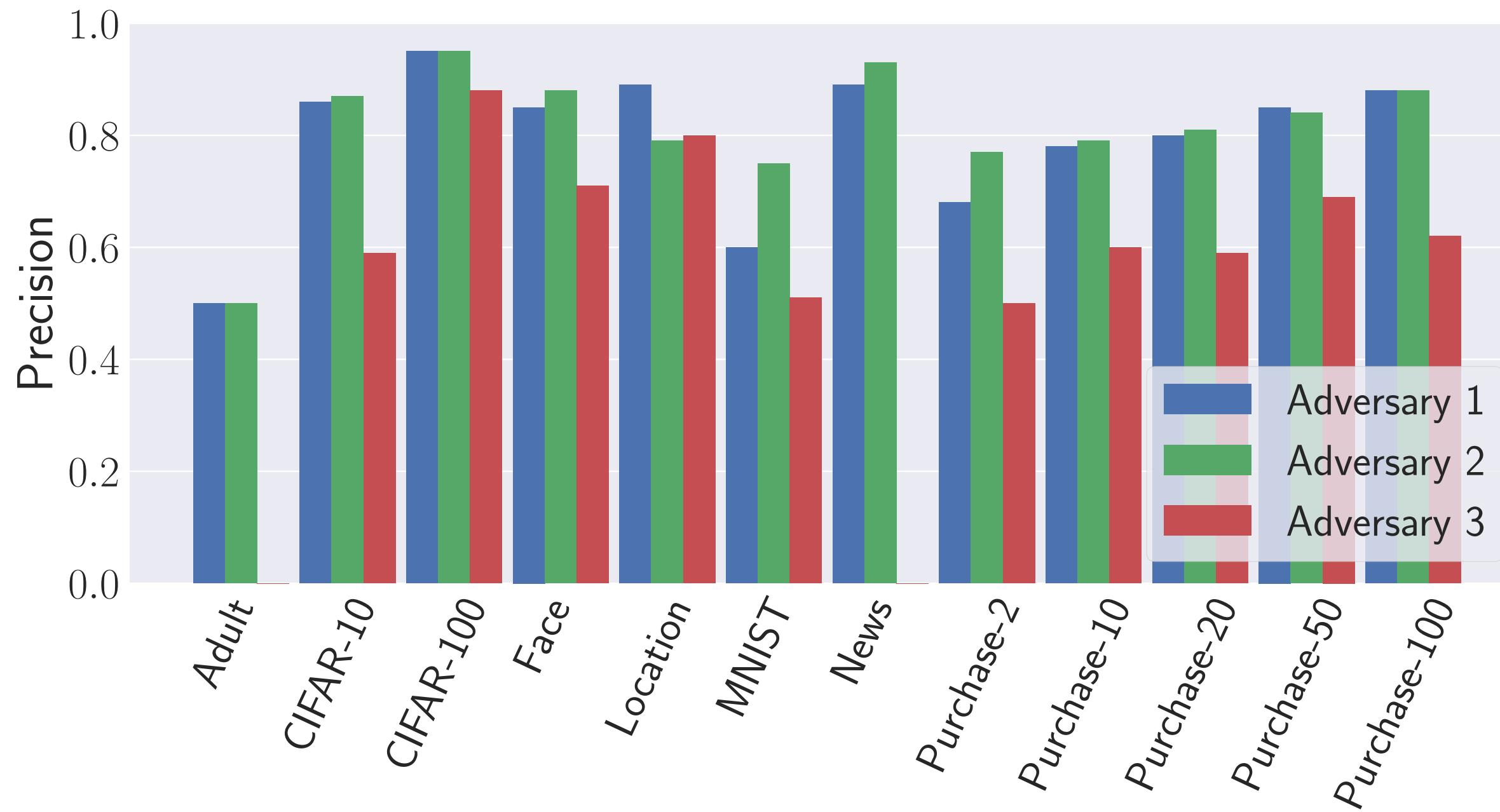
# Attack 3



Threshold  
Picking

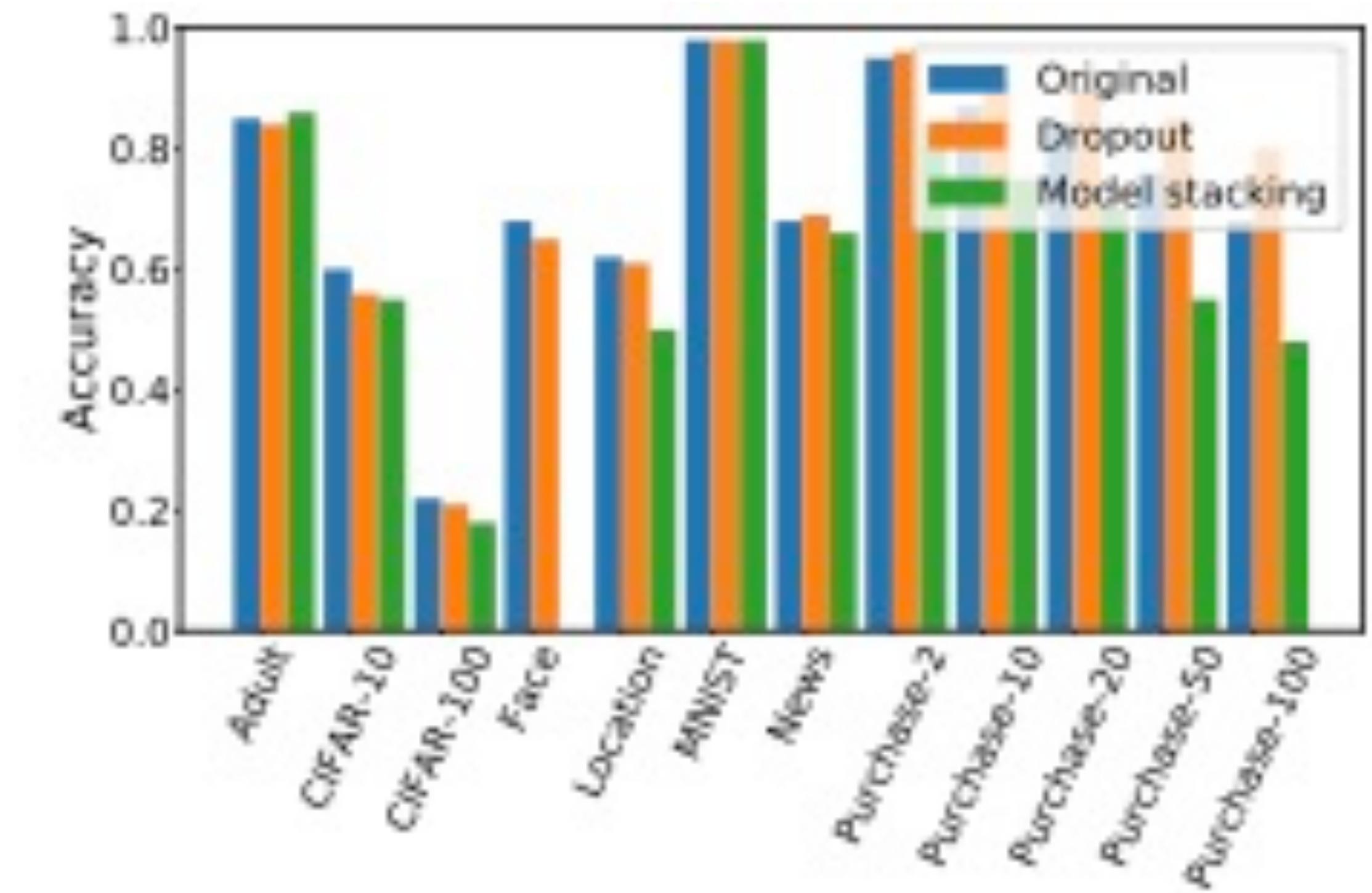
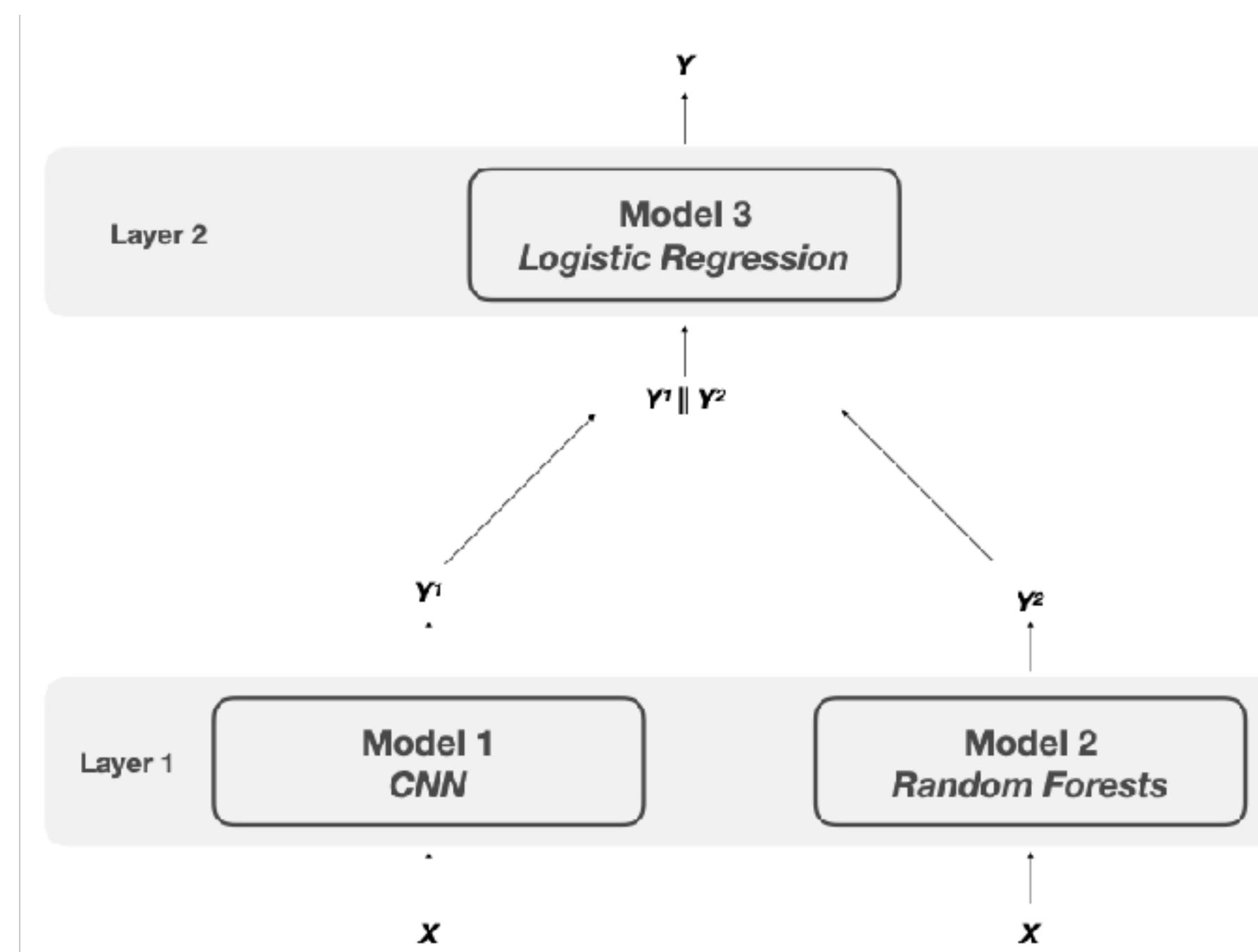
member or  
non-member?

# All Together



# How To Defend the Attack?

- Reduce overfitting
  - E.g. regularize: Dropout
  - Combine models: Model Stacking
  - Differential Privacy (later)



(c) Accuracy.



- Membership Inference can tell if a data point was in training set
- Works for classification model
- How about other models?

# Segmentations Leak: Membership Inference Attacks and Defenses in Semantic Image Segmentation

Yang He, Shadi Rahimian, Bernt Schiele and Mario Fritz

ECCV2020



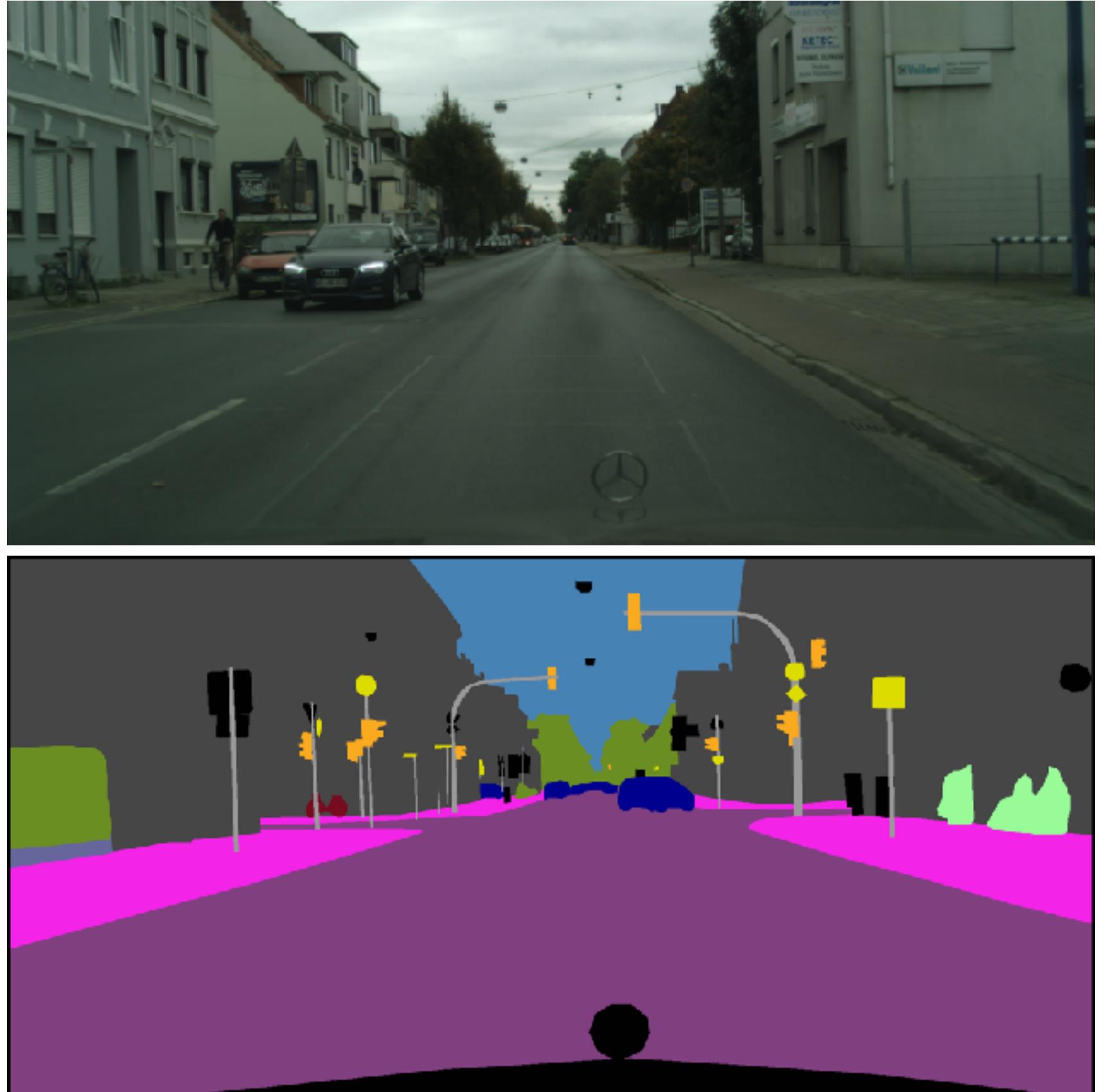
**mpii**  
max planck institut  
informatik



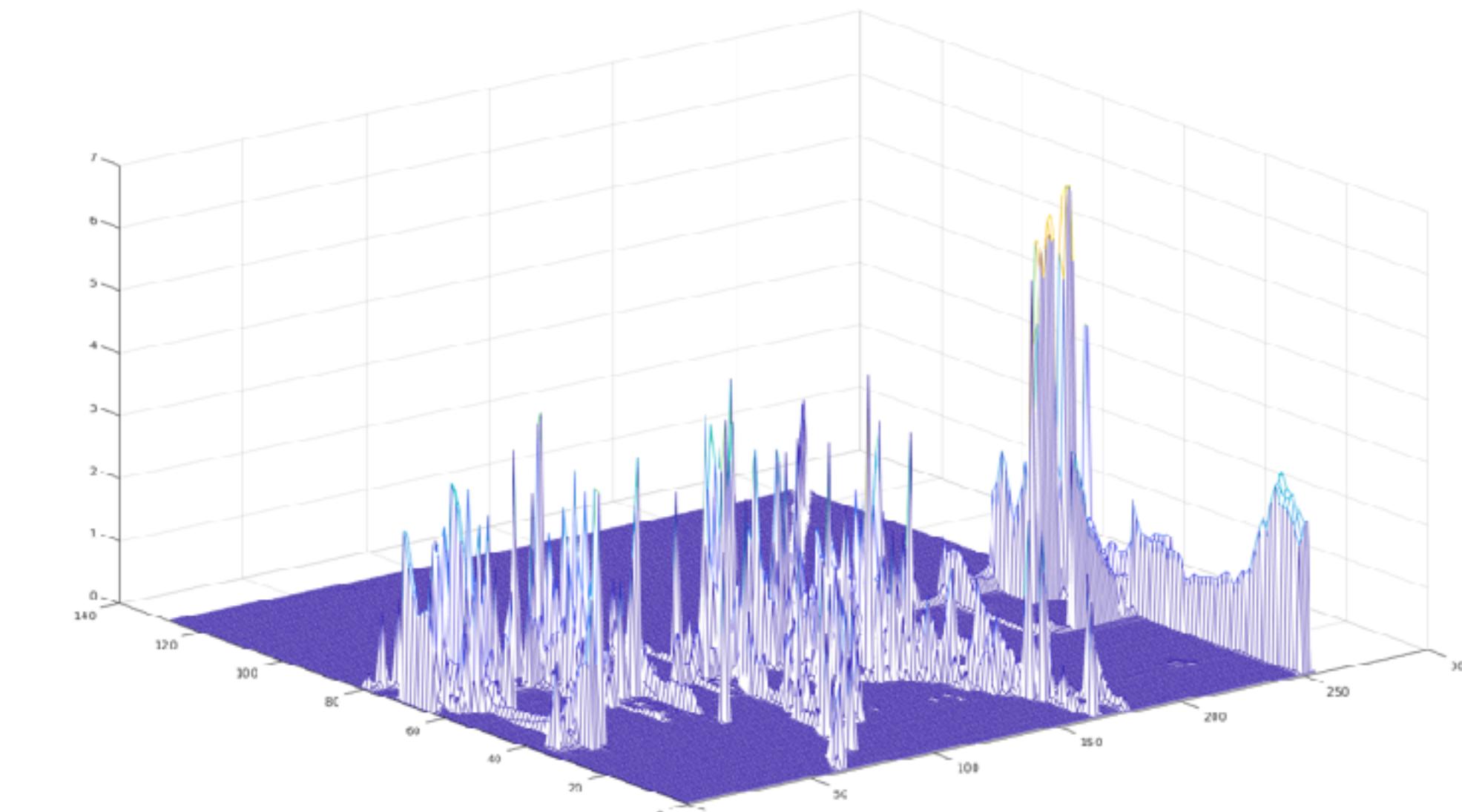
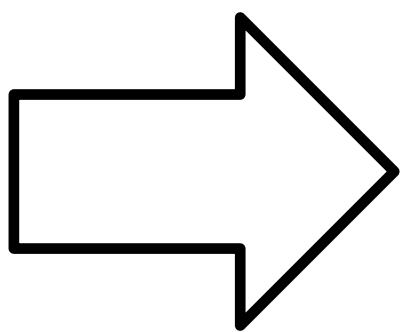
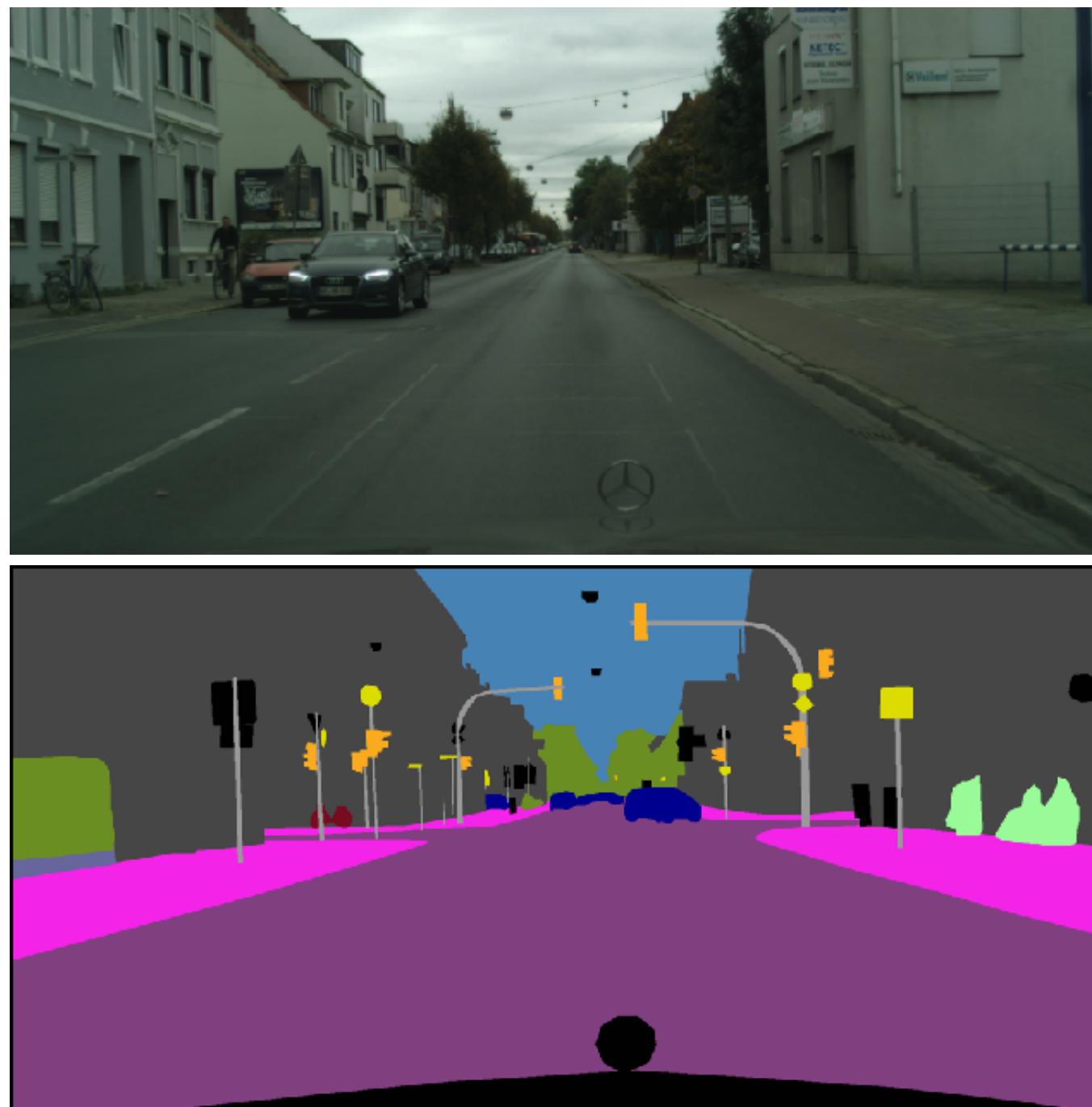
**SAARLAND**  
UNIVERSITY The logo of Saarland University features a detailed illustration of an owl perched on a branch, symbolizing wisdom and knowledge.

# Semantic Segmentation

- Deep Fully Convolution Networks can predict highly accurate segmentation masks
- One ingredient for autonomous driving
- Large datasets play key role
- Membership inference:
  - Privacy violation
  - Model forensics

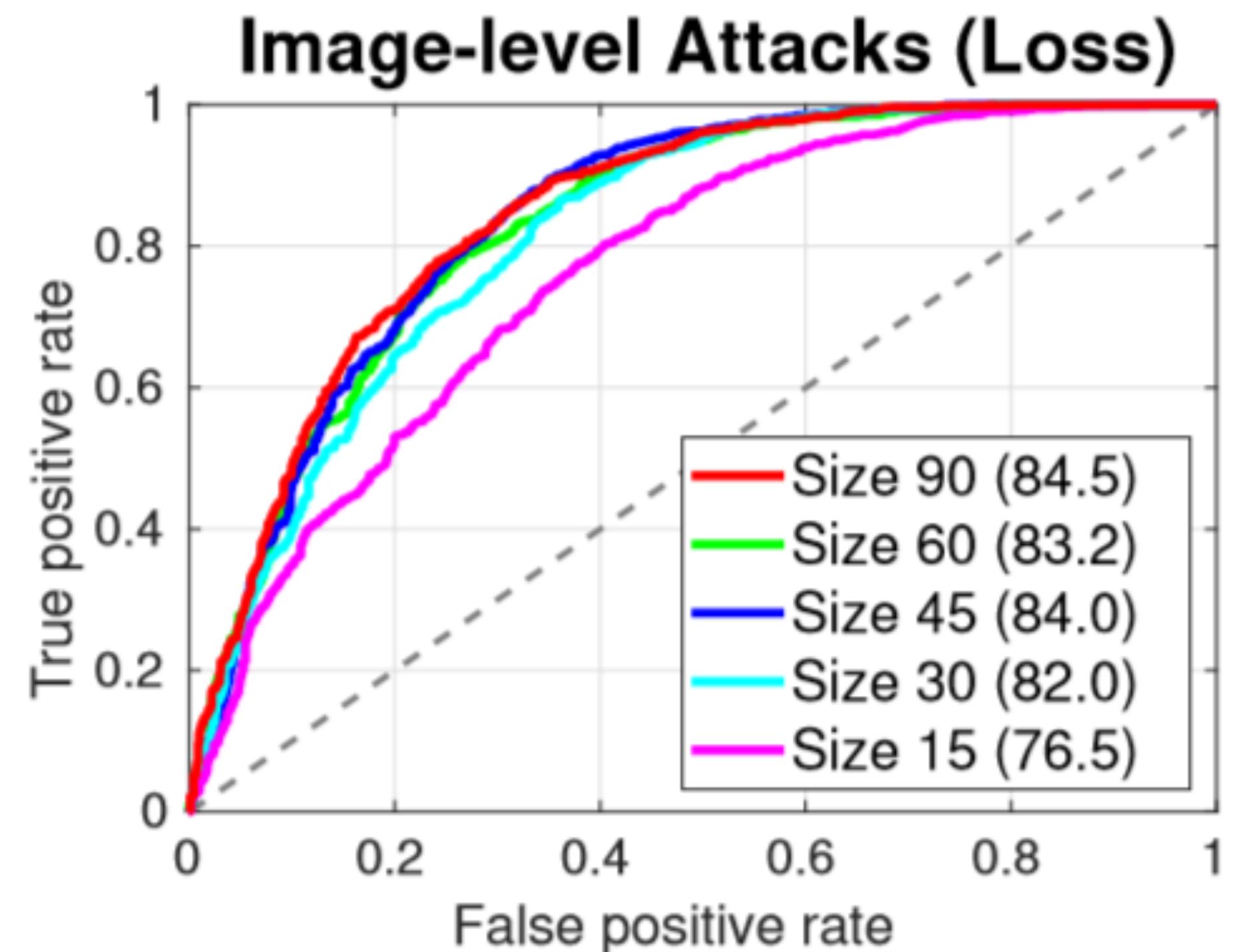


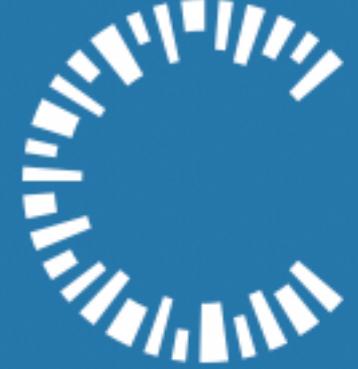
# Test Loss Map as Membership “Signal”



Structured loss map

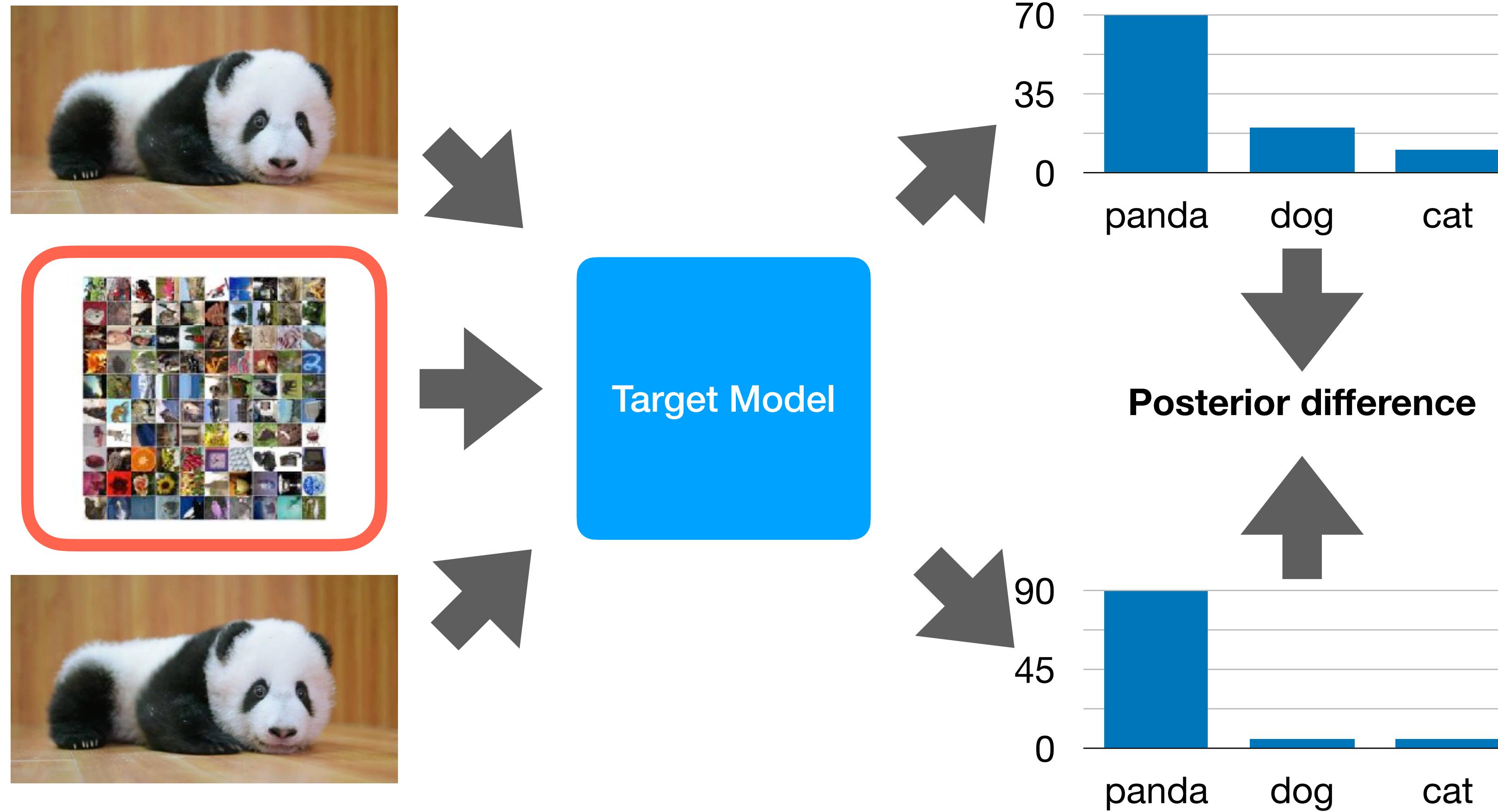
# Results





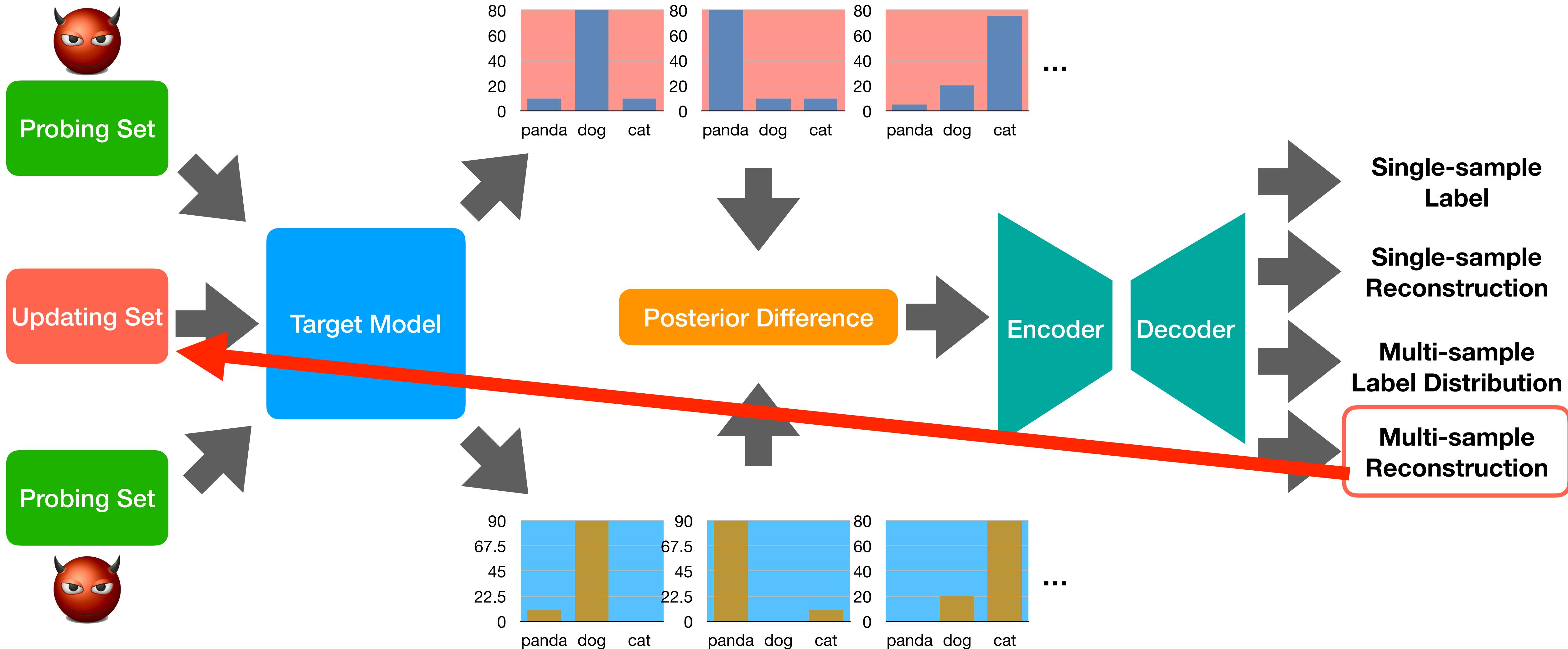
- Membership Inference can tell if a data point was in training set
- Classification, Segmentation, Generators
- ... but maybe it's not so bad. Or is it?

# A New Attack Surface for Online Learning

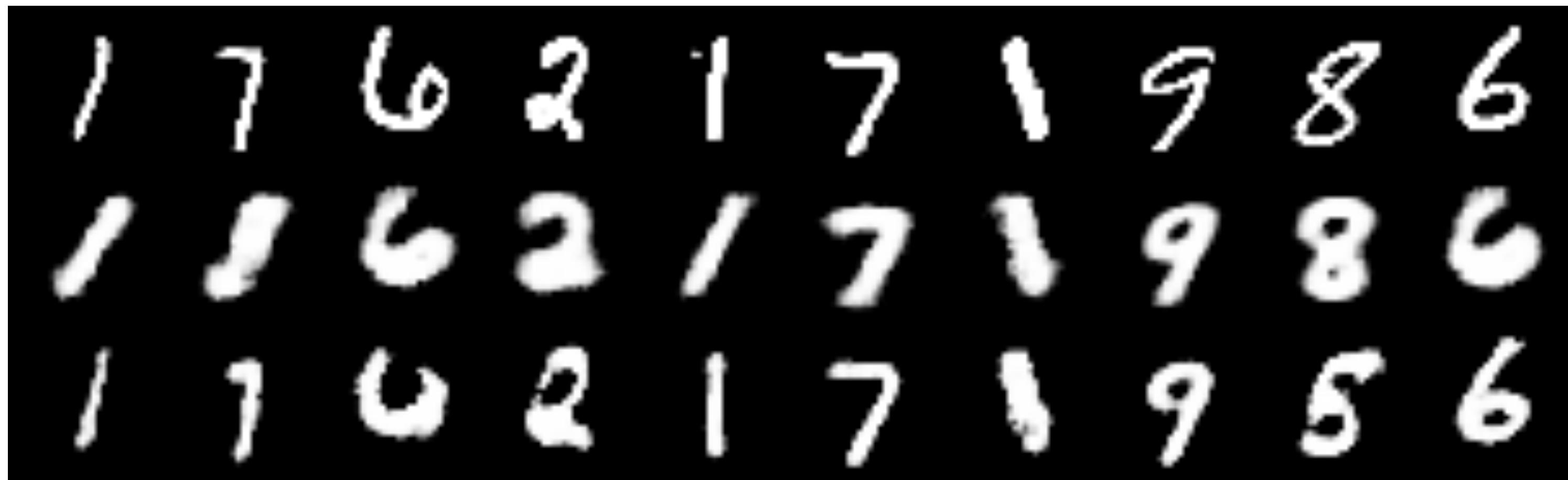
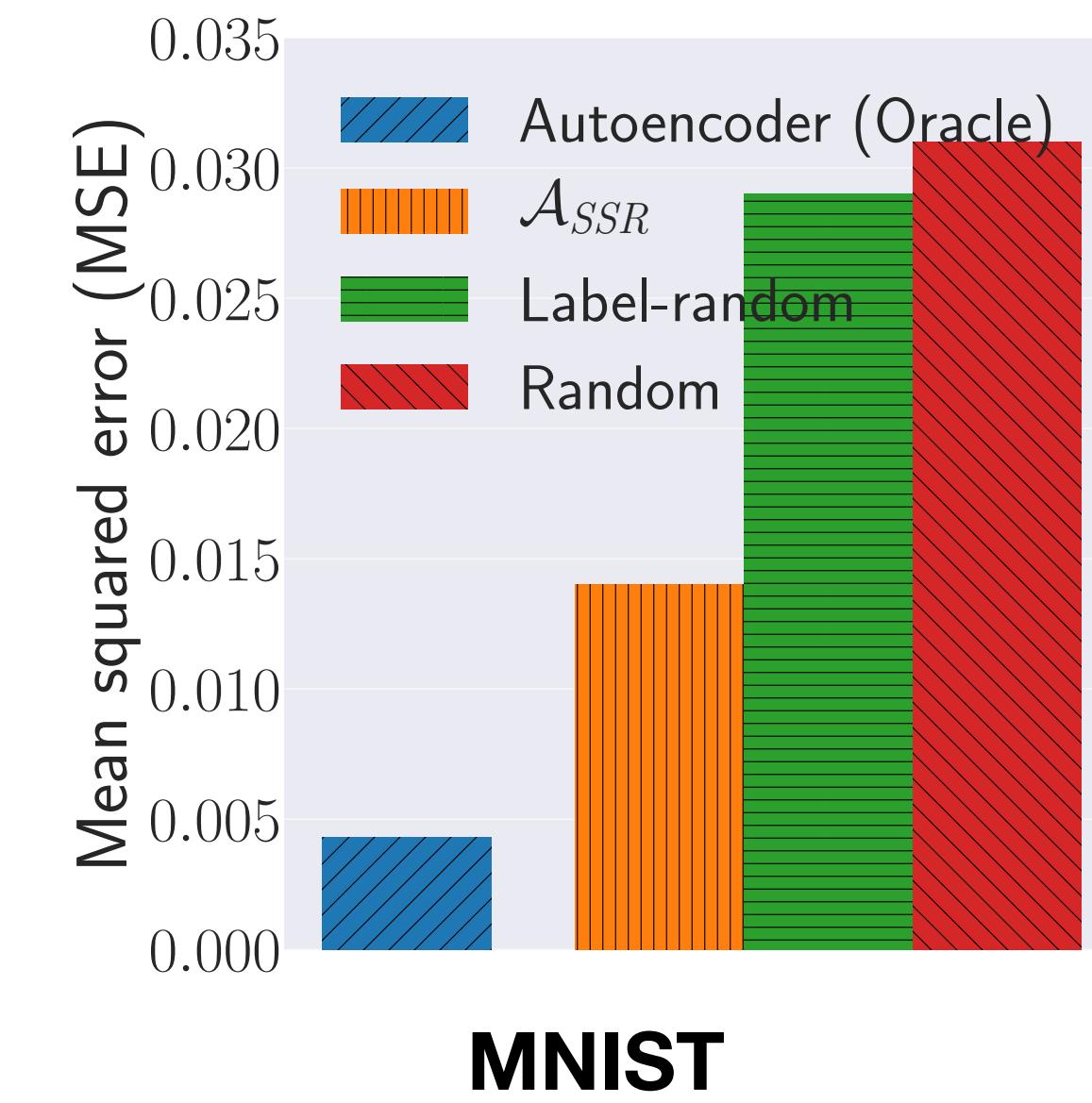
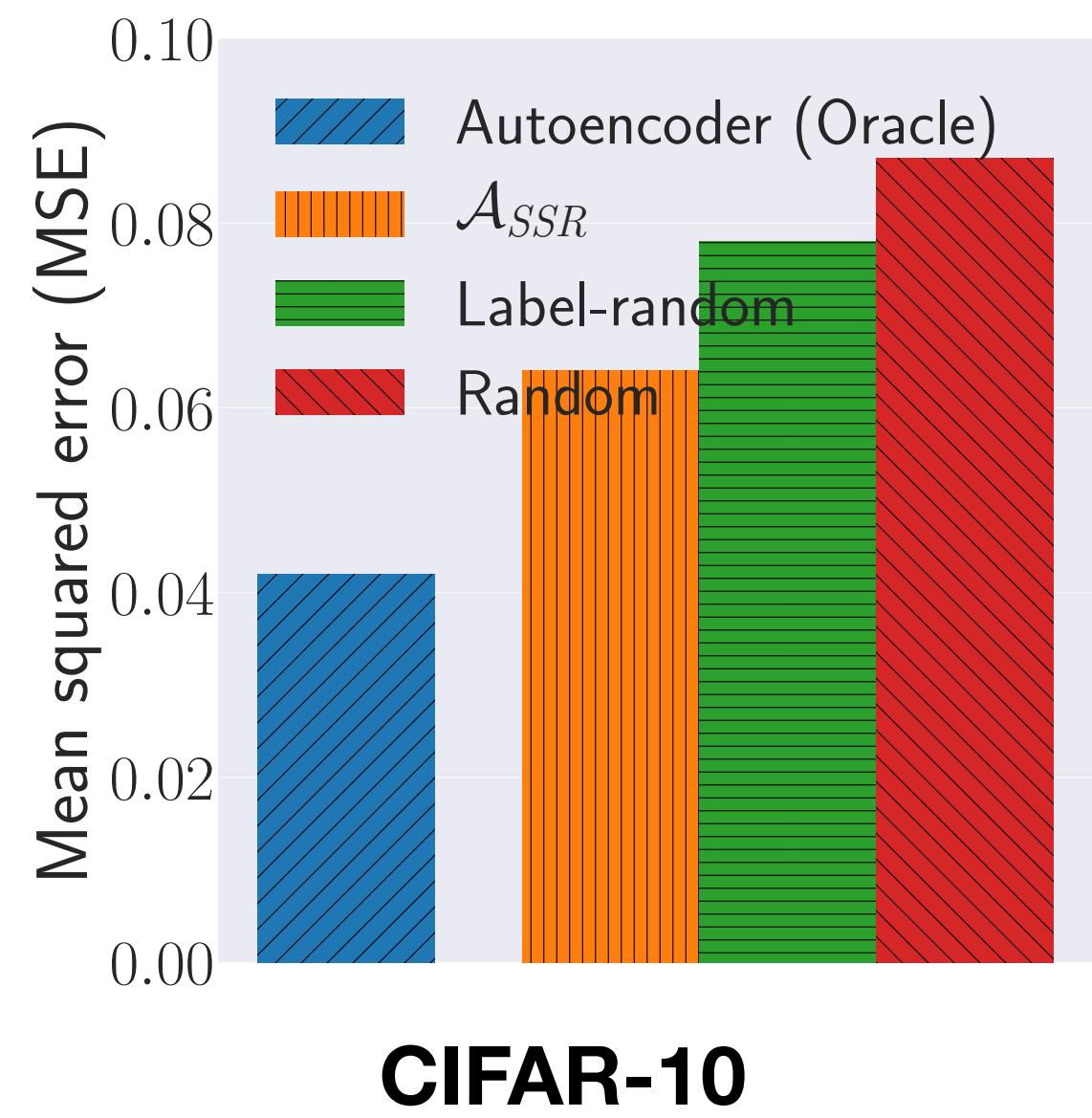


**Research Question:** Can this posterior difference be used to infer information of the updating set?

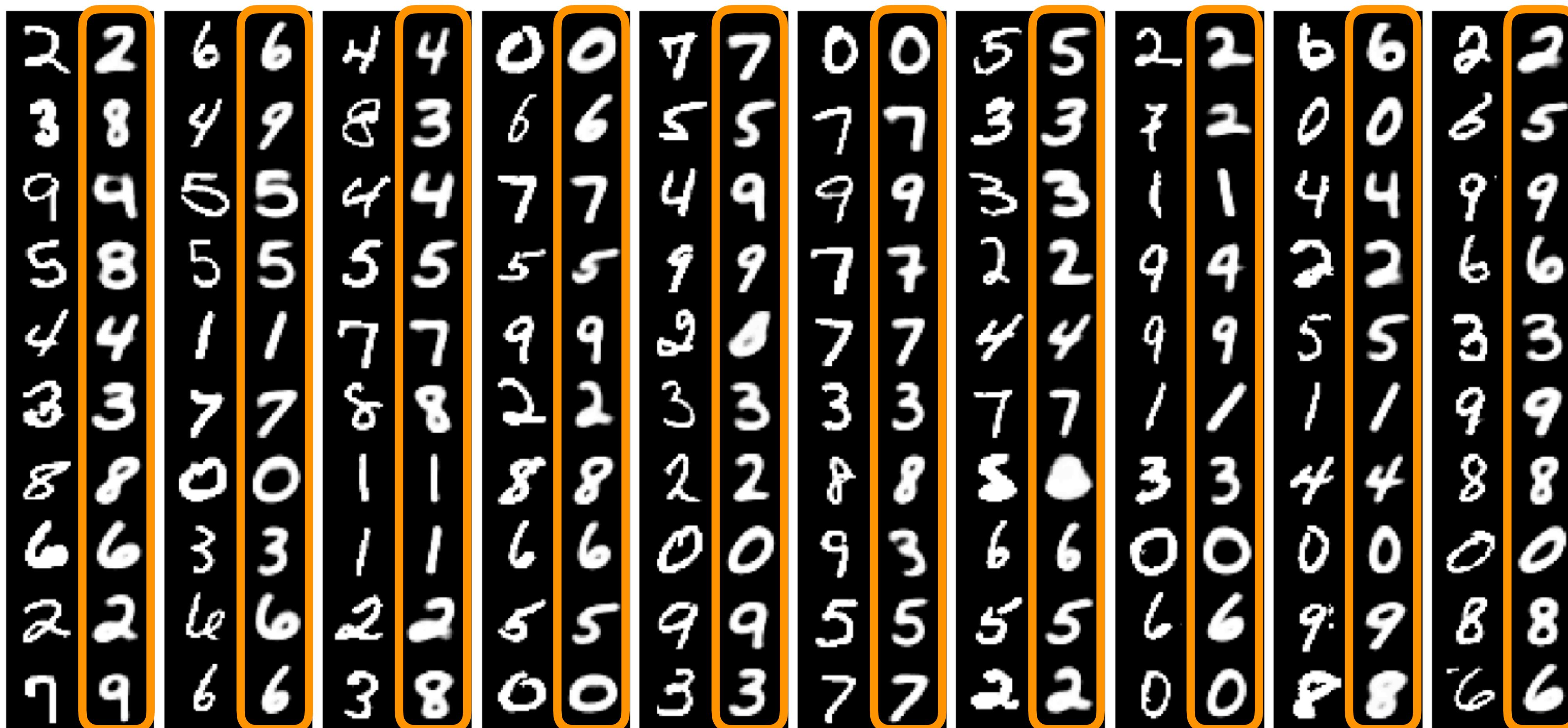
# General Attack Pipeline

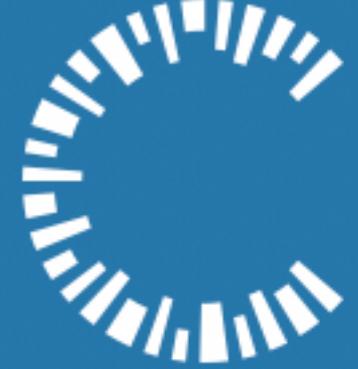


# Single-Sample Reconstruction



# Reconstruction of MNIST

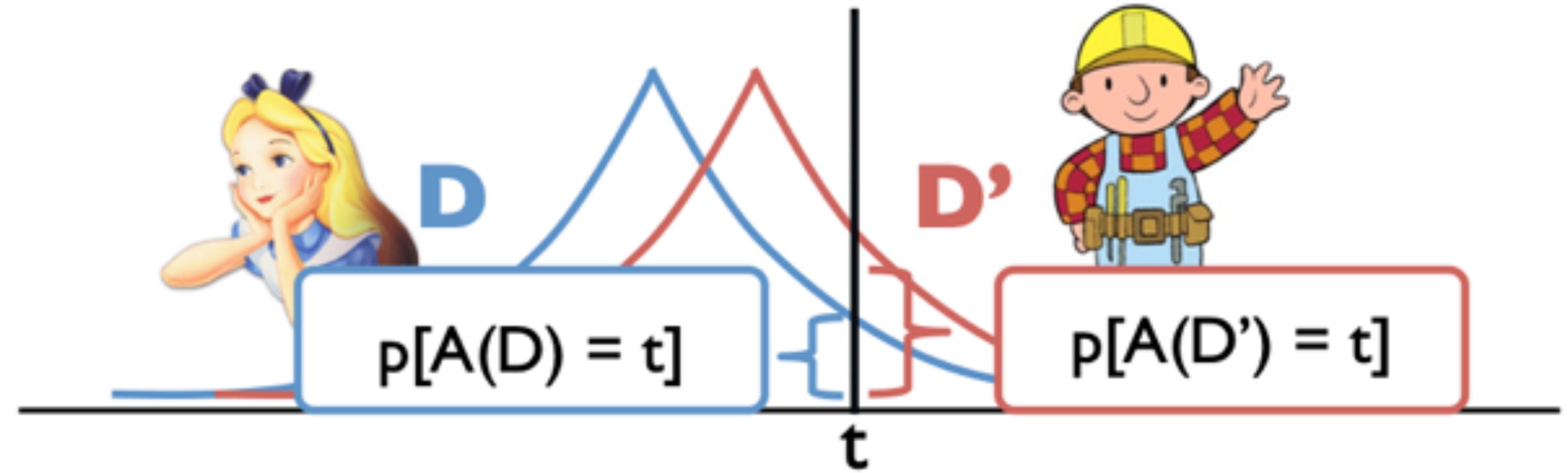




- Under some assumptions - reconstructing the training is possible.
- How can we prevent privacy leakage?

# Defenses - Privacy Preserving Machine Learning

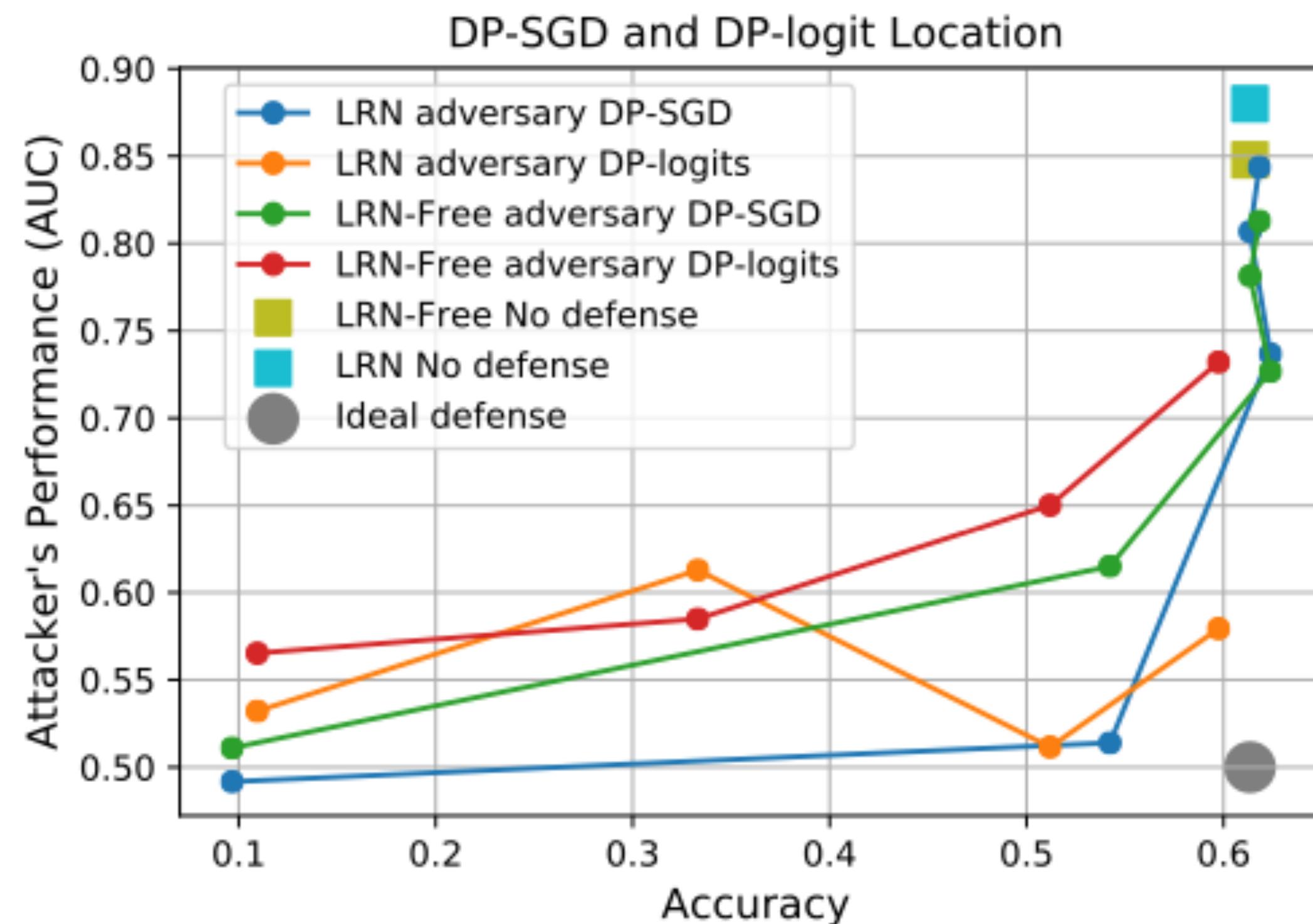
- Membership Inference relies on overfitting
- Defenses by drop-out, ensembles, proper model training, return only label
- No guarantees
- Principled defenses by differential privacy (DP)
- Algorithm can be transformed to satisfy DP, e.g. DP-SGD



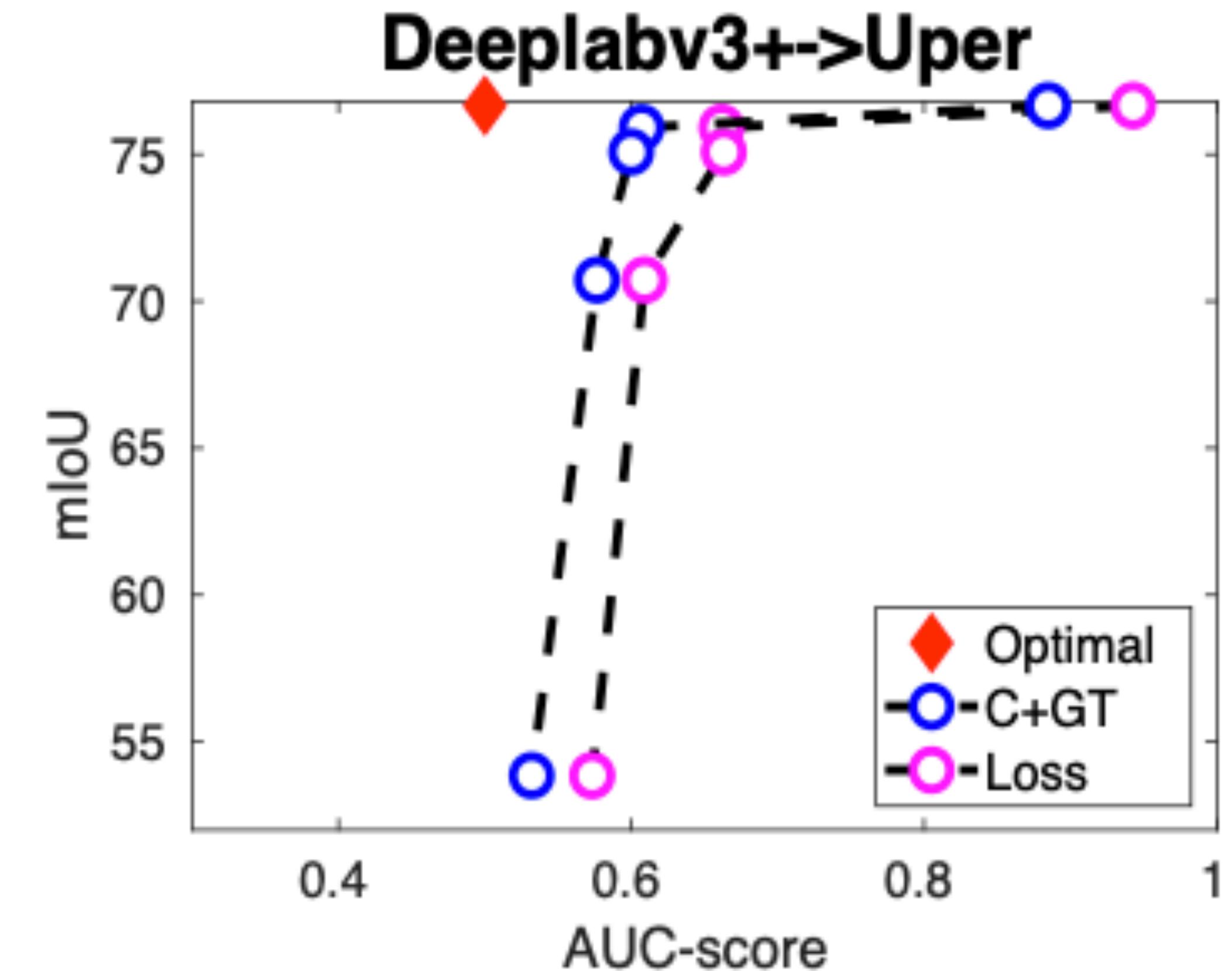
For all  $D, D'$  that differ in one person's value,  
If  $A = \epsilon$ -differentially private randomized algorithm, then:

$$\sup_t \left| \log \frac{p(A(D) = t)}{p(A(D') = t)} \right| \leq \epsilon$$

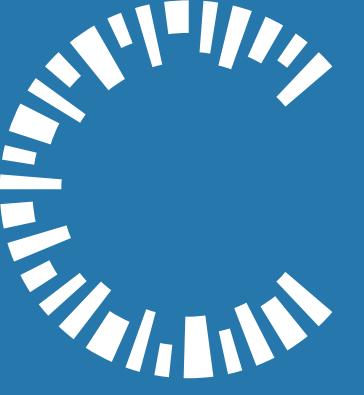
# Differential Privacy defends against Membership Inference



Classification



Segmentation



# Machine Learning in Cybersecurity

- Adversarial Perturbations & Game Theory
- Membership Inference Attacks

Prof. Dr. Mario Fritz | 28.11.2019