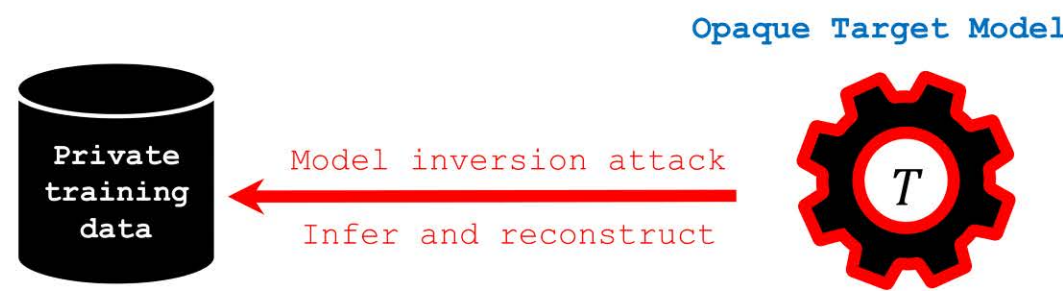


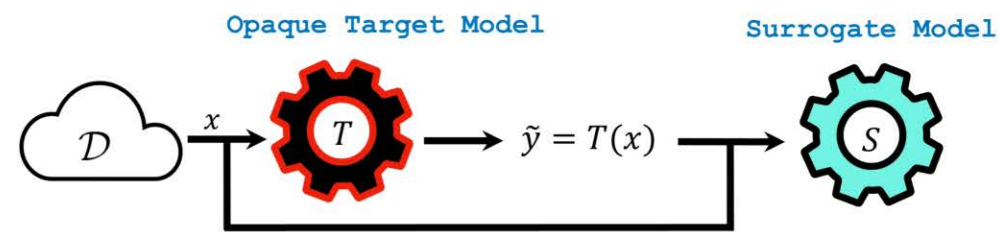
Label-only Model Inversion (MI)

Label-only Model inversion (MI) attacks aim to infer and reconstruct private training data by abusing access to a model's predicted label (hard label) without confidence scores nor any other model information.

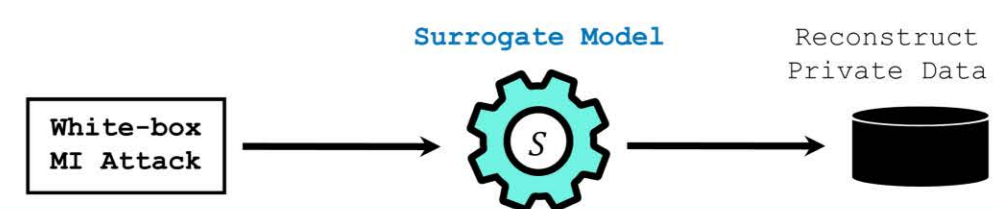


Label-only MI via Knowledge Transfer (LOKT)

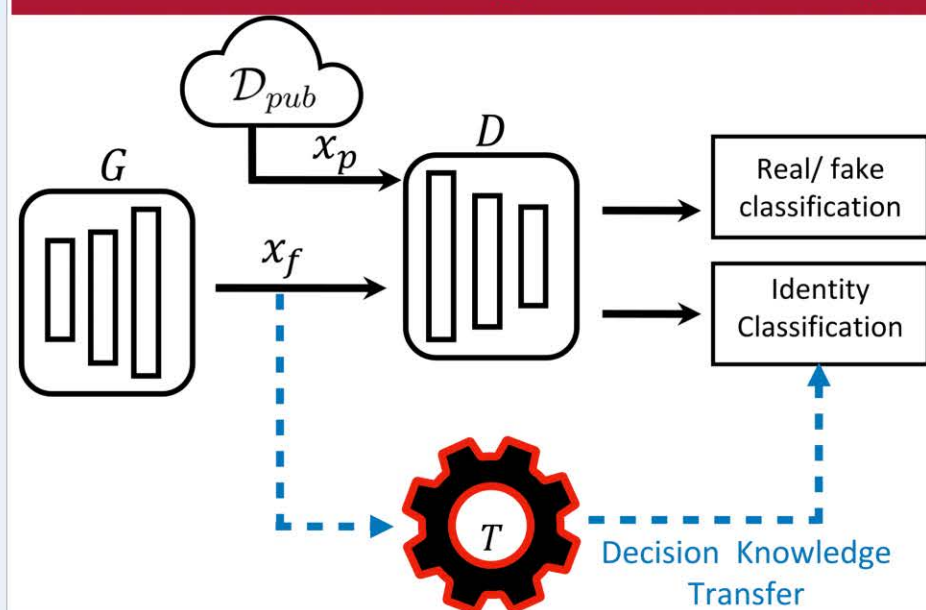
Stage 1: Decision Knowledge Transfer



Stage 2: Casting Label-only MI Attack as a White-box MI Attack



Decision Knowledge Transfer using our Target model-assisted ACGAN



We propose a new T-ACGAN to leverage generative modeling and the target model for effective knowledge transfer.

$$\mathcal{L}_{D,C} = -E[\log P(s = Fake|x_f)] - E[\log P(s = Real|x_p)] - E[\log P(c = \hat{y}|x_f)]$$

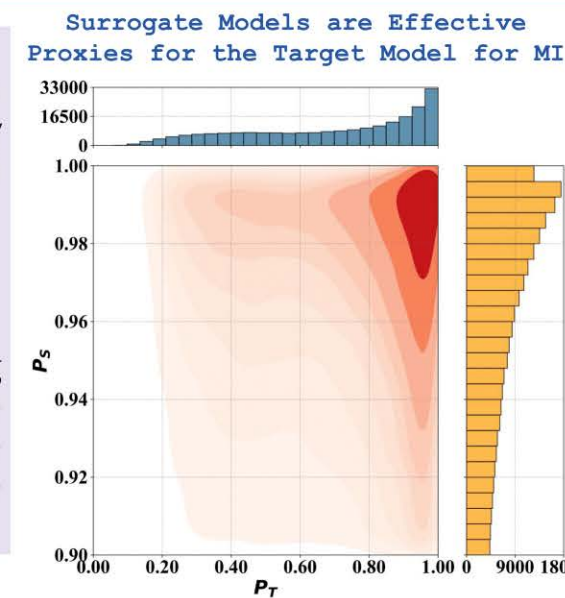
With T-ACGAN capturing the data manifold of public samples, synthetic data is diverse and abundant. We hypothesize that such rich synthetic data could lead to improved decision knowledge transfer.

Analysis for justification of surrogate models

Property P1:

For high-likelihood samples under S, it is likely that they also have high likelihood under T.

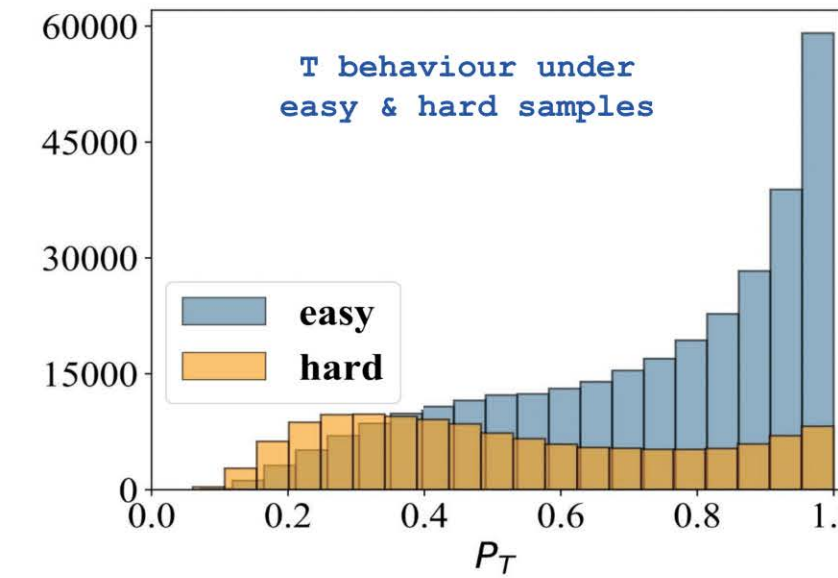
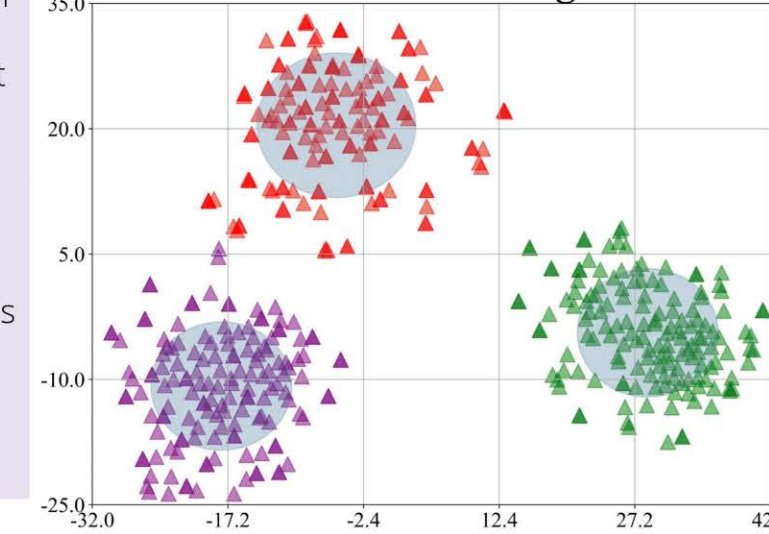
P1 is essential for performing MI attacks on S because during the inversion, we are finding samples that have high probability under T.



The general learning dynamics of DNNs:

- "Easy samples" are ones that fit better some patterns in the data (and correspondingly "hard samples").
- In DNNs learning, the models learn simple and general patterns of the data first in the training stage to fit the easy samples.

Face Embeddings



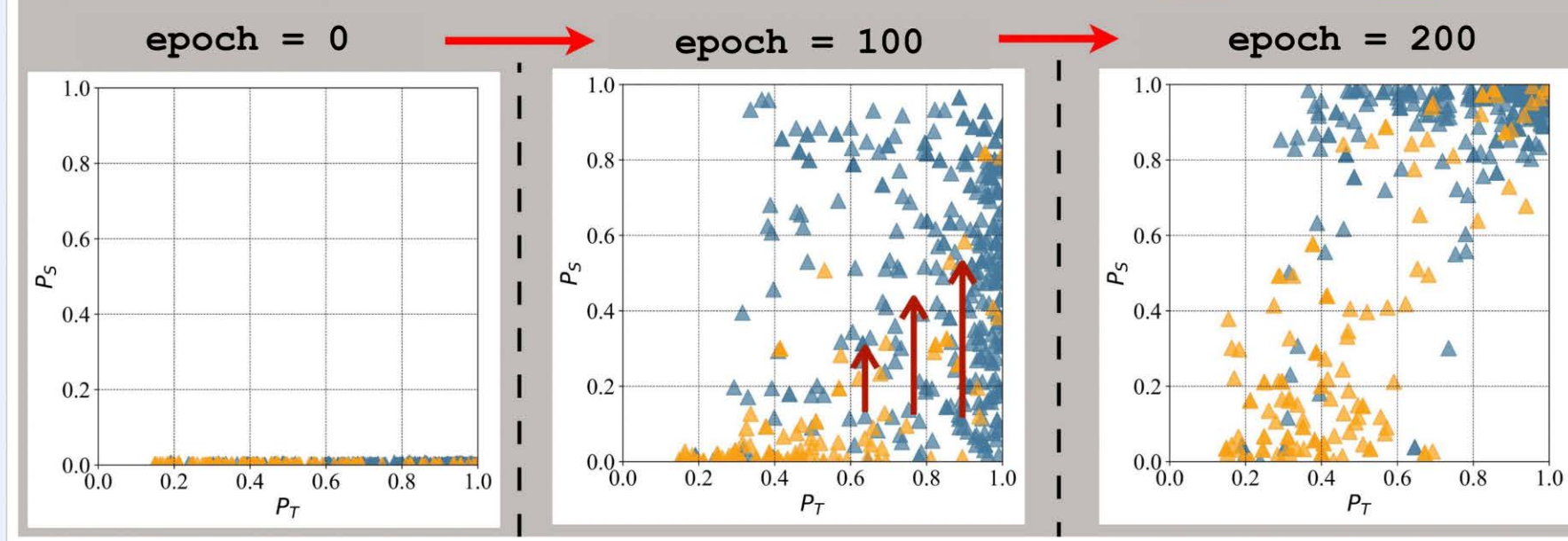
Easy samples



Hard samples



Training Dynamics of Easy & Hard samples



Main Results

Private Training Data	Attack	Attack Acc. (↑)	KNN dt. (↓)
Existing SOTA	BREPMI	73.93%	1284.41
Our LOKT	LOKT	93.93%	1181.72

Setup	Attack	Attack acc. ↑	KNN dt. ↓
$T = \text{FaceNet64}$	BREPMI	73.93 ± 4.98	1284.41
$D_{priv} = \text{CelebA}$	$C \circ D$	81.00 ± 4.79	1298.63
$D_{pub} = \text{CelebA}$	LOKT	92.80 ± 2.59	1207.25
	S_{en}	93.93 ± 2.78	1181.72
$T = \text{IR152}$	BREPMI	71.47 ± 5.32	1277.23
$D_{priv} = \text{CelebA}$	$C \circ D$	72.07 ± 4.03	1358.94
$D_{pub} = \text{CelebA}$	LOKT	89.80 ± 2.33	1220.00
	S_{en}	92.13 ± 2.06	1206.78
$T = \text{VGG16}$	BREPMI	57.40 ± 4.92	1376.94
$D_{priv} = \text{CelebA}$	$C \circ D$	71.33 ± 4.39	1364.47
$D_{pub} = \text{CelebA}$	LOKT	85.60 ± 3.03	1252.09
	S_{en}	87.27 ± 1.97	1246.71
$T = \text{FaceNet64}$	BREPMI	43.00 ± 5.14	1470.55
$D_{priv} = \text{CelebA}$	$C \circ D$	43.27 ± 3.53	1516.18
$D_{pub} = \text{FFHQ}$	LOKT	59.13 ± 2.77	1437.86
	S_{en}	62.07 ± 3.89	1428.04
$T = \text{BiDO}$	BREPMI	37.40 ± 3.66	1500.45
$D_{priv} = \text{CelebA}$	$C \circ D$	45.73 ± 5.94	1493.48
$D_{pub} = \text{CelebA}$	LOKT	58.53 ± 4.87	1427.22
	S_{en}	60.73 ± 3.07	1395.93

Contributions

- We propose a new approach for **Label-Only MI attack using Knowledge Transfer (LOKT)** by transferring decision knowledge from the target model to surrogate models and performing white-box attacks on the surrogate models. Our proposed approach is the first to address label-only MI via white-box attacks.
- We propose a new **Target model-assisted ACGAN (T-ACGAN)** to leverage generative modeling and the target model for effective knowledge transfer.
- We perform **analysis** to support that our surrogate models are effective proxies for the target model for MI.

