

Label-Only Model Inversion Attacks via Knowledge Transfer

Ngoc-Bao Nguyen^{*1}

Keshigeyan Chandrasegaran^{*2‡}

Milad Abdollahzaden¹




Ngai-Man Cheung¹

¹ Singapore University of Technology and Design (SUTD)

² Stanford University

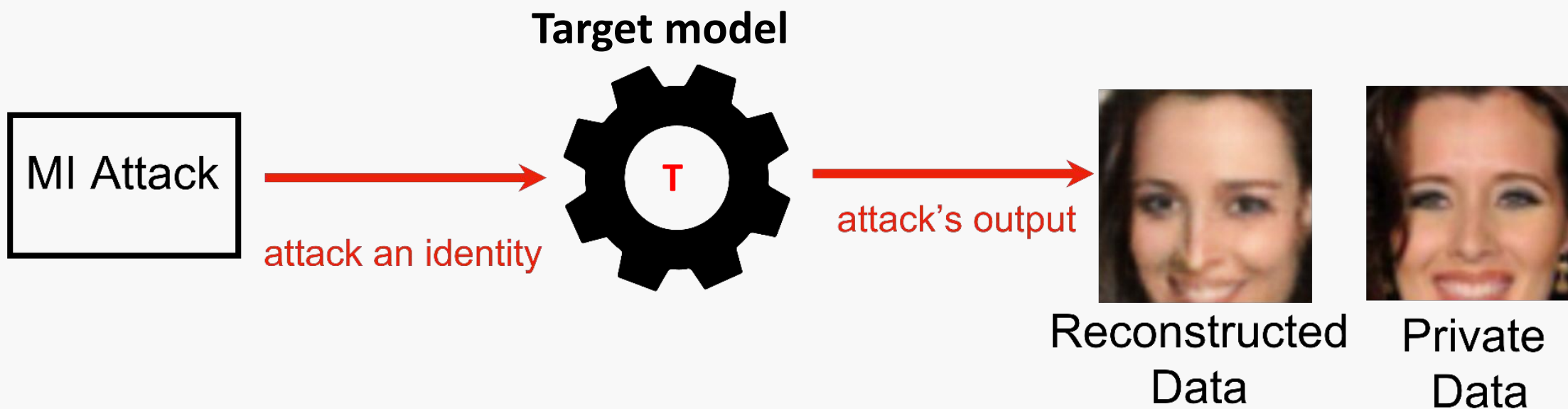
Our contributions

- We propose **Label-only Model inversion via Knowledge Transfer (LOKT)**, a new label-only MI by transferring decision knowledge from the target model to surrogate models and performing white-box attacks on the surrogate models.
- We propose a new 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.

<i>Private Training Data</i>		Attack Acc. (↑)
<i>Existing SOTA</i>		73.93%
Our Reconstruction Results		93.93%

Model Inversion (MI)

Model inversion (MI) attacks aim to infer and reconstruct private training data by abusing access to a model.

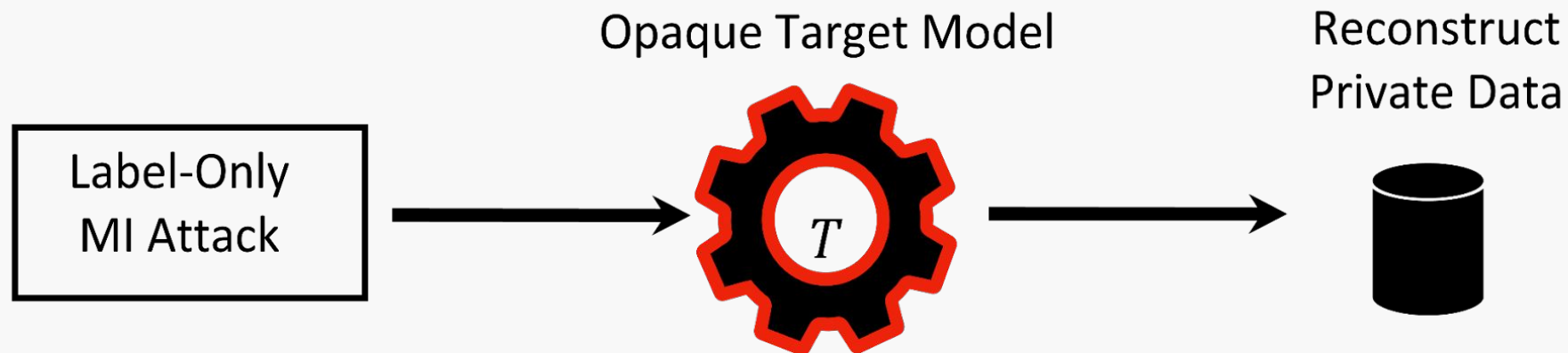


Model Inversion (MI)

We focus on **label-only model inversion attack** which is the most challenging MI Attack.

Criteria	Architecture / Parameters	Soft-labels	Hard-labels	Concern reg. Queries
White-box	✓	✓	✓	Low
Black-box	✗	✓	✓	High
Label-only	✗	✗	✓	High

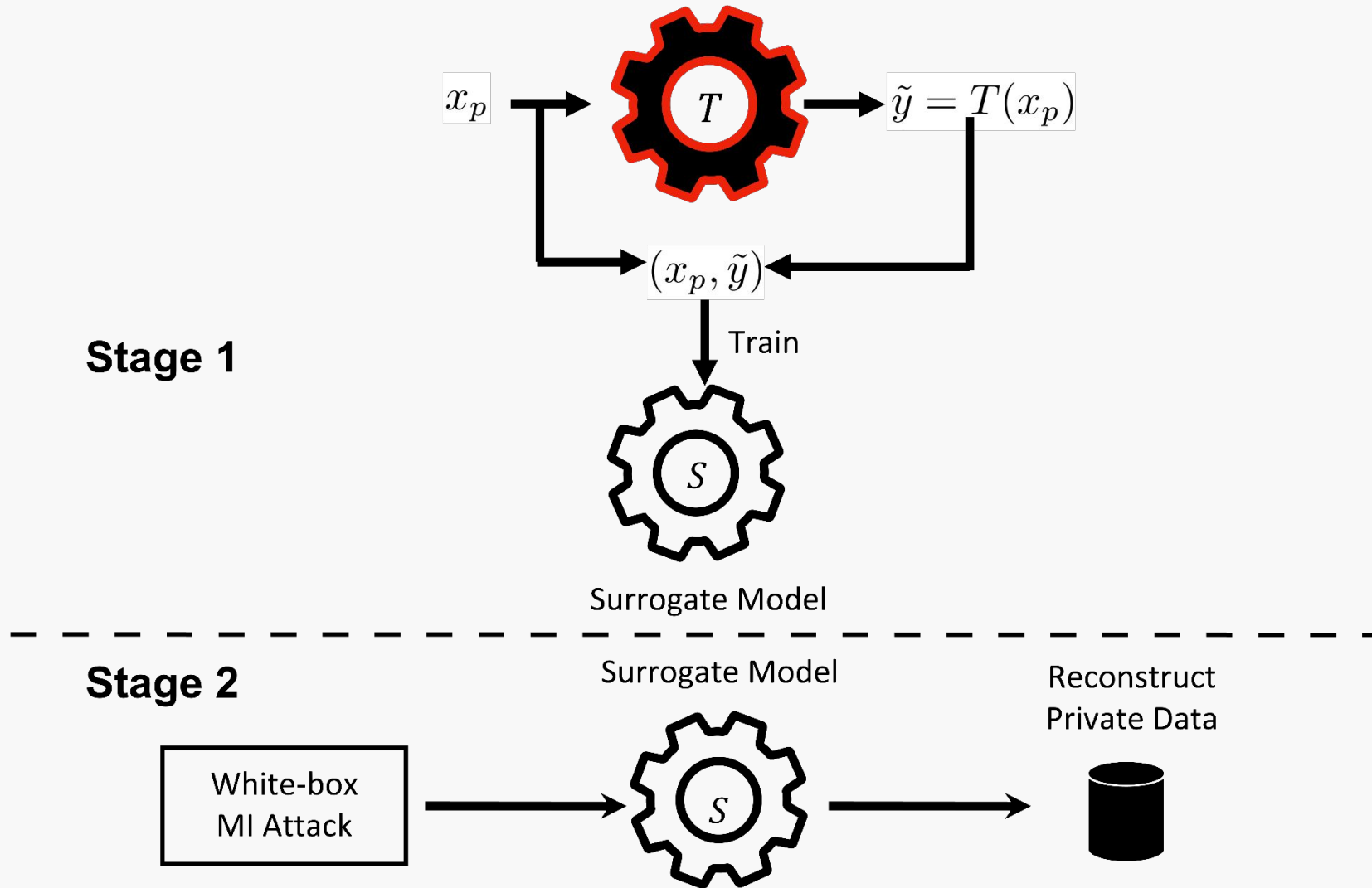
Existing work on Label-only Model Inversion Attack



SOTA Label-only Model Inversion attacks employ **black-box search on the target model T** to reconstruct training data.

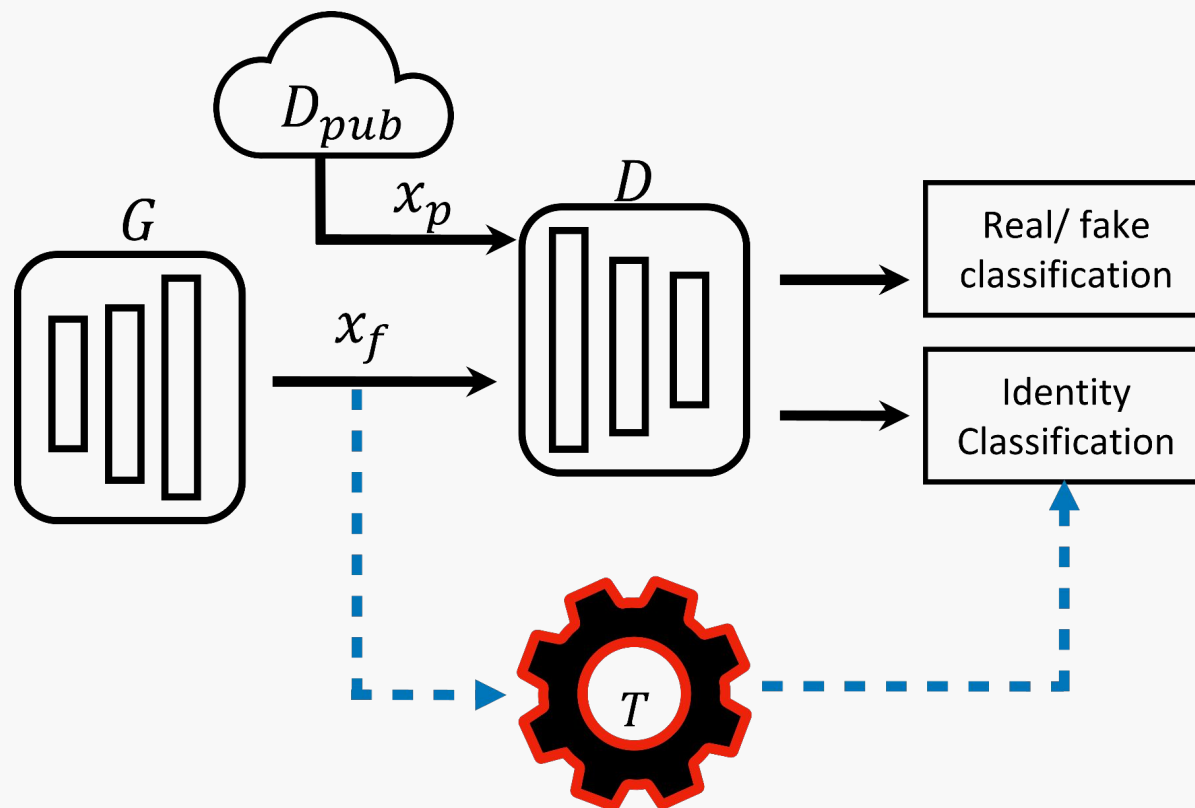
Label-only Model inversion via Knowledge Transfer (LOKT)

Decision Knowledge Transfer



Casting Label-only MI Attack as a White-box MI Attack

Decision Knowledge Transfer using our T-ACGAN



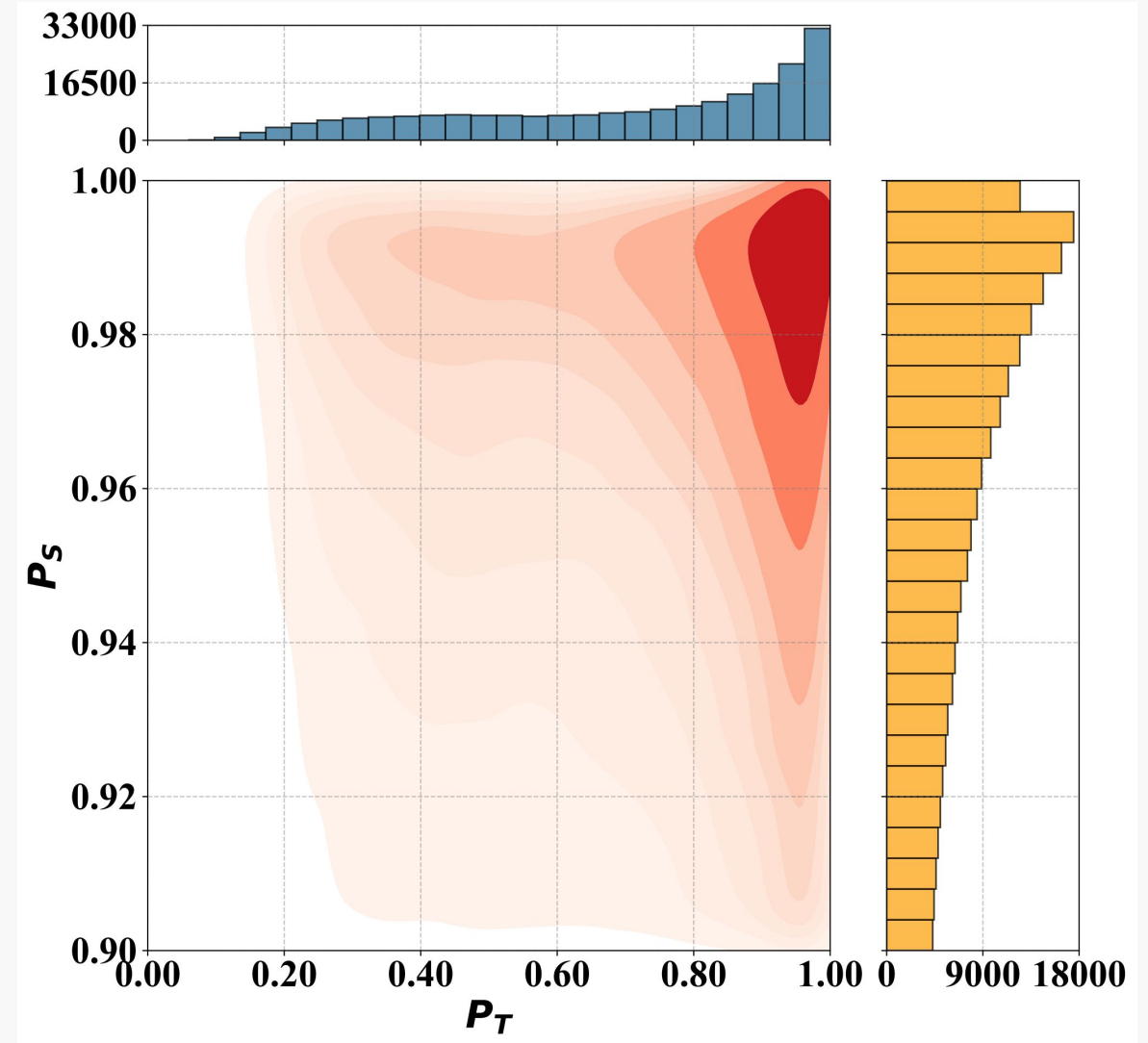
Decision Knowledge Transfer

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

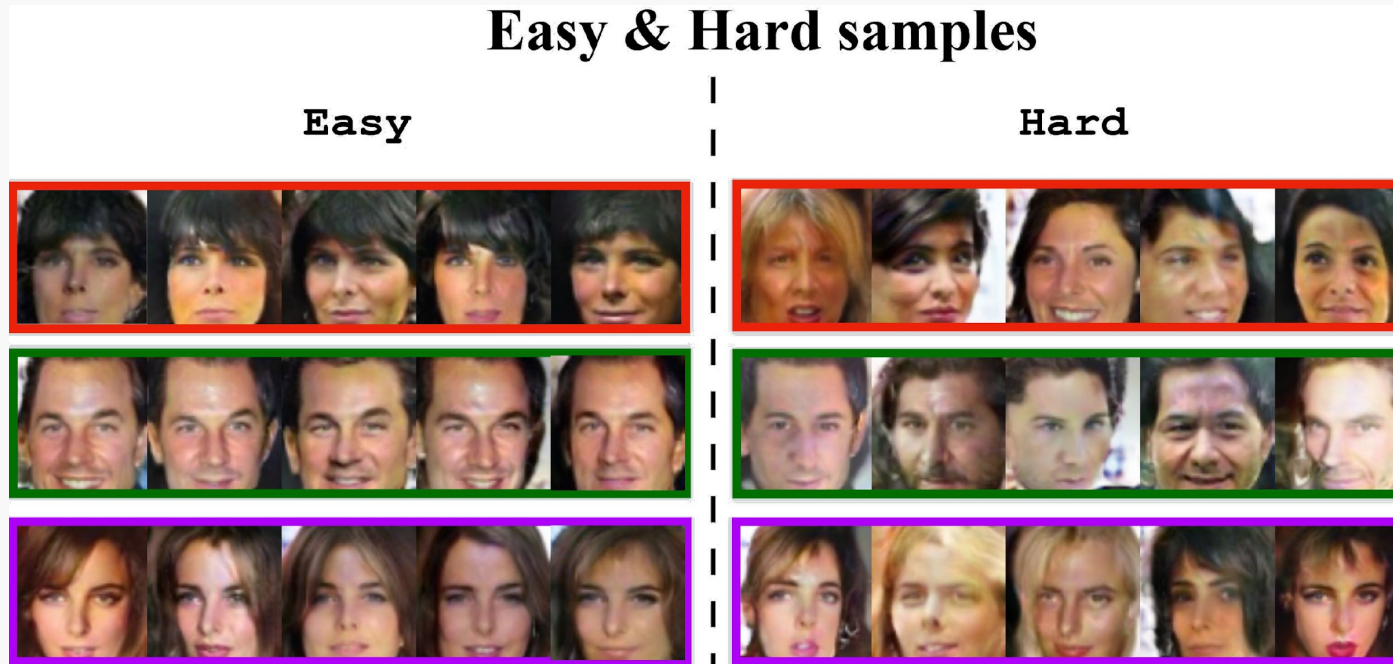
Analysis for justification of surrogate models

Property P1:

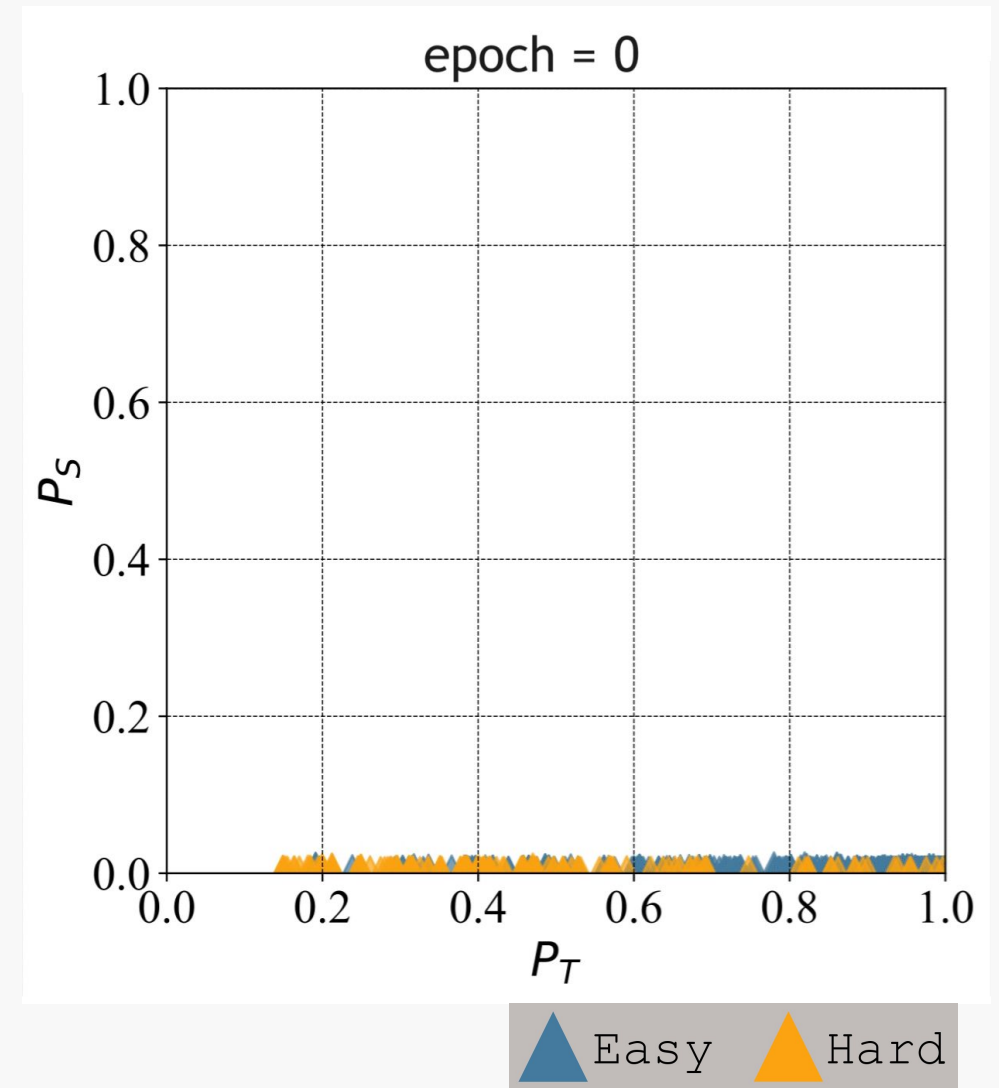
For high-likelihood samples under S , it is likely that they also have high likelihood under T .



Analysis for justification of surrogate models



DNNs Learn Patterns First



Model inversion attack results

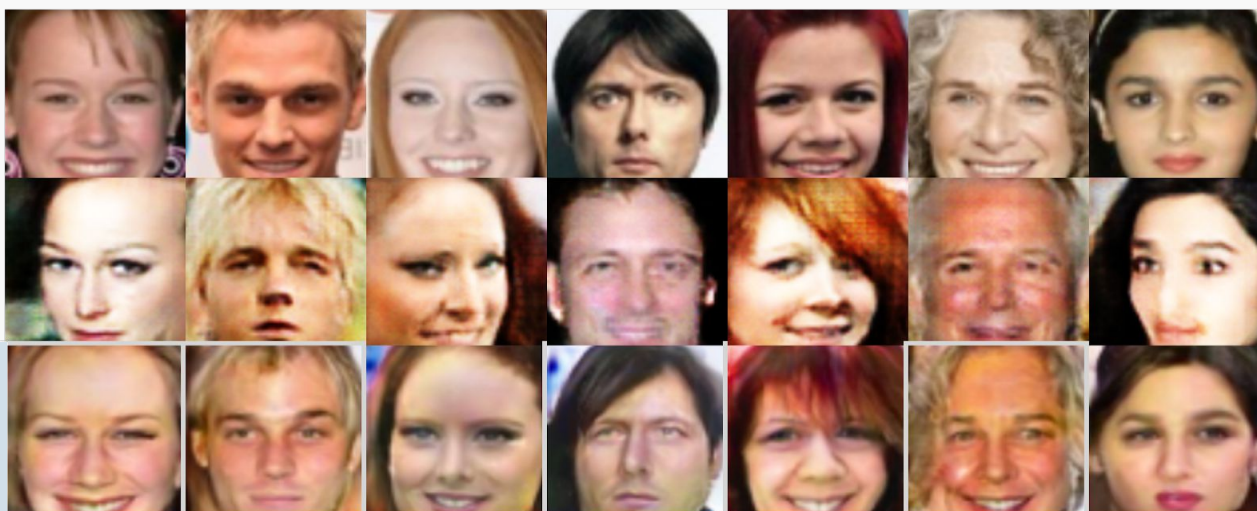
Setup		Attack	Attack acc. \uparrow	KNN dt. \downarrow
T = FaceNet64 \mathcal{D}_{priv} = CelebA \mathcal{D}_{pub} = CelebA		BREPMI	73.93 ± 4.98	1284.41
		$C \circ D$	81.00 ± 4.79	1298.63
		LOKT S	92.80 ± 2.59	1207.25
		S_{en}	93.93 ± 2.78	1181.72
T = IR152 \mathcal{D}_{priv} = CelebA \mathcal{D}_{pub} = CelebA		BREPMI	71.47 ± 5.32	1277.23
		$C \circ D$	72.07 ± 4.03	1358.94
		LOKT S	89.80 ± 2.33	1220.00
		S_{en}	92.13 ± 2.06	1206.78

Setup		Attack	Attack acc. \uparrow	KNN dt. \downarrow
T = VGG16 \mathcal{D}_{priv} = CelebA \mathcal{D}_{pub} = CelebA		BREPMI	57.40 ± 4.92	1376.94
		$C \circ D$	71.33 ± 4.39	1364.47
		LOKT S	85.60 ± 3.03	1252.09
		S_{en}	87.27 ± 1.97	1246.71
T = FaceNet64 \mathcal{D}_{priv} = CelebA \mathcal{D}_{pub} = FFHQ		BREPMI	43.00 ± 5.14	1470.55
		$C \circ D$	43.27 ± 3.53	1516.18
		LOKT S	59.13 ± 2.77	1437.86
		S_{en}	62.07 ± 3.89	1428.04

Private
Training
Data

Existing
SOTA

Our
Reconstruction
Results



Attack
Acc. (\uparrow)

73.93%

93.93%

Conclusion

- We propose **Label-only Model inversion via Knowledge Transfer (LOKT)**, a new label-only MI by transferring decision knowledge from the target model to surrogate models and performing white-box attacks on the surrogate models.
- We propose a new 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.

Thank you!

Poster Session

Wed 13 Dec 10:45 a.m. CST — 12:45 p.m. CST

Great Hall & Hall B1+B2

#1517

Project page



<https://ngoc-nguyen-0.github.io/lokt/>