

## Introduction

- Transfer learning has become an increasingly popular technique in private machine learning as a way to leverage a model trained for one task to assist with building a model for a related task.
- Because pretraining data is often considered to be public, it provides a good initialization for sensitive downstream tasks.
- Unfortunately, large, public datasets are typically scraped from the Web indiscriminately, raising concerns about the sensitivity of this data.

Thus, the central question we attempt to understand in this work is: How much sensitive information does a finetuned model reveal about the data that was used for pretraining?

# Membership Inference Threat Model



# **TMI! Finetuned Models Spill** Secrets from Pretraining

John Abascal, Stanley Wu, Alina Oprea, Jonathan Ullman Northeastern University



Results

## Conclusion





### • Finetuning leaves pretraining examples at risk of membership inference attacks.

• If an individual's data was included in a public pretraining dataset, finetuning on this individual's data with DP will **not** protect them from membership inference attacks.

## nbership

ized hip by making abels in the

TMI TMI Adapt Adapt LiRA

TMI - Coa TMI - Co Adapted Adapted DP Upper DP Upper





**Logit Distributions of Finetuned Models Queried on Pretraining Data** 



#### Finetuned on Coarse-Labeled CIFAR-100 with **DP** ( $\varepsilon = 0.5, \delta = 10^{-5}$ ) when $D_{PT} \cap D_{FT} \neq \emptyset$

Task	TPR @ 0.1% FPR	TPR @ 1% FPR
(Coarse CIFAR-100)	5.7%	16.1%
(CIFAR-10)	2.0%	8.0%
ted LiRA (Coarse CIFAR-100)	0.7%	3.1%
ted LiRA (CIFAR-10)	0.3%	1.5%
Directly on Pretrained Model	15.6%	22.9%

Task	TPR @ 0.1% FPR	TPR @ 1% FPR
parse CIFAR-100; ( $\epsilon = 0.5, \delta = 10^{-5}$ )	2.6%	8.5%
parse CIFAR-100; ( $\epsilon = 1, \delta = 10^{-5}$ )	4.2%	12.6%
LiRA - Coarse CIFAR-100 ( $\varepsilon = 0.5, \delta = 10^{-5}$ )	0.17%	1.8%
LiRA - Coarse CIFAR-100 ( $\varepsilon = 1, \delta = 10^{-5}$ )	0.35%	2.8%
r Bound ( $\epsilon = 0.5, \ \delta = 10^{-5}$ )	0.16%	1.6%
r Bound ( $\varepsilon = 1, \ \delta = 10^{-5}$ )	0.27%	2.7%

#### **Attack Performance**