

# A Method to Facilitate Membership Inference Attacks in Deep Learning Models

Zitao Chen, Karthik Pattabiraman University of British Columbia



THE UNIVERSITY OF BRITISH COLUMBIA

## Membership inference attacks (MIAs)



#### Which data point was used to train a model?

#### MIAs as a privacy threat



health status

## A common thread of many studies



## ML democratization



# Code poisoning attacks are realistic threats

The **A**Register<sup>®</sup>

O PyTorch



PyTorch dependency poisoned with malicious code

#### The Hacker News

**TensorFlow CI/CD Flaw Exposed Supply Chain to** 

**Poisoning Attacks** 

#### The Hacker News

New Evidence Suggests SolarWinds' Codebase Was Hacked to Inject Backdoor

#### Privacy risk of using untrusted ML codebase in development

# Threat model



## Prior attacks: How to increase privacy leakage





Tramer et al.,  $CCS'22 \longrightarrow Data$  poisoning Chen et al, NeurIPS'22  $\longrightarrow$  Code poisoning Song et al., AsiaCCS'21  $\longrightarrow$  Code poisoning

#### Prior attacks

#### Trade-off between privacy and utility



Tramer et al., CCS'22

#### Prior attacks

How to overcome the trade-off between privacy and utility



#### This work: A new direction to construct high-power MIAs

## Our indirect attack: Divide and conquer



## Encode membership of *D*<sub>train</sub> via *D*<sub>secret</sub>





# How to make the (indirect) attack easy?

**Outlier** data are easy to memorize





# Challenge

#### Norm functions expect data from a similar distribution



# Solution: Divide and conquer (again)

A secondary norm func to separately process *D<sub>secret</sub>* ReLu 2<sup>nd</sup> BN BN **High model utility** High privacy leakage U conv





Tramer et al., CCS'22

Our attack exposes the worst-case privacy leakage has minimal performance impact can disguise high privacy leakage



Artifact Evaluated

Available

Functional

Reproduced

## Conclusion

Zitao Chen zitaoc@ece.ubc.ca

Using third-party **ML codebase** has hidden privacy risk

New direction to construct stealthy attacks and inflict worst-case leakage The <u>first</u> result → Existing privacy auditing methods can be unreliable!

