

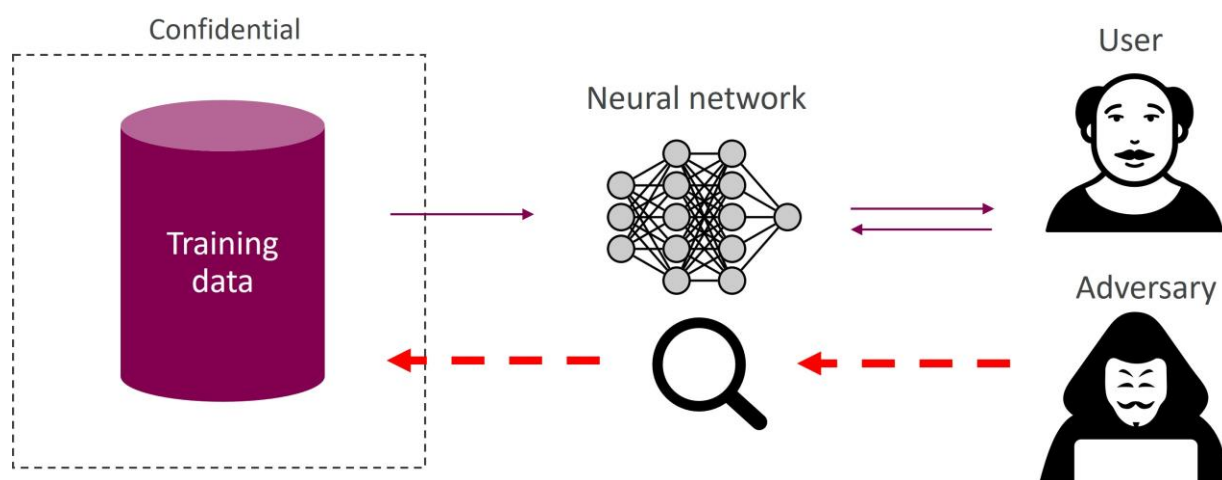
# Publishing Neural Networks in Drug Discovery Might Compromise Training Data Privacy

Fabian Krüger<sup>1,2,3</sup>, Johan Östman<sup>4</sup>, Lewis Mervin<sup>1</sup>, Igor Tetko<sup>3</sup>, Ola Engkvist<sup>1,5</sup>

AstraZeneca R&D<sup>1</sup>, Technical University of Munich<sup>2</sup>, Helmholtz Munich<sup>3</sup>, AI Sweden<sup>4</sup>, Chalmers University of Technology<sup>5</sup>

## Abstract

This study investigates the risks of exposing confidential chemical structures when machine learning models trained on these structures are made publicly available. We use membership inference attacks, a common method to assess privacy that is largely unexplored in the context of drug discovery, to examine neural networks for molecular property prediction in a black-box setting. Our results reveal significant privacy risks across all evaluated datasets and neural network architectures.

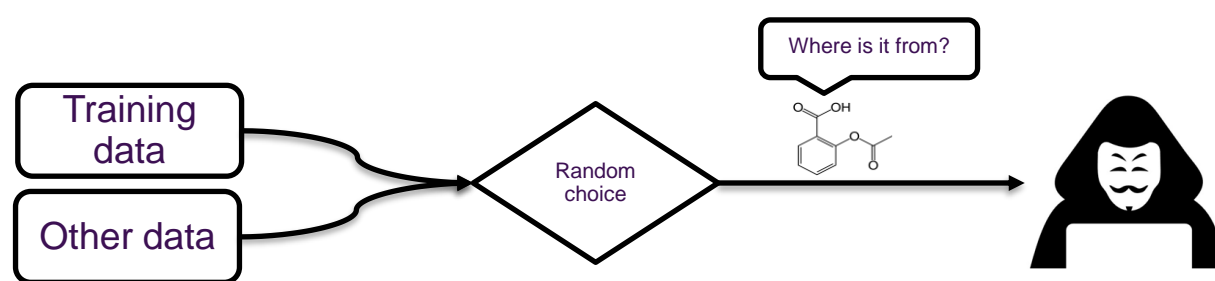


## Introduction

- Datasets in drug discovery are expensive to generate. Leaking information about proprietary data can severely harm an organization.
- Organizations need to balance benefits from open science and collaboration with the scientific community with their privacy concerns.
- There is a lack of studies on how much training data information can be inferred from neural networks in a drug discovery context.

## Methods

- Developed and evaluated neural networks trained on diverse molecular representations (fingerprints, graphs, SMILES) across four drug discovery datasets.
- Applied state-of-the-art membership inference attacks (LiRA<sup>1</sup> and RMIA<sup>2</sup>) in a black-box setting to measure how well attackers can identify molecules from training data.

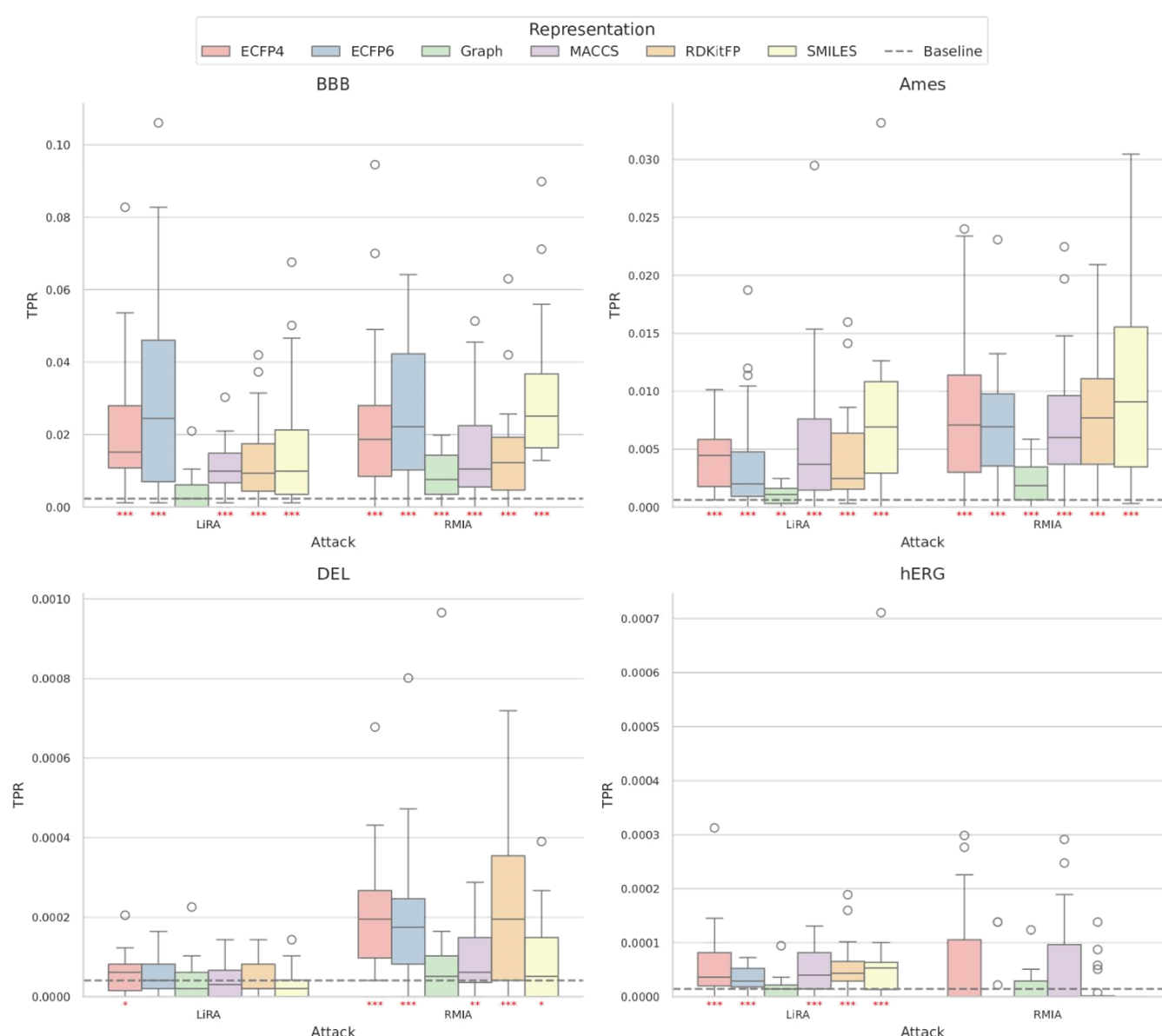


Scheme of how membership inference attacks are evaluated.

## Conclusions

- It is consistently possible to identify parts of the training data, even at false positive rates as low as 0 (under some assumptions).
- Combining both attacks allows getting even more information about the training data.
- Minority class molecules are easier to identify.
- Message passing neural network has the least information leakage.

## Results



True positive rates for identifying training data molecules at a false positive rate of 0. The distributions of 20 experimental repetitions are shown for each representation and dataset, for both the likelihood ratio attack (LiRA<sup>1</sup>) and the robust membership inference attack (RMIA<sup>2</sup>). Distributions with significantly higher true positive rates (information leakage) than the baseline (random guessing) are indicated by red stars. Training dataset sizes are: 859 molecules for the blood-brain barrier permeability dataset<sup>3</sup>; 3,264 for the Ames mutagenicity prediction dataset<sup>4</sup>; 48,837 for the DNA-encoded library enrichment dataset<sup>5</sup>; and 137,853 for the hERG channel inhibition dataset<sup>6</sup>.

Paper:



## References

1. Membership inference attacks from first principles, Carlini et al., IEEE, 2022.
2. Low-cost high-power membership inference attacks, Zarifzadeh et al., ICML, 2024
3. A bayesian approach to in silico blood-brain barrier penetration modeling, Martins et al., JCIM, 2012
4. Benchmark data set for in silico prediction of ames mutagenicity, Hansen et al., JCIM, 2009
5. Machine learning on dna-encoded library count data using an uncertainty-aware probabilistic loss function, Lim et al., JCIM 2022
6. Hergcentral: a large database to store, retrieve, and analyze compound-human ether-a-go-go related gene channel interactions to facilitate cardiotoxicity assessment in drug development, Du et al., Assay and drug development technologies, 2011

## Acknowledgements

This study was partially funded by the Horizon Europe funding programme under the Marie Skłodowska-Curie Actions Doctoral Networks grant agreement "Explainable AI for Molecules - AiChemist", no. 101120466. The work of Johan Östman was funded by Vinnova, the Swedish innovation agency, under grant 2023-03000.