

# An Analysis Of Protected Health Information Leakage In Deep-Learning Based De-Identification Algorithms

Salman Seyedi<sup>1,\*</sup>, Li Xiong<sup>2</sup>, Shamim Nemati<sup>3</sup>, Gari Clifford<sup>1,4</sup>

<sup>1</sup>Department of Biomedical Informatics, Emory University, Atlanta, GA, <sup>2</sup>Department of Computer Science, Emory University, Atlanta, GA,

<sup>3</sup>Department of Biomedical Informatics, University of California San Diego Health, La Jolla, CA, <sup>4</sup>Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA

\*sseyedi@dbmi.emory.edu

<https://arxiv.org/abs/2101.12099>

## Goal

Investigating potential leakage of sensitive information from a de-identification algorithm

## Investigated Model (NeuroNER)

NeuroNER  
(State of the art De-identification Model)

Layer 1: Text tokenizer

Layer 2: Neural Networks

Layer 3: Conditional Random Field

## Re-identification Attacks

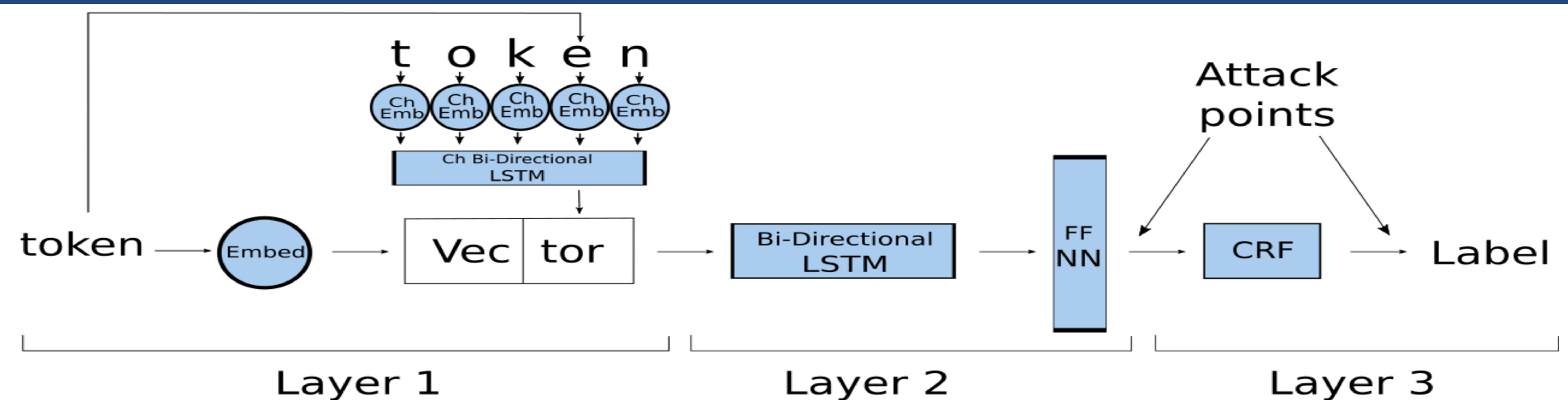
White Box Attempts:

- Naive cut-off
- Brute-force cut-off
- Membership inference attack

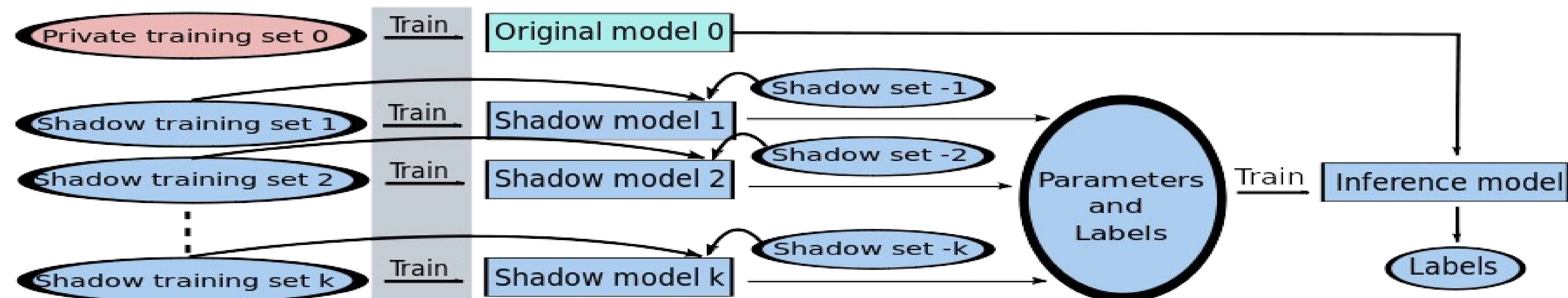
## Results and Conclusions

- Despite different distributions, zero successful re-identification
- Model not prone to several implemented attacks
- Statistically different distributions but with overwhelming overlap for successful cut-off attacks

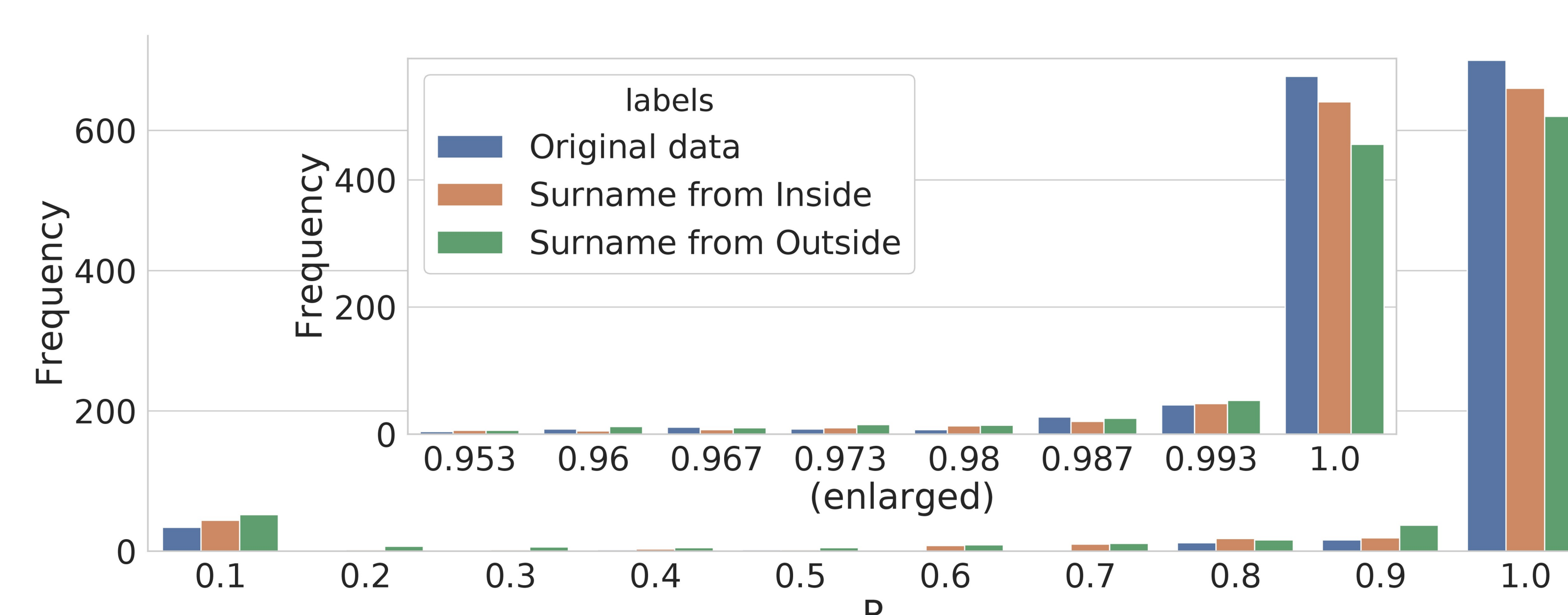
## De-Identification Deep Neural Network Model (NeuroNER) And Investigated Attack Points



## Membership Inference Attack



## Histogram of Probabilities for Surnames



## Acknowledgments

National Science Foundation, grant # 1822378  
'Leveraging Heterogeneous Data Across International Borders in a Privacy Preserving Manner for Clinical Deep Learning',  
The National Center for Advancing Translational Sciences of the National Institutes of Health under Award Number UL1TR002378.

