





Partially Encrypted Machine Learning using Functional Encryption

THÉO RYFFEL^{1,2}, EDOUARD DUFOUR-SANS¹, ROMAIN GAY^{1,3}, FRANCIS BACH^{2,1}, DAVID POINTCHEVAL^{1,2} ¹ENS, CNRS, PSL University, Paris, France ²INRIA, Paris, France ³University of California, Berkeley



Functional Encryption in Machine Learning

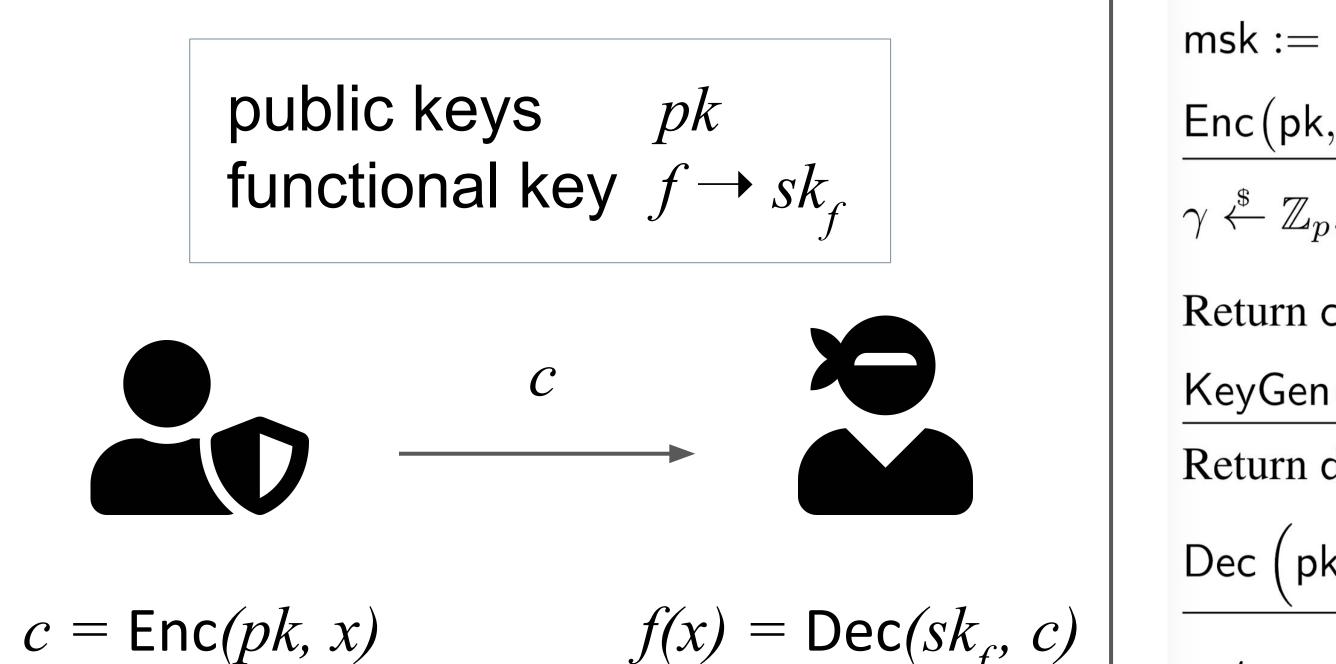
Build non-interactive encrypted machine learning as a service (MLaaS) using Functional Encryption

Use cases:

- Spam filtering: encrypted emails be classified as spam or not by the email server
- Privacy-preserving enforcement of content policies: maintaining user privacy while filtering out illegal content in messaging apps

What is Functional Encryption?

Functional Encryption allows a non-trusted party to learn the output of a specific function over encrypted data without interacting with the data owner.



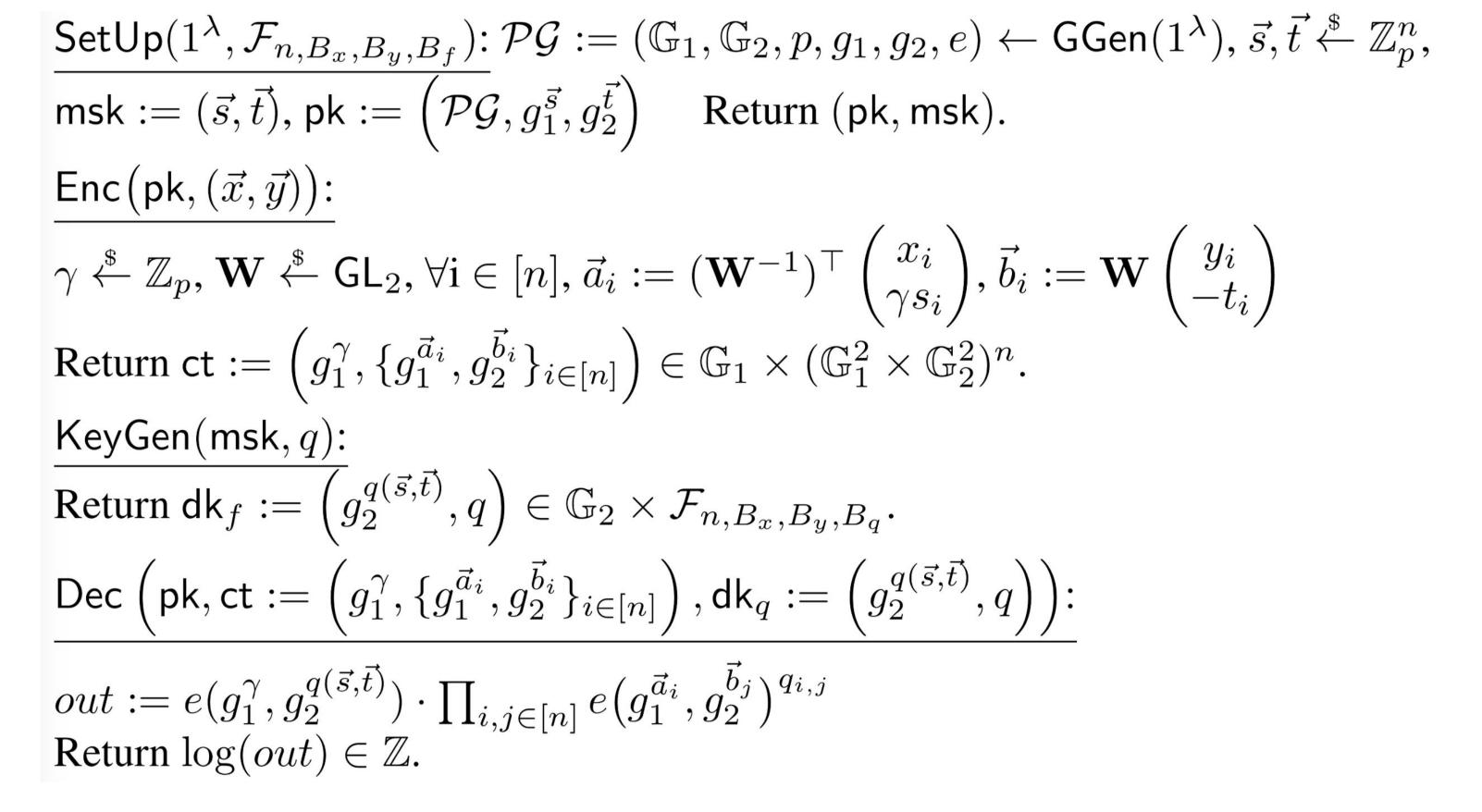
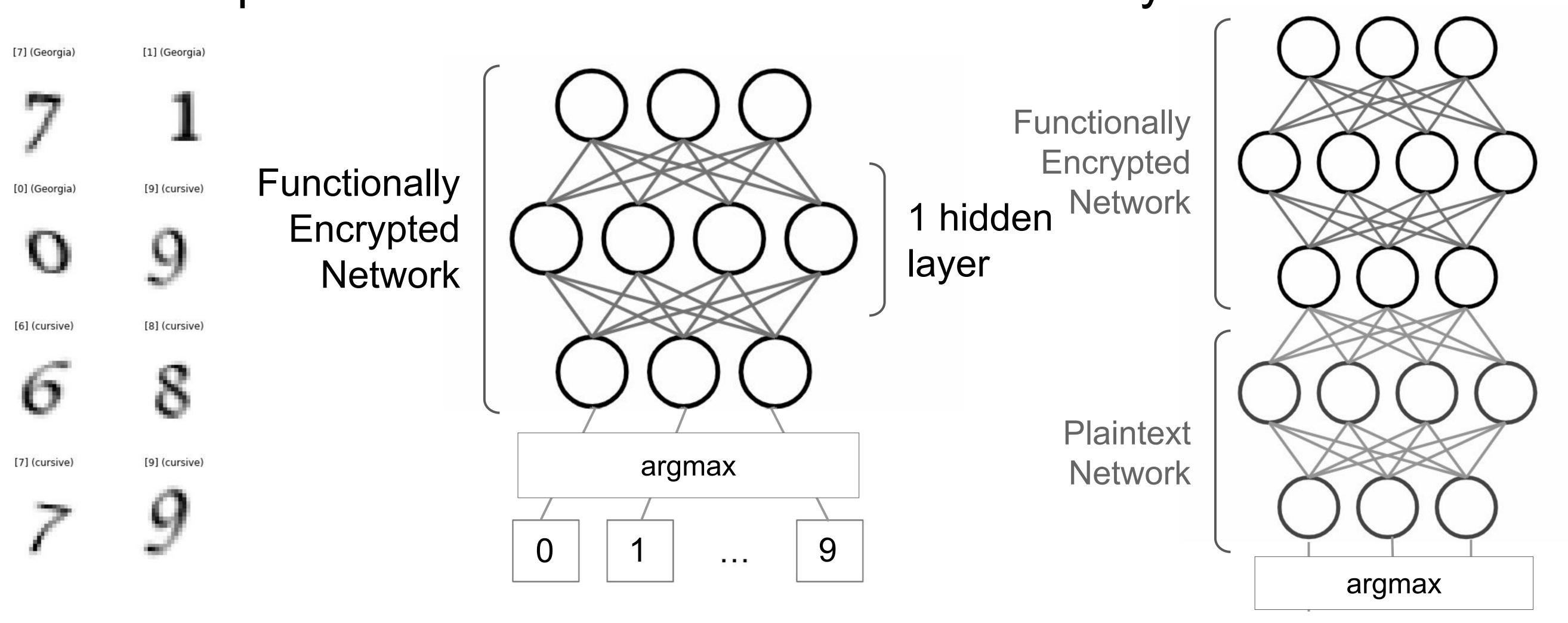


Figure: Our quadratic functional encryption scheme

Application to Neural Network

We propose an efficient quadratic functional encryption scheme which is equivalent to a network with one hidden layer.



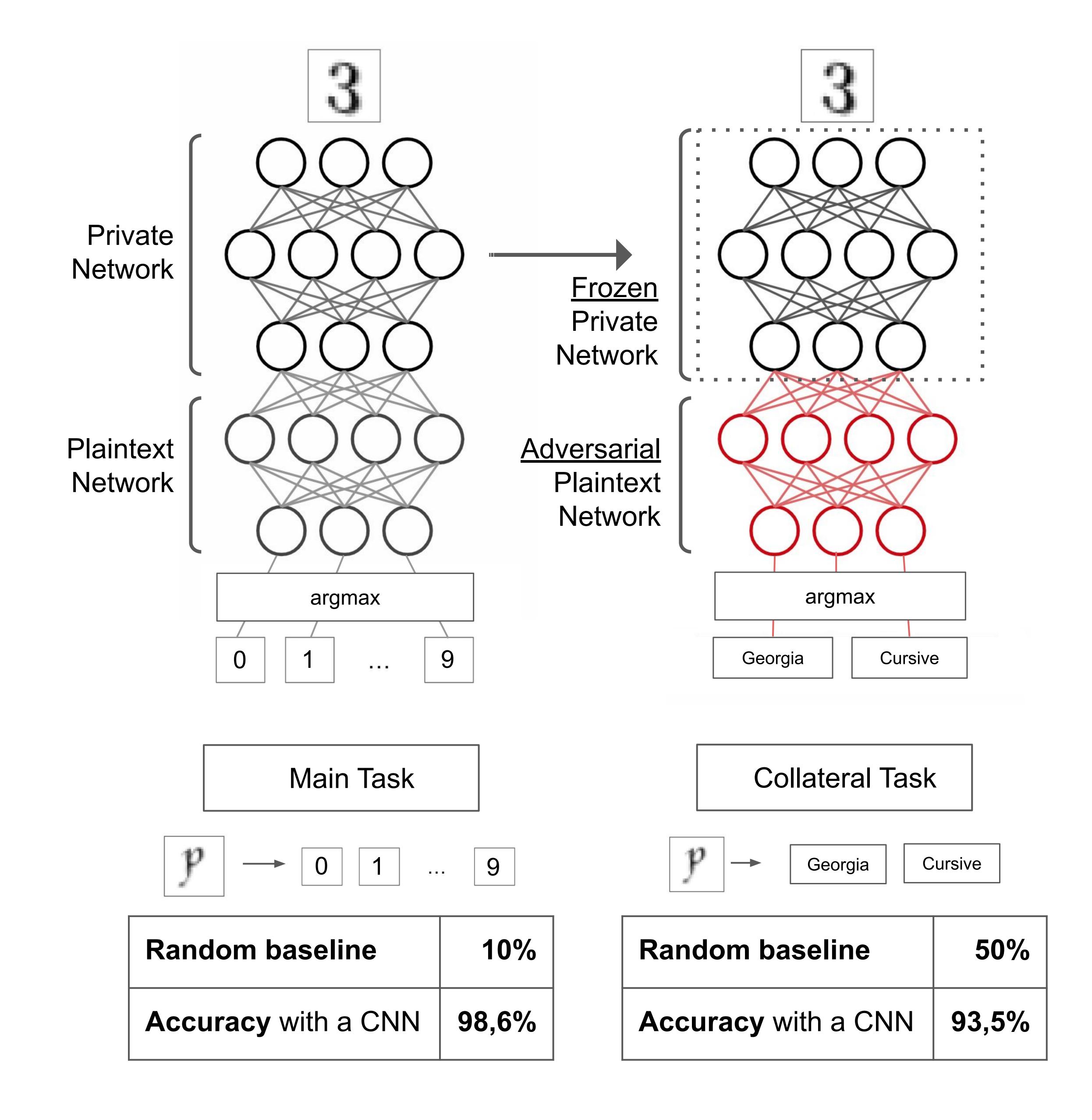
Collateral Learning, a threat to

The neural network trained to perform a specific task can also be used to perform another different and highly sensitive task

An adversary can misuse the private network to perform an unexpected task and extract private information about the data. This is done below by changing the plaintext network and the target labels.

Example:

Instead of predicting the digit on the image, the adversary tries to detect the font used.



Defeat Collateral learning with adversarial learning

We make the private model robust by simulating adversaries. We train it simultaneously on the main objective and against the collateral objective of the simulated adversary.

Semi-adversarial training protocol:

- 1. Optimize the main plaintext network with the private one freezed
- 2. Optimize the adversarial plaintext one with the private net freezed
- 3. Optimize the private network with the two others freezed, on the joint objective:

$$L_{joint} = L_{main}$$
 - $\alpha \cdot L_{collateral}$

where L_{main} and $L_{collateral}$ respectively are the loss of the main and adversarial networks

Robustness is further increased by simulating stronger adversaries and by reducing the output size of private network:

Main Task		
Random baseline	10%	Ra
Accuracy with a CNN	98,4%	Ac

Random baseline	50%
Accuracy	
- with a CNN	55,3%
- with best sklearn model	58,9%

Collateral Task

Get in touch



github.com/LaRiffle/collateral-learning



9 detailed tutorials





References

Baltico et al. Practical FE for quadratic functions, 2017 Carpov et al. Illuminating the dark, 2018

Louppe et al. Learning to pivot with adversarial networks, 2017