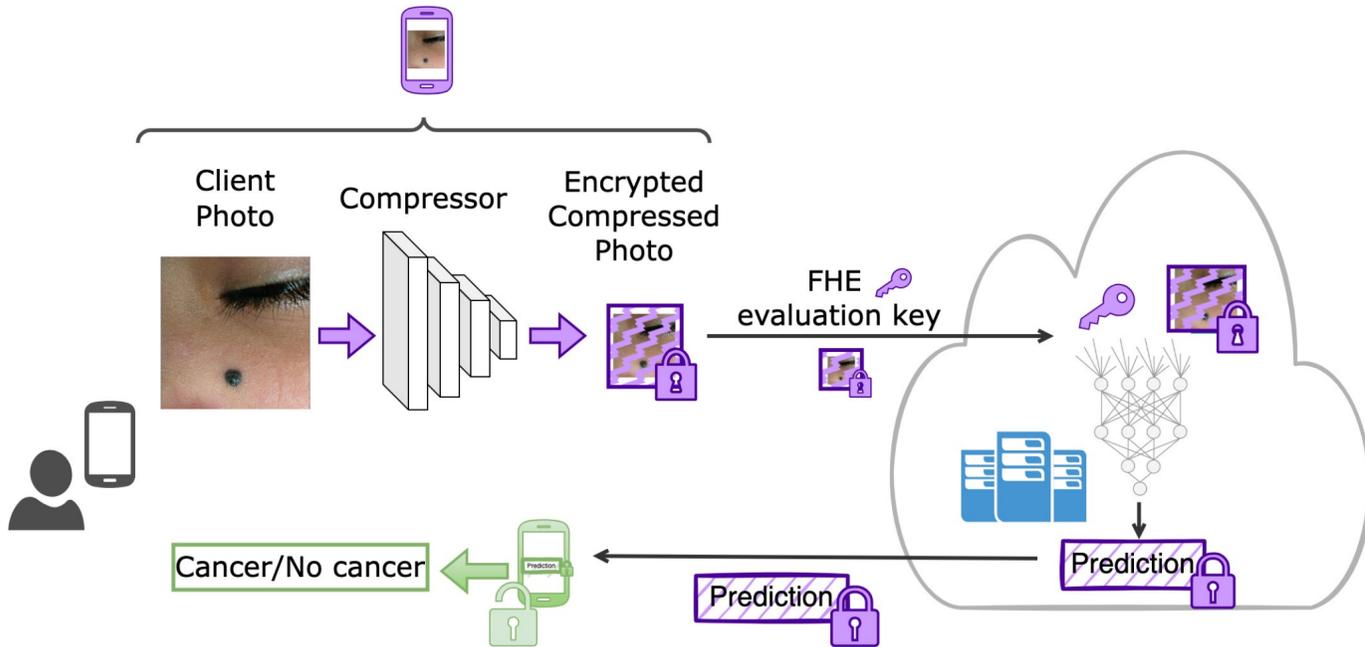


# VEIL: A FHE Framework for PPML

Abdelrahman Aly   Chiara Marcolla   Victor Mateu   Pradeep Mishra   Aryan Raj   Sergi Rovira   Victor Sucasas  
Luis Zamarripa   Floyd Zweyding



## Tools

Scheme used: CKKS  
Polynomial degree:  $N = 2^{15}$   
Ciphertext modulus:  $\log_2 q = 590$   
Security:  $\lambda = 128$

Approximate AF: Chebyshev polynomial of degree 15

### Library used

- CPU: OpenFHE v1.2.3
- GPU: FidesLib v1.0

### CNN models:

- EfficientNet (5.5M – 60M)
- BERT (14M – 100M)

## How to Add Privacy to a CNN

VEIL is a privacy-preserving machine learning framework that enables encrypted inference on neural networks by combining pretrained representation learning with FHE. The framework operates in two stages:

- 1 A **teacher model** is trained on a target dataset: each input sample is transformed into a low-dimensional vector. The plaintext never leaves the user's device.
- 2 VEIL trains a FHE-friendly **student model** which *learns* to replicate the behavior of the original encoder-classifier pipeline.

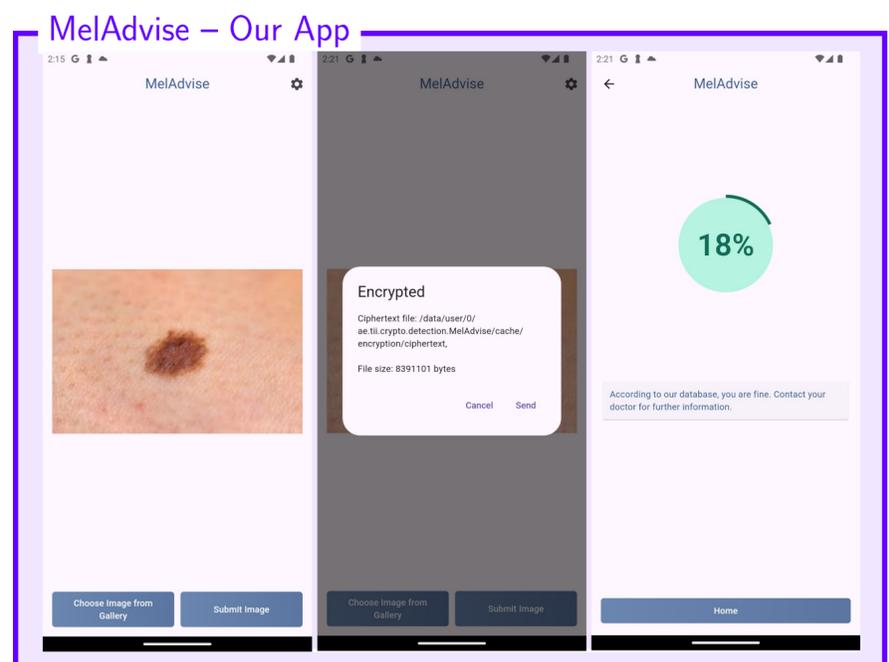
Only the student model is evaluated under encryption!

## Client Benchmarks

- Communication costs are *not* included.
- Client keys: 282 MB memory.
- Image size: 1MB (128 bits after compression)

Device	iPhone 6s	Android Pixel XL	Ubuntu 24.04
Key gen	13	6	2
Enc	0.4	0.3	0.1

Client pre-processing (key generation and encryption) time (in seconds)



## Encrypted Inference Benchmarks

- Latency for a **single query**, with **one encrypted image**.
- Times include the **server startup phase** – worst-case scenario  $\implies$  GPU initialization phase  $\approx 0.6$  s.

Dataset / Machine	A100	A100 + GPU	Aldebaran
ISIC Melanoma	$7.32 \pm 0.05$	$2.08 \pm 0.06$	$11.25 \pm 0.44$

Inference time (s) on encrypted ISIC skin cancer dataset, verified against plaintext (max. error  $5 \times 10^{-4}$ , recall **0.94**).

**A100** 512 GB RAM, AMD EPYC 7742 @ 2.25 GHz, 4x NVIDIA A100 (80 GB).

**A100 + GPU** Same as A100, with FidesLib GPU backend (CUDA 12.8, Ubuntu 24.04).

**Aldebaran** 512 GB RAM, dual Intel Xeon Gold 6250L @ 3.9 GHz (2x8 cores, SMT), Docker on Ubuntu 20.04.

### Contact

- C. Marcolla:  
chiara.marcolla@tii.ae
- S. Rovira:  
sergi.rovira@tii.ae