# HE-MAN
# Homomorphically Encrypted MAchine learning with oNnx models

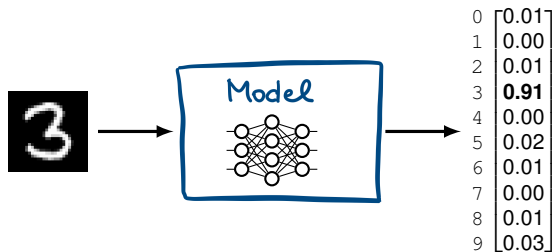**Martin Nocker**, David Drexel, Michael Rader,
Alessio Montuoro, Pascal Schöttle

# Machine Learning Services

# Machine Learning Services

# Classical Cryptosystems

# Classical Cryptosystems

**Client side**

**Server side**

4

# Fully Homomorphic Encryption (FHE)



- $+$ <u>and</u> $\times \Rightarrow$ **Fully** Homomorphic Encryption (FHE)

# Fully Homomorphic Encryption (FHE)



Further challenges:

- FHE operations are orders of magnitude more complex
- Only additions and multiplications of ciphertexts are possible

# Fully Homomorphic Encryption (FHE)



**Client side**

**Server side**

NO
cleartext
data

```
0  ⎡0.01⎤
1  ⎢0.00⎥
2  ⎢0.01⎥
3  ⎢0.91⎥
⋮  ⎢ ⋮  ⎥
9  ⎣0.03⎦
```

# HE-MAN
## High-level architecture



Previous work

- FHE-implementation of specific NNs [BGBE19]
- Individual ML framework support: TensorFlow [RRK+20], PyTorch [KVH+21]

# HE-MAN
## High-level architecture



Why ONNX?

- Format definition
- Framework independence
- Broad language support

# ONNX in HE-MAN

So far implemented
- AddOperator
- AveragePoolOperator
- ConstantOperator
- ConvOperator
- FlattenOperator
- GemmOperator
- MatMulOperator
- MulOperator
- PadOperator
- ReluOperator
- ReshapeOperator
- SubOperator

# Evaluation

- Performance
  - Classification accuracy
  - Latency
- Handwritten digit classification
  - MNIST

  

- Face recognition
  - LFW (Labeled Faces in the Wild)

# Results

| Dataset | baseline | HE-MAN-Concrete | | HE-MAN-TenSEAL | |
|---|---|---|---|---|---|
| Network | accuracy | accuracy | latency | accuracy | latency |
| **MNIST** | | | | | |
| CryptoNets | **.975** | **.968** | 279 s | .924 | 8 s |
| LeNet-5 | **.991** | **.984** | 1672 s | .789 | 237 s |
| **LFW** | | | | | |
| MobileFaceNets (classifier) | **.990** | .970 | 69 s | **.972** | 208 s |

**Key Result:** accuracy on par with plaintext result, at increased runtime

# Conclusion

HE-MAN

- Neural network inference on homomorphically encrypted data
- Preserves privacy of model and data
- Accuracies close to cleartext results
- Broad model support via ONNX format

Future work:

- Full ONNX operator set implementation

# Thank you!



https://github.com/smile-ffg/he-man-concrete

https://github.com/smile-ffg/he-man-tenseal

Paper:



https://arxiv.org/abs/2302.08260

# References I

[BGBE19] Alon Brutzkus, Ran Gilad-Bachrach, and Oren Elisha.
Low latency privacy preserving inference.
In *International Conference on Machine Learning*, pages 812–821. PMLR, 2019.

[KVH+21] Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, and Laurens van der Maaten.
Crypten: Secure multi-party computation meets machine learning.
In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 4961–4973. Curran Associates, Inc., 2021.

[RRK+20] Deevashwer Rathee, Mayank Rathee, Nishant Kumar, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul Sharma.
Cryptflow2: Practical 2-party secure inference.
In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, CCS '20, page 325–342, New York, NY, USA, 2020. Association for Computing Machinery.
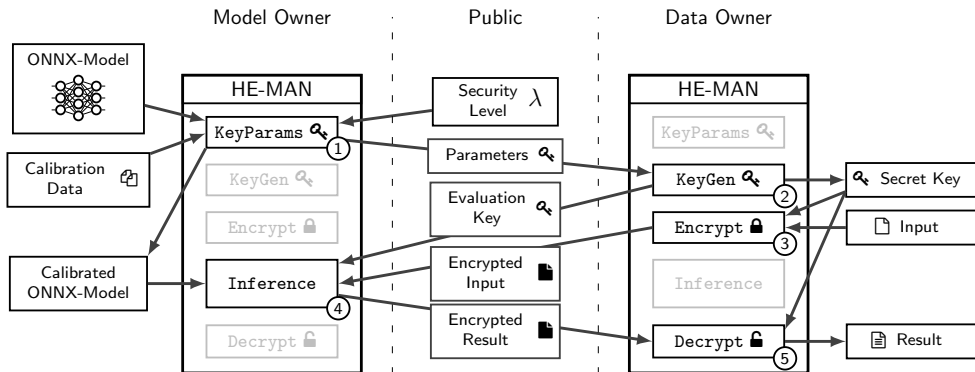
# RSA

$$\text{Encrypt: } c = m^e \mod N$$

$$\prod_i c_i = \prod_i m_i^e = \left( \prod_i m_i \right)^e \mod N$$

$$\sum_i c_i = \sum_i m_i^e \neq \left( \sum_i m_i \right)^e \mod N$$

# HE-MAN Architecture

# Crypto Parameters in FHE
**TenSEAL**

| $N$ | $\log_2 N$ | $\log_2 q$ | | |
|---|---|---|---|---|
| | | $\lambda = 128$ | $\lambda = 192$ | $\lambda = 256$ |
| 2048 | 11 | 54 | 37 | 29 |
| 4096 | 12 | 109 | 75 | 58 |
| 8192 | 13 | 218 | 152 | 118 |
| 16384 | 14 | 438 | 300 | 237 |
| 32768 | 15 | 881 | 600 | 476 |