





Keep it Unsupervised: Horizontal Attacks Meet Simple Classifiers

Ninon Calleja Albillos – CEA LETI – France

Sana Boussam – THALES – France

With the guidance of Gabriel Zaid, Ange Martinelli, Eleonora Cagli, Simon Abelard, Antoine Loiseau, Julien Eynard and Guénaël Renault



THALES



Keep it Unsupervised: Horizontal Attacks Meet THIS PAPER Simple Classifiers

Sana Boussam^{2,3} and Ninon Calleja Albillos^{1,4}

¹ CEA, LETI, MINATEC Campus, F-38054 Grenoble, France ² INRIA and LIX, Institut Polytechnique de Paris, Palaiseau, France ³ Thales ITSEF, Toulouse, France ⁴ Univ. Grenoble Alpes, F-38000, Grenoble, France sana.boussam@inria.fr bits are located in each separate randomized private key. Horizontal Attacks . **Keywords:** Side-channel Analysis - Public-key Algorithms Deep Learning ninon.callejaalbillos@cea.fr key bits, the proposed iterative framework improves the *V*, ked where 90% on average and to 100% for at least one of the need to know where remains fully unsupervised and excludes the need to know where 90% on average and to 100% for at least one of the attacted where the need to know where the rend to know where the need to know where the render the need to know where the render the need to know where the render the need to know where the need to know the need When a simple horizontal attack can recover around the key bits, the proposed iterative framework improves on e of the attack $QQ_{0\%}^{\circ}$ on average and to 100% for at least one of the strategies. (or exponents). borkontal attack can recover around When a simple borkontal attack ramework immediate key bits, the proposed iterative framework immediate scalar multiplication (or exponent the in-network can significantly reduce the in-(or ernoments) resulting from a custering of exponent a deep learning vaso a deep learning bases

become imprat. many with imp

52% accuracy key recovery

10

del

oisy

100% accuracy key recovery

Reproducibility?

Benefits?

Limitotions? ...

remains fully unsupervised and excludes the need to know bits are located in each separate randomized private key. Deep Learning



1. Introduction

2. Overview of Perin et al. CHES 2021 [PCBP20]

3. In-depth analysis of the iterative framework

Simplification of the neural networks Correction capability of the iterative framework

4. Deep Learning Free Approach

5. Conclusion



Side-Channel Attacks





Keep it Unsupervised: Horizontal Attacks Meet Deep Learning

Target



Algorithm 1 Montgomery ladder with cswap and coordinate re-randomization

```
\begin{array}{l} \text{...initialization omitted...}\\ bprev \leftarrow 0\\ \textbf{for } i = 254...0 \ \textbf{do}\\ \mid & RE\_RANDOMIZE\_COORDS(work\_state)\\ & b \leftarrow \text{bit } i \ \text{of the secret scalar}\\ & s \leftarrow b \oplus bprev\\ & bprev \leftarrow b\\ & CSWAP(work\_state,s)\\ & LADDERSTEP(work\_state)\\ \textbf{end}\\ & \dots \text{return omitted...} \end{array} \triangleright \ \textbf{Leakage (see [NCOS16, NC17])}\\ \end{array}
```

Two datasets o CSWAP_ARITH o CSWAP_POINTER

Targets: protected ECC software implementations on Curve25519 from μNaCl https://munacl.cryptojedi.org/curve25519-cortexm0.shtml



Target – EM acquisition



Time

Selection of subtraces corresponding to CSWAP executions during an ECSM.

Target – Dataset construction



- A trace corresponds to the 255 CSWAPs concatenation. Each subtrace *ti* corresponds to the processing of the *i*-th bit of the secret scalar **s**.
- Attack bit by bit proposed in [PCBP20] to recover the 255-bit secret scalar s.

Overall attack methodology presented in [PCBP20]

CSWAP

<u>cea</u>

Applying Horizontal Clustering Side-Channel Attacks on Embedded ECC Implementations Nascimento, Chmielewski [NC17] Keep It Unsupervised : Horizontal Attacks Meet Deep Learning *Perin, Chmielewski, Batina, Picek* [PCBP20]



Overall attack methodology

Corrective phase - Presentation of the Iterative Framework: Preliminaries

Traces processing bit 1



Creation of the dataset.



Phase 1 - Split of \mathcal{D} into \mathcal{D}_1 (blue) and \mathcal{D}_2 (yellow).

Keep it Unsupervised: Horizontal Attacks Meet Simple Classifiers





Phase 2 – Training on \mathcal{D}_{l} (blue) in order to relabel \mathcal{D}_{2} (yellow).

Keep it Unsupervised: Horizontal Attacks Meet Simple Classifiers





Phase 3 – Training on \mathcal{D}_2 (yellow) in order to relabel \mathcal{D}_1 (blue).



Phase 4 – Merging of \mathcal{D}_1 (blue) and \mathcal{D}_2 (yellow) into \mathcal{D} .

Keep it Unsupervised: Horizontal Attacks Meet Simple Classifiers

Results achieved in [PCBP20]

	cswap_pointer		cswap_arith		
	Average	Maximum	Average	Maximum	
After prelabelling phase	52.24%	59.22%	52.44%	59.22%	
CNN with no regularisation	85%	97.64%	52%	76.07%	
CNN with regularisation (dropout + data augmentation)	91%	100%	83%	100%	

Average and maximum single trace accuracies obtained in [PCBP20] for fixed CNN hyperparameters (50 framework iterations).

Our hypothesis: the model complexity is too high and leads to overfitting

3 In-depth analysis of the iterative framework

Data analysis: linear decision boundaries



Welch's t-tests obtained on the two datasets using true labels.

K-means on interesting regions : High average accuracy



Simplification of the neural networks

Layers	_	<u> </u>							
Conv1D 1	8 filters			I	_ayers		-		
				S	igmoid	1 ne	euron		
Conv1D_2	16 filters				-				
Conv1D_3	32 filters	Architecture proposed fo <mark>r both datasets -</mark> Perceptron, a linear classifier .							
Dense_1	100 neurons	Nb. of parameters Attack duration							
Dense_2	100 neurons			(a)	(b)	(c)	(a)	(b)	(c)
<u> </u>	0		CNN [PCBP20]	45,978	3,069,820	669,820	75 min.	1,080 min.	300 min.
Softmax	2 neurons		Perceptron	1,001	8,001	2,001	8 min.	35 min.	10 min.
too complex.									

CNN architecture for both datasets [PCBP20].

cea

Keep it Unsupervised: Horizontal Attacks Meet Simple Classifiers

Simplification of the neural networks: results



Average accuracy obtained using iterative framework with our neural network after 50 iterations (no regularization)

	(a)	(b)	(c)
Perceptron, no regularisation	98.9 %	70.93%	98 %
CNN [PCBP20], no regularisation	85%	52%	58%
CNN [PCBP20], dropout $+$ data augmentation	91%	83 %	98 %

Correction capability of the iterative framework – Case study: noisy traces



Visualization of cswap_pointer dataset's clusters

Correction capability of the iterative framework – Case study: noisy traces – Results

	$\sigma = 15$	$\sigma = 30$	$\sigma = 90$
CNN [PCBP20]	99.87%	98.83%	79.43%
CNN [PCBP20]	97.11%	71.25%	50.78%
Perceptron	98.87%	96.66%	77.1%
Perceptron	98.21%	67.15%	56.43%

Average single trace accuracy achieved for cswap pointer dataset in a **supervised/unsupervised** setting.

Deep Learning Free Approach



Deep Learning free attack

- \Box Features extraction \rightarrow PCA
- Component selection issue
 - Visual selection
 - Natural order (decreasing eigenvalues)
 - Explained Local Variance (ELV [CDP15])
 - Denoising strategy
- \Box Clustering \rightarrow K-means



cswap_arith



(a) cswap_pointer

Deep Learning Free Results

	Cswap_pointer			Cswap_arith			
	Average	Maximum	PC/Pols	Average	Maximum	PC/Pols	
PCA+Virtual	94.1%	97.7%	8th	60.8%	72.2%	6th	
PCA+Natural order	94.1%	97.7%	8th	76.1%	82.4%	18th	
PCA + ELV + Denoising	98% 🗖	100% 📕	8th	90.6% 📕	96.1%	18th	
SOST [PCBP20]	52.24%	55.22%	20 Pols	52.44%	55.22%	20 Pols	
SOST+Framework [PCBP20]	91% 📕	100% 📕	20 Pols	83% 📕	100% 🗖	20 Pols	

91% \rightarrow 23 wrong bits on average \rightarrow 2¹⁰⁹ operations to end the attack

96,1% \rightarrow 10 wrong bits on average \rightarrow 2⁵⁸ operations to end the attack

98% \rightarrow 5 wrong bits on average \rightarrow 2³³ operations to end the attack

Deep Learning Free approach leads to **at least the same successfull attacks** than [PCBP20] in the studied case.



5 Conclusion



Conclusion

- The study proposed in [PCBP20] was biased by the simplicity of the used datasets (two linear separable classes)
- Such a simplicity was hidden by the huge complexity of the proposed deeplearning-based solution
- When moving to a more complex scenario, the iterative framework is strongly impacted.

Open problem: finding a good strategy to correct mislabeled data, filling the supervised/unsupervised gap

- Is it suitable to use very complex machine-learning solution for general problems?
- Instead: focus on the simplest possible attacks → better insight on the actual threats
- In case machine-learning is needed \rightarrow size it carefully!





Thank you