

Don't Learn What You already Know

Scheme-Aware Modeling for Profiling Side-Channel Analysis against Masking

Loïc Masure, Valence Cristiani, Maxime Lecomte, François-Xavier Standaert

Prague, September 12th







Content

Introduction: SCA & Masking

Deep Learning Against Masking

Scheme-Aware Approach

The Elephant in the Room

Conclusion

Context : Side-Channel Analysis (SCA)

"Cryptographic algorithms don't run on paper, they run on physical devices" Msg





Black-box cryptanalysis: \rightarrow Exponential with N

Side-Channel Analysis: \rightarrow Exponential with bit-size \rightarrow Linear with N

Trace : power, EM, acoustics, runtime, ...

The Counter-Measure: Masking

 $Y(\mathsf{secret})$

The Counter-Measure: Masking



The Counter-Measure: Masking



The Counter-Measure: Masking



The Counter-Measure: Masking



4 / 23

Content

Introduction: SCA & Masking

Deep Learning Against Masking

Scheme-Aware Approach

The Elephant in the Room

Conclusion

Deep Learning (DL) for SCA



General way to modelize, *i.e.*, to convert leakage into probabilities

$$\begin{array}{cccc} F : & \mathcal{L} & \longrightarrow & \mathcal{P}(\mathcal{Y}) \\ \mathbf{I} & \longmapsto & \mathbf{y} = F(\mathbf{I}) \approx \Pr\left(Y \mid \mathbf{L} = \mathbf{I}\right) \end{array}$$
(1)

F(I): output of a Directed Acyclic Graph (DAG) of computation:

Each node: elementary function $f_i(\cdot, \theta_i)$

 θ_i : *parameters* fully describing f_i

Shape of the DAG, nature of the classes of functions: architecture of the DNN.

Training a DNN for Profiled SCA



(Open sample)

Training a DNN for Profiled SCA



(Open sample)

Training a DNN for Profiled SCA



Training a DNN for Profiled SCA



Training a DNN for Profiled SCA



Training a DNN for Profiled SCA



 \mathcal{L} (): loss function to minimize, with gradient descent

Training a DNN for Profiled SCA



 \mathcal{L} (): loss function to minimize, with *gradient descent* **Uninformed adversary**: no knowledge of random shares during profiling

Training a DNN for Profiled SCA



\mathcal{L} (): loss function to minimize, with *gradient descent* Worst-case adversary: knowledge of random shares during profiling

How to profile as Worst-Case Adversary?



The natural way: divide & conquer $\rightarrow \Pr(Y \mid \mathbf{L})$ decomposed as collection of $\Pr(Y_i \mid \mathbf{L}_i)$

How to profile as Worst-Case Adversary?



The natural way: divide & conquer $\rightarrow \Pr(Y \mid L)$ decomposed as collection of $\Pr(Y_i \mid L_i)$ \rightarrow Each, modeled by m_{θ_i}

ightarrow Separately trained with \mathcal{L}_{y_i}

How to profile as Worst-Case Adversary?



The natural way: divide & conquer $\rightarrow \Pr(Y \mid L)$ decomposed as

collection of $\Pr(Y_i \mid \mathbf{L}_i)$

- \rightarrow Each, modeled by m_{θ_i}
- ightarrow Separately trained with \mathcal{L}_{y_i}
- \rightarrow Then use \circledast to recombine

How to profile as Uninformed Adversary?



The End-to-End Way:

 $ightarrow \mathsf{Pr}\left(\mathrm{Y} \mid \mathbf{L}
ight)$ directly modeled by $\mathsf{m}_{ heta}, \, \mathsf{trained} \, \mathsf{with} \, \mathcal{L}_y$

Simulated Experiments



Figure: Learning curves: MI estimation vs. data complexity.

What Kind of Adversary for Evaluation ?

Adversary	Worst-case	Uninformed
Access to shares	Yes	No
Knowledge of scheme	Yes	No

Access to shares during profiling: Easy to reach optimal attacks \checkmark

Too conservative X Not realistic X

What Kind of Adversary for Evaluation ?

Adversary	Worst-case	Scheme-Aware	Uninformed
Access to shares	Yes	No	No
Knowledge of scheme	Yes	Yes	No

Access to shares during profiling: Easy to reach optimal attacks \checkmark

Too conservative X Not realistic X

How to leverage the knowledge of the masking scheme, without relying on the knowledge of the shares?

Content

Introduction: SCA & Masking

Deep Learning Against Masking

Scheme-Aware Approach

The Elephant in the Room

Conclusion

Don't Learn what You Already Know !

Can we find a trade-off between both approaches ?



 \rightarrow Model still decomposed as collection of $\Pr(\mathbf{Y}_i \mid \mathbf{L}_i)$

Don't Learn what You Already Know !

Can we find a trade-off between both approaches ?



- $\rightarrow \text{ Model still decomposed as} \\ \text{ collection of } \Pr\left(\mathbf{Y}_i \mid \mathbf{L}_i\right)$
- \rightarrow Still recombined with \circledast but ...

Don't Learn what You Already Know !

Can we find a trade-off between both approaches ?



- \rightarrow Model still decomposed as collection of Pr (Y_i | L_i)
- \rightarrow Still recombined with \circledast but ...
- ightarrow ... Training done *jointly* with \mathcal{L}_y
- \rightarrow Need to backprop gradients through \circledast

Pytorch code available at github.com/uclcrypto/Scheme-Aware-Architectures

Back to our Simulated Experiments



Figure: Learning curves: MI estimation vs. data complexity.

Scheme-Aware spares some data complexity Loic Masure Don't Learn What You already Know

Application to Experimental Data



Figure: Learning curves: MI estimation vs. data complexity.

Don't Learn What You already Know

What about Higher Order Masking?

Affine masking: positive results on simulation, but not on experimental data. Why ? The Plateau Effect



Don't Learn What You already Know

Content

Introduction: SCA & Masking

Deep Learning Against Masking

Scheme-Aware Approach

The Elephant in the Room

Conclusion

The Elephant in the Room



(a) Timon, Ches'19

0.6

0.4 8 0.2

0.0

200 400



(b) Perin & Picek, SAC'20

Architecture variant 1

Architecture variant 3



(c) Cristiani et al., JoC'23

-f = Id

f = HW

Loïc Masure

(a) First order masking $(\sigma = 1)$

Epoch

200 400 600 800 1000

0.4

-0.2

Bits 0.0

Don't Learn What You already Know

How Masking Affects the Plateau Length

Simulation with HW leakage model and *exhaustive* dataset (no profiling error)



An Explanation

Theorem $(INFORMAL^1)$

Assume that each \mathbf{L}_i is i.i.d. standard Gaussian in \mathbb{R}^p . Define the target function $h_{\boldsymbol{u}}(\boldsymbol{I}) = \prod_{i=1}^d \operatorname{sign}(\boldsymbol{u}^{\mathsf{T}}\boldsymbol{I}_i)$, for some normalized hyperplane \boldsymbol{u} . Let \mathfrak{m}_{θ} be a model, such that $\mathbb{E}\left[\|\nabla_{\theta} \mathfrak{m}_{\theta}\|^2 \right] \leq G(\theta)^2$. Then,

$$\mathbb{E}_{oldsymbol{u}}\left[\left\|
abla_{ heta}\mathcal{L}\left(heta
ight)-\mathbb{E}_{oldsymbol{u}}\left[
abla_{ heta}\mathcal{L}\left(heta
ight)
ight]^{2}
ight]\leq G(heta)^{2}\cdot\mathcal{O}\left(\sqrt{rac{d\log(p)}{p}}
ight)^{d}$$

The gradient almost takes the same direction, regardless of u !

¹Shalev-Shwartz, Shamir, and Shammah, "Failures of Gradient-Based Deep Learning", p. ICML 2017. Loïc Masure Don't Learn What You already Know

Content

Introduction: SCA & Masking

Deep Learning Against Masking

Scheme-Aware Approach

The Elephant in the Room

Conclusion

Open Problems

How to tackle higher orders with DL remains unclear: GD really not suitable ?

- Efficient surrogate to gradient descent ?

 \rightarrow Then current evaluator run suboptimal attacks

- No efficient surrogate to GD (reduction to hard learning problem) ?

 $\rightarrow Then$ intrinsic gap between worst-case approach and others

Worth investigating, no matter the answer !

References I

Shalev-Shwartz, S., O. Shamir, and S. Shammah. "Failures of Gradient-Based Deep Learning". In: Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017. Ed. by D. Precup and Y. W. Teh. Vol. 70. Proceedings of Machine Learning Research. PMLR, 2017, pp. 3067–3075. URL: http://proceedings.mlr.press/v70/shalev-shwartz17a.html.