

# A Comprehensive Study of Deep Learning for Side-Channel Analysis

Loïc Masure<sup>1,3</sup>   Cécile Dumas<sup>1</sup>   Emmanuel Prouff<sup>2, 3</sup>

<sup>1</sup>Univ. Grenoble Alpes, CEA, LETI, DSYS, CESTI, F-38000 Grenoble  
loic.masure@cea.fr

<sup>2</sup>ANSSI, France

<sup>3</sup>Sorbonne Université, UPMC Univ Paris 06, POLSYS, UMR 7606, LIP6, F-75005, Paris, France

17/09/2020, CHES

## Outline

### 1. Context

2. SCA Optimization Problem versus Deep Learning Based SCA

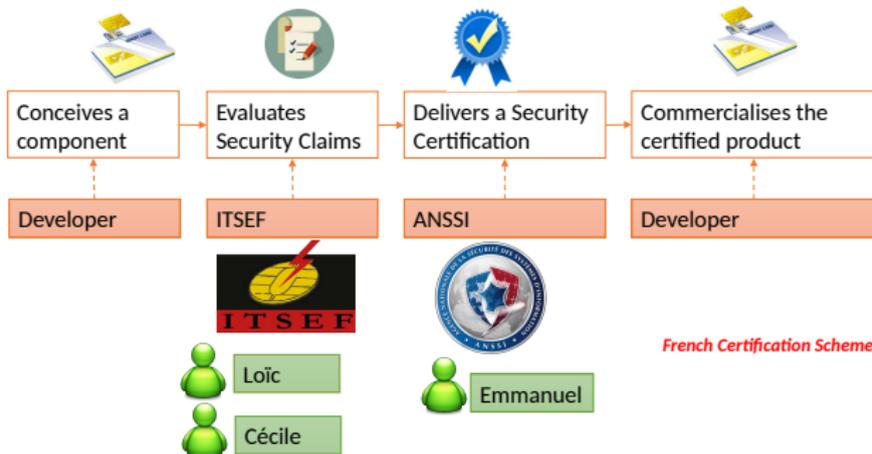
3. NLL Minimization is PI Maximization

4. Simulation results

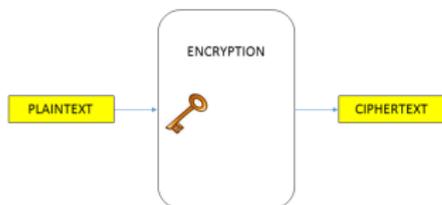
5. Experimental results

## Who am I

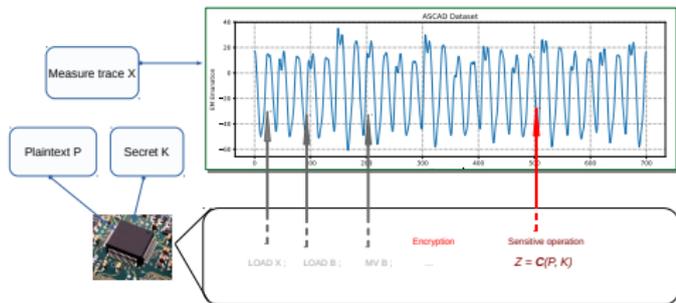
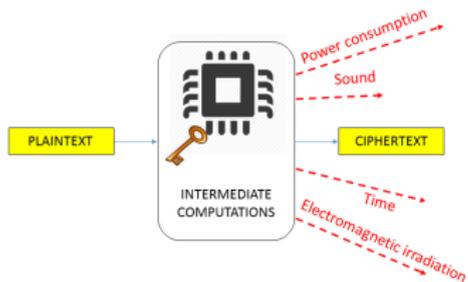
- PhD student, studying Deep Learning (DL) for Side-Channel Analysis (SCA)



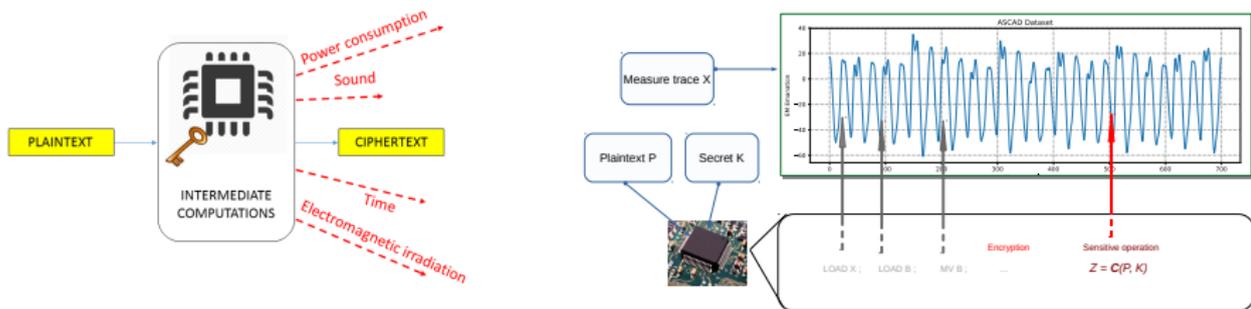
## What is SCA?



## What is SCA?



## What is SCA?



## Profiling Attack

Attack using *open samples* similar to the target device – same code, same chip, etc. – with full knowledge of the secret key

Two steps:

- ▶ Profiling phase:  $P, K$  known  $\implies Z$  known,  $X$  acquired on an open sample
- ▶ Attack phase:  $P$  known,  $X$  acquired on the *target* device,  $K$  guessed

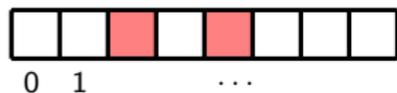
## Outline

1. Context
2. SCA Optimization Problem versus Deep Learning Based SCA
3. NLL Minimization is PI Maximization
4. Simulation results
5. Experimental results

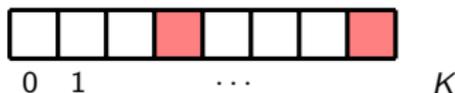
## Profiling Attacks

Key Recovery (*i.e.* attack step)

Given  $N_a$  attack traces  $\mathbf{x}_i$  with plaintext  $p_i$ , calculate scores  $\mathbf{y}_i = F(\mathbf{x}_i)$



$$Z_i = \mathbf{C}(p_i, k^*)$$

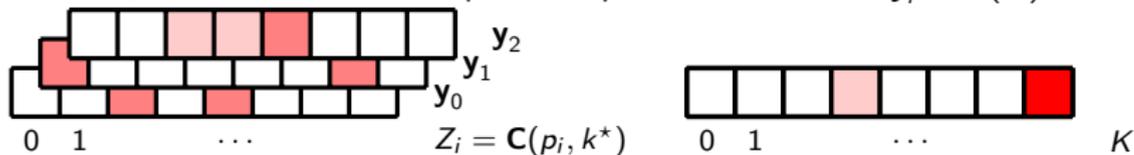




## Profiling Attacks

Key Recovery (*i.e.* attack step)

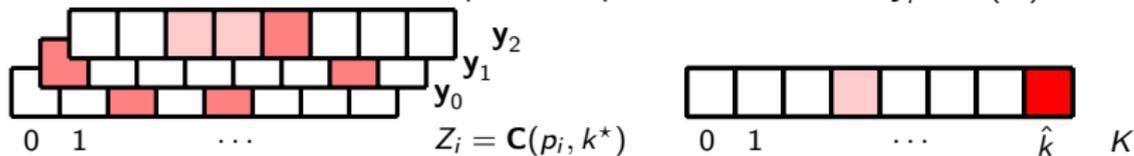
Given  $N_a$  attack traces  $x_i$  with plaintext  $p_i$ , calculate scores  $y_i = F(x_i)$



## Profiling Attacks

Key Recovery (*i.e.* attack step)

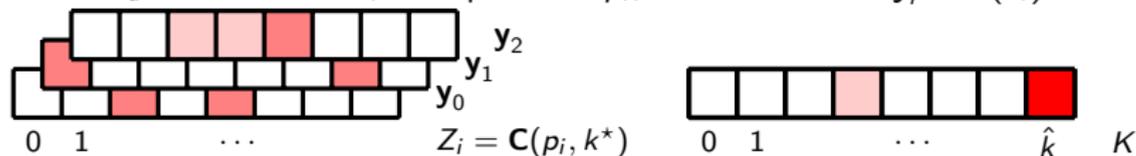
Given  $N_a$  attack traces  $x_i$  with plaintext  $p_i$ , calculate scores  $y_i = F(x_i)$



## Profiling Attacks

Key Recovery (*i.e.* attack step)

Given  $N_a$  attack traces  $\mathbf{x}_i$  with plaintext  $p_i$ , calculate scores  $\mathbf{y}_i = F(\mathbf{x}_i)$

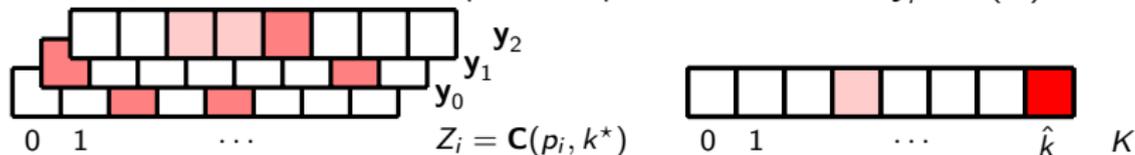


Goal: find  $F$  that minimizes  $N_a$  s.t.  $\hat{k} = k^*$  with probability  $\geq \beta$  (e.g. 0.9)

## Profiling Attacks

Key Recovery (*i.e.* attack step)

Given  $N_a$  attack traces  $\mathbf{x}_i$  with plaintext  $p_i$ , calculate scores  $\mathbf{y}_i = F(\mathbf{x}_i)$



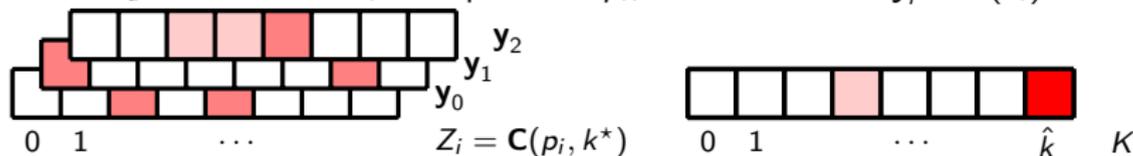
Goal: find  $F$  that minimizes  $N_a$  s.t.  $\hat{k} = k^*$  with probability  $\geq \beta$  (e.g. 0.9)

Optimal model:  $F^*$ , with  $N_a^*$  traces

## Profiling Attacks

Key Recovery (*i.e.* attack step)

Given  $N_a$  attack traces  $\mathbf{x}_i$  with plaintext  $p_i$ , calculate scores  $\mathbf{y}_i = F(\mathbf{x}_i)$



Goal: find  $F$  that minimizes  $N_a$  s.t.  $\hat{k} = k^*$  with probability  $\geq \beta$  (e.g. 0.9)

Optimal model:  $F^*$ , with  $N_a^*$  traces

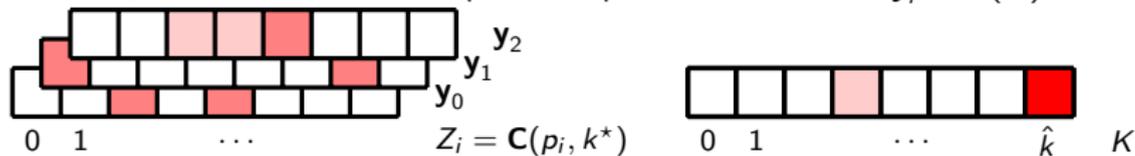
How to find  $F^* \implies$  profiling step

Requires to know the probability distribution  $F^* = \Pr[Z|\mathbf{X}]$

## Profiling Attacks

## Key Recovery (i.e. attack step)

Given  $N_a$  attack traces  $\mathbf{x}_i$  with plaintext  $p_i$ , calculate scores  $\mathbf{y}_i = F(\mathbf{x}_i)$



Goal: find  $F$  that minimizes  $N_a$  s.t.  $\hat{k} = k^*$  with probability  $\geq \beta$  (e.g. 0.9)

Optimal model:  $F^*$ , with  $N_a^*$  traces

How to find  $F^* \implies$  profiling step

Requires to know the probability distribution  $F^* = \Pr[Z|\mathbf{X}]$

Reality: unknown for the evaluator/attacker. Estimation with parametric models  $F(\cdot, \theta)$ :



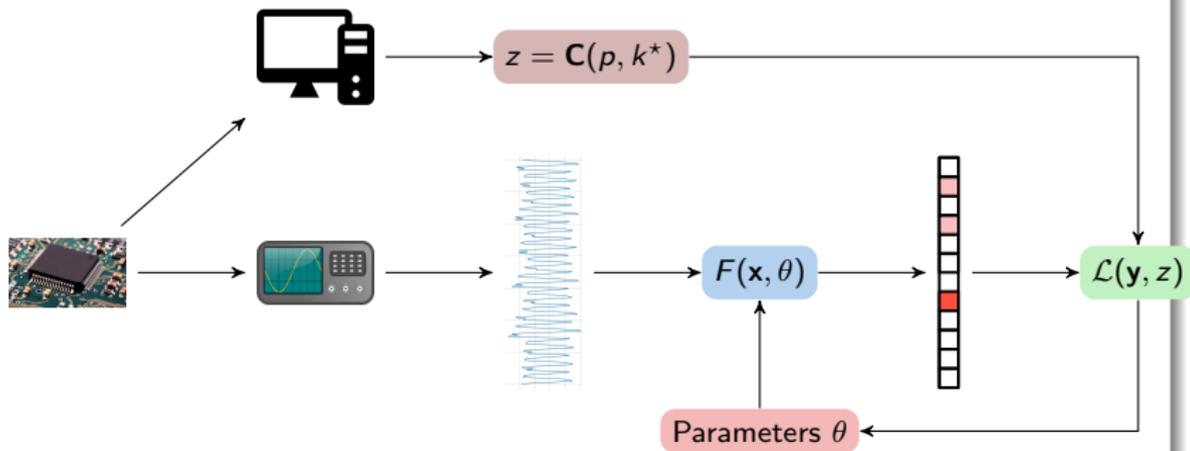
## Deep Learning (DL) based SCA is a hot topic currently

Recent milestones about its effectiveness: more robust against counter-measures like masking [MPP16], jitter (misalignment) [CDP17], whether on software or FPGA [Kim+19]

## Deep Learning (DL) based SCA is a hot topic currently

Recent milestones about its effectiveness: more robust against counter-measures like masking [MPP16], jitter (misalignment) [CDP17], whether on software or FPGA [Kim+19]

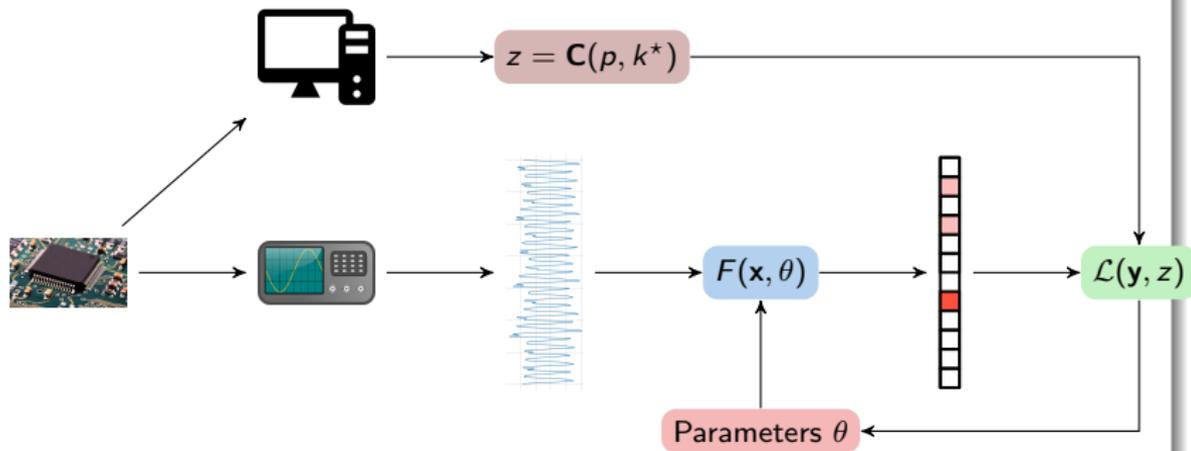
## Training a Neural Network



## Deep Learning (DL) based SCA is a hot topic currently

Recent milestones about its effectiveness: more robust against counter-measures like masking [MPP16], jitter (misalignment) [CDP17], whether on software or FPGA [Kim+19]

## Training a Neural Network



$\mathcal{L}$ : performance metric (accuracy, recall, ...) or loss function (Mean Square Error, NLL, ...)

## Open issue with Machine Learning based SCA<sup>1</sup>

**“How to evaluate the quality of a model during training?”**

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy*  $\implies$  successful key recovery

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy*  $\implies$  successful key recovery
- ▶ *Low accuracy*  $\implies$  nothing

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy*  $\implies$  successful key recovery
- ▶ *Low accuracy*  $\implies$  nothing, problem: often happens (e.g. highly noisy leakages)

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

#### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy*  $\implies$  successful key recovery
- ▶ *Low accuracy*  $\implies$  nothing, problem: often happens (e.g. highly noisy leakages)
- ▶ *Apparently, no other machine learning metric related to SCA metrics*

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

#### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy  $\implies$  successful key recovery*
- ▶ *Low accuracy  $\implies$  nothing, problem: often happens (e.g. highly noisy leakages)*
- ▶ *Apparently, no other machine learning metric related to SCA metrics*

Accuracy: find  $\beta$  s.t.  $N_a^* = 1$

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

#### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy  $\implies$  successful key recovery*
- ▶ *Low accuracy  $\implies$  nothing, problem: often happens (e.g. highly noisy leakages)*
- ▶ *Apparently, no other machine learning metric related to SCA metrics*

Accuracy: find  $\beta$  s.t.  $N_a^* = 1$

$\neq$  SCA: fix  $\beta$  and find  $N_a^*$  instead

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Open issue with Machine Learning based SCA<sup>1</sup>

### “How to evaluate the quality of a model during training?”

- ▶ Accuracy: probability to recover the secret key with one trace

#### Their observations

*“Accuracy does not seem to be the right performance metric in SCA”*

- ▶ *High accuracy  $\implies$  successful key recovery*
- ▶ *Low accuracy  $\implies$  nothing, problem: often happens (e.g. highly noisy leakages)*
- ▶ *Apparently, no other machine learning metric related to SCA metrics*

Accuracy: find  $\beta$  s.t.  $N_a^* = 1$   $\neq$  SCA: fix  $\beta$  and find  $N_a^*$  instead

**Our claim:** we can accurately estimate  $N_a^*$  with DL !

---

<sup>1</sup>Picek *et al.*, CHES 2019 [Pic+18]

## Outline

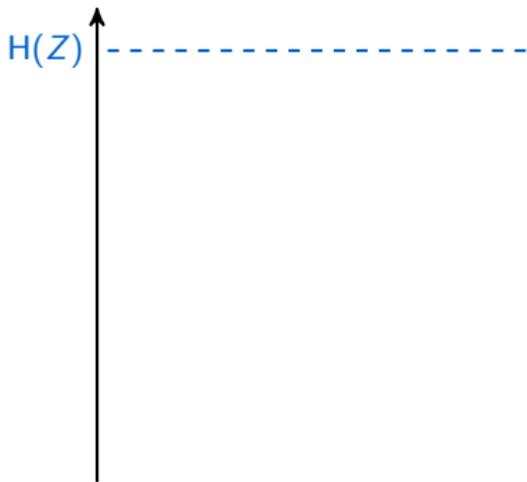
1. Context
2. SCA Optimization Problem versus Deep Learning Based SCA
3. NLL Minimization is PI Maximization
4. Simulation results
5. Experimental results

## Bridging the gap between the loss function and the SCA metric

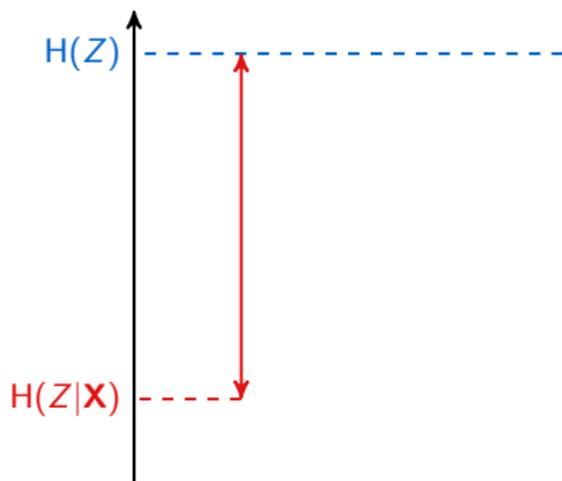
## Bridging the gap between the loss function and the SCA metric



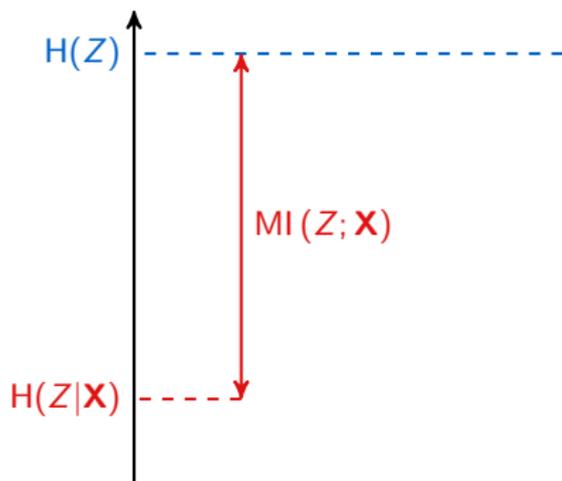
## Bridging the gap between the loss function and the SCA metric



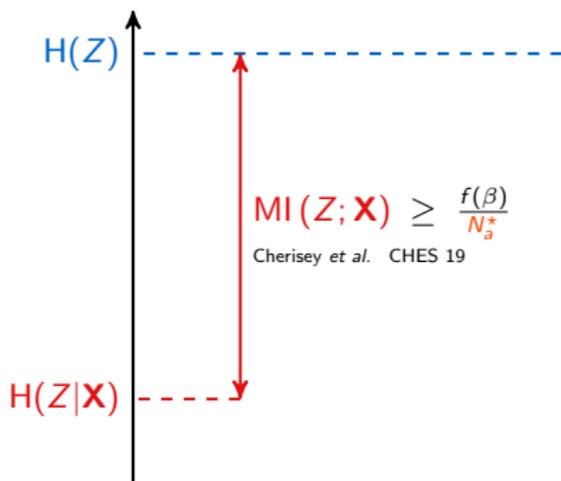
## Bridging the gap between the loss function and the SCA metric



## Bridging the gap between the loss function and the SCA metric

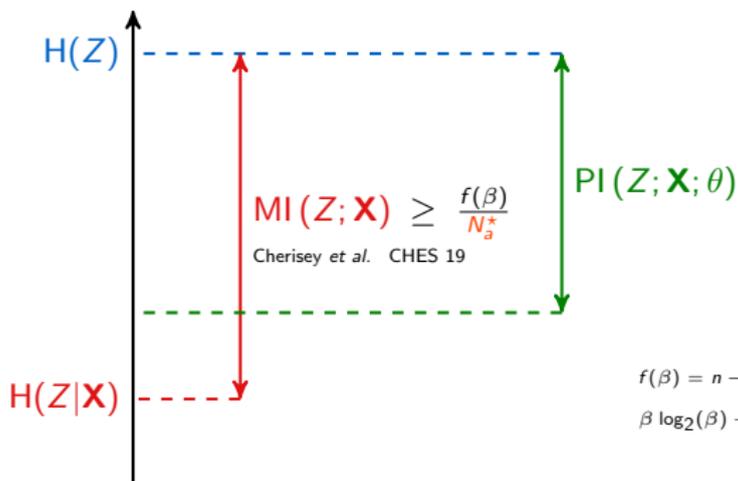


## Bridging the gap between the loss function and the SCA metric

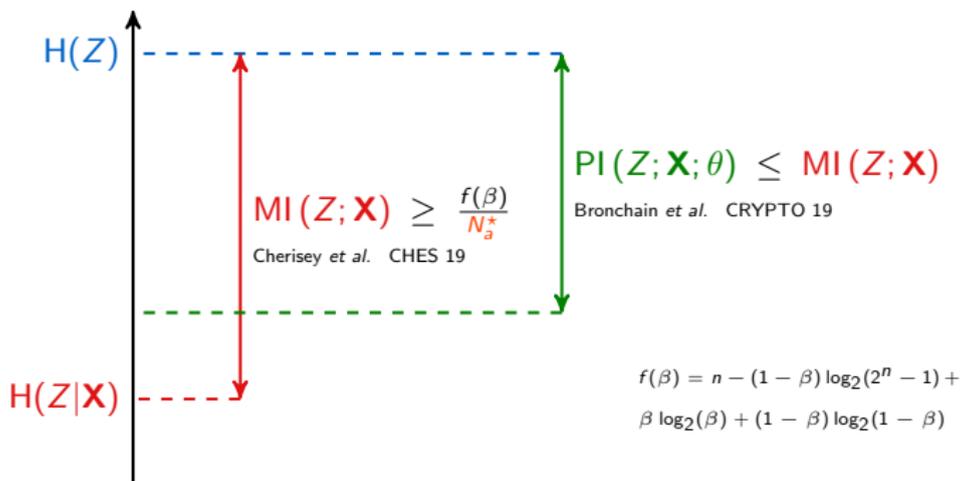


$$f(\beta) = n - (1 - \beta) \log_2(2^n - 1) + \beta \log_2(\beta) + (1 - \beta) \log_2(1 - \beta)$$

## Bridging the gap between the loss function and the SCA metric



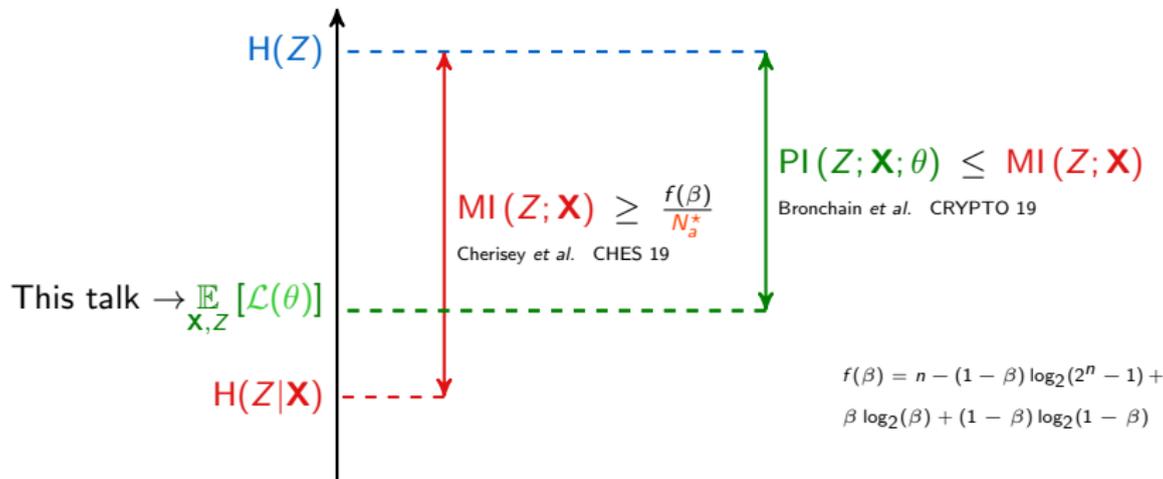
## Bridging the gap between the loss function and the SCA metric



## Bridging the gap between the loss function and the SCA metric

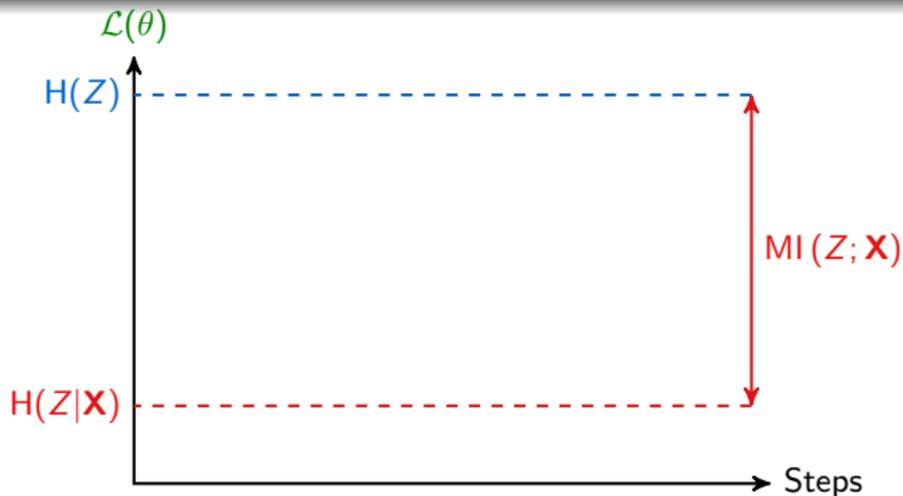
Training: minimization of the NLL a.k.a. Cross Entropy

$$\mathcal{L}(\theta) = \frac{1}{N_p} \sum_{i=1}^{N_p} -\log_2 F(\mathbf{x}_i, \theta)[z_i] = H(Z) - \widehat{PI}_{N_p}(Z; \mathbf{X}; \theta)$$



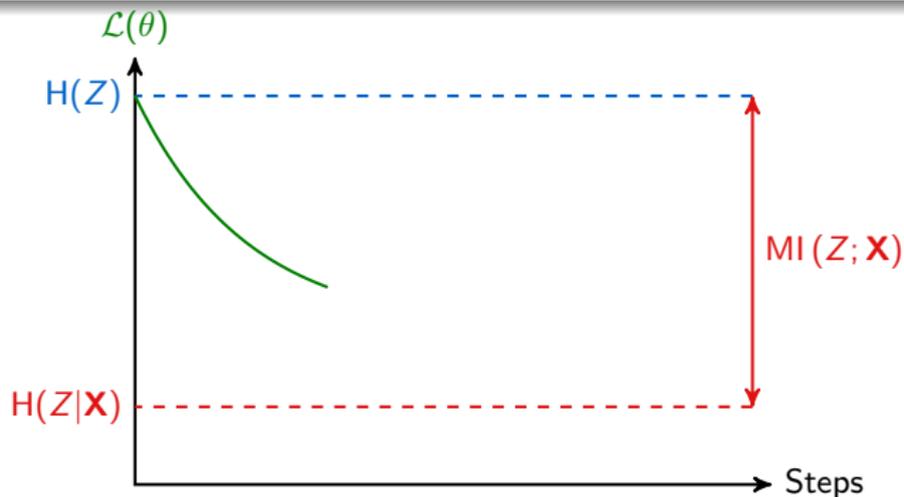
## Main Result

## Proposition



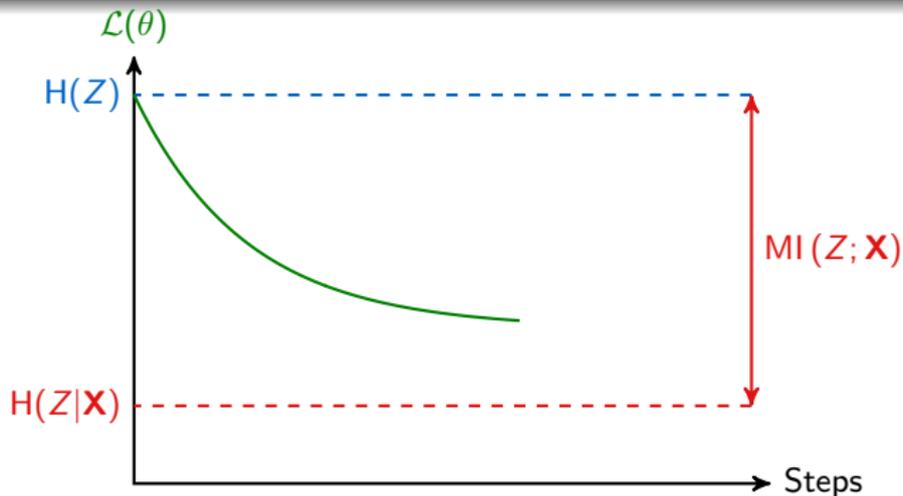
## Main Result

## Proposition



## Main Result

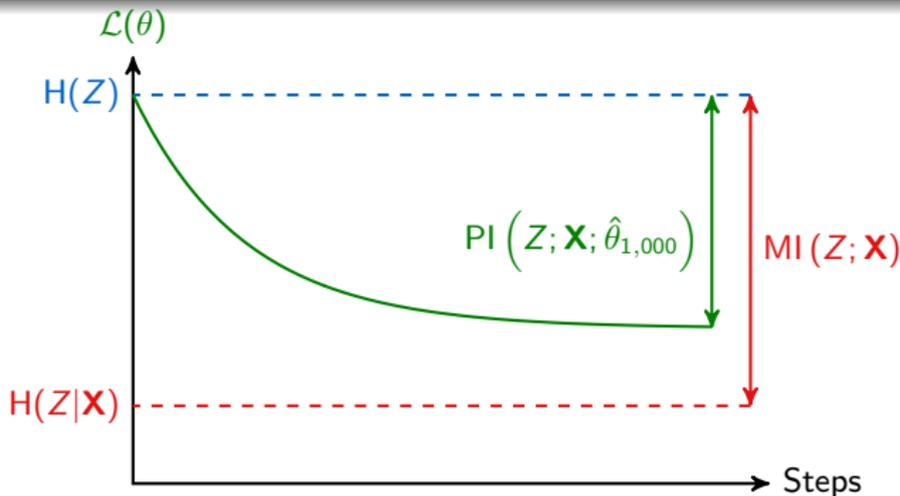
## Proposition



## Main Result

## Proposition

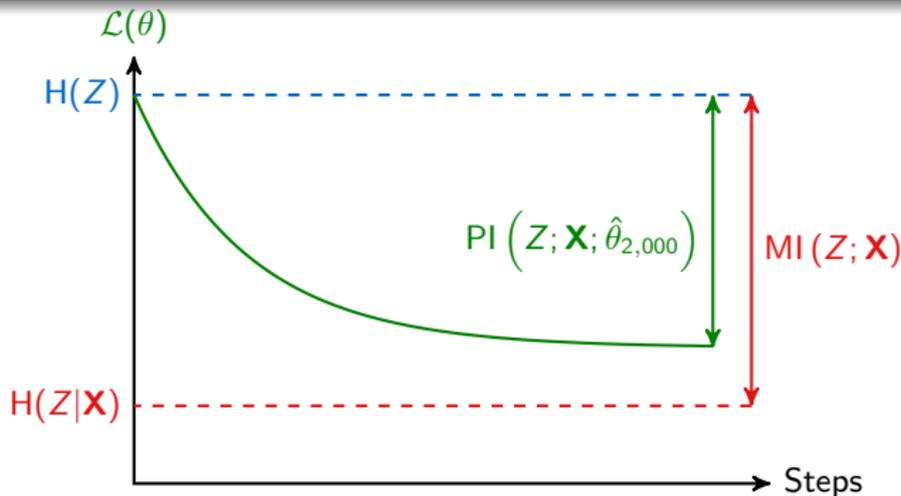
Let  $\hat{\theta}_{N_p} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) = \operatorname{argmax}_{\theta} \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta)$ .



## Main Result

## Proposition

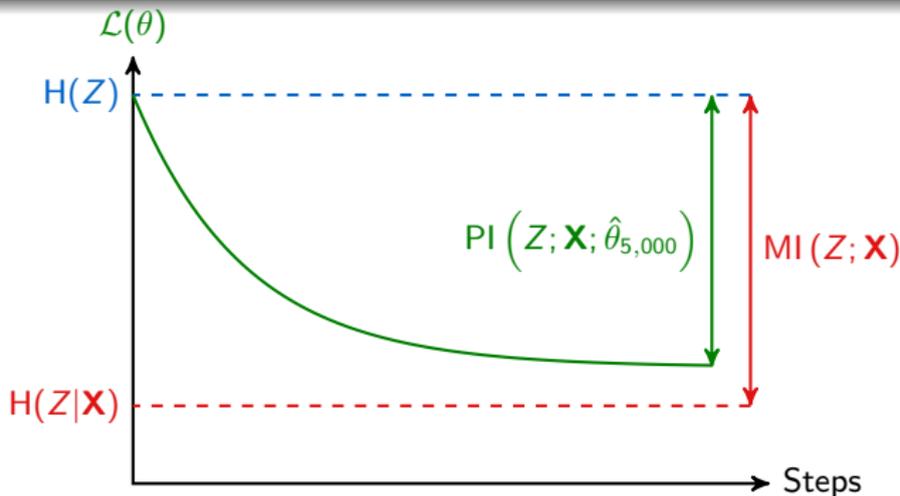
Let  $\hat{\theta}_{N_p} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) = \operatorname{argmax}_{\theta} \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta)$ .



## Main Result

## Proposition

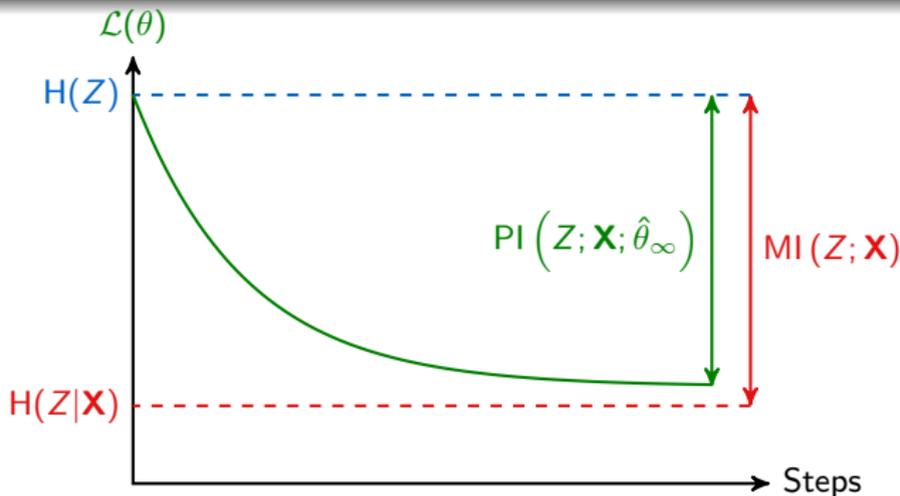
Let  $\hat{\theta}_{N_p} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) = \operatorname{argmax}_{\theta} \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta)$ .



## Main Result

## Proposition

Let  $\hat{\theta}_{N_p} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) = \operatorname{argmax}_{\theta} \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta)$ .

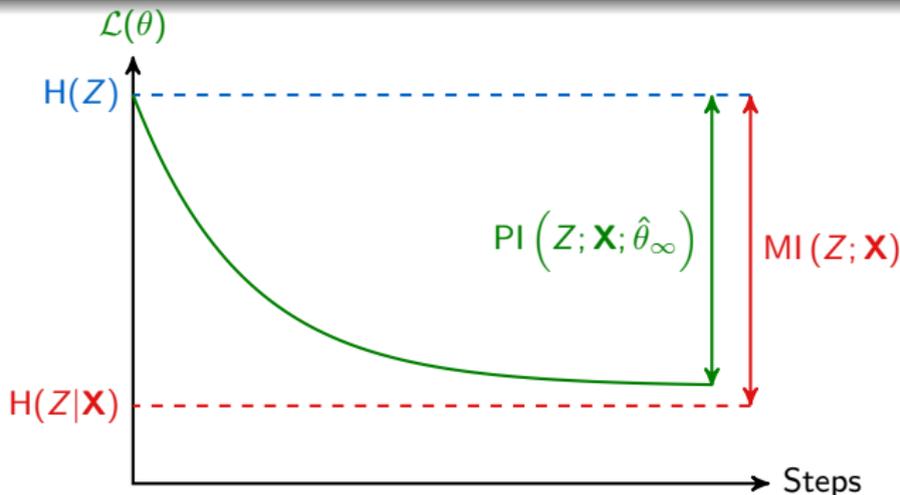


## Main Result

## Proposition

Let  $\hat{\theta}_{N_p} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta) = \operatorname{argmax}_{\theta} \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta)$ . Then:

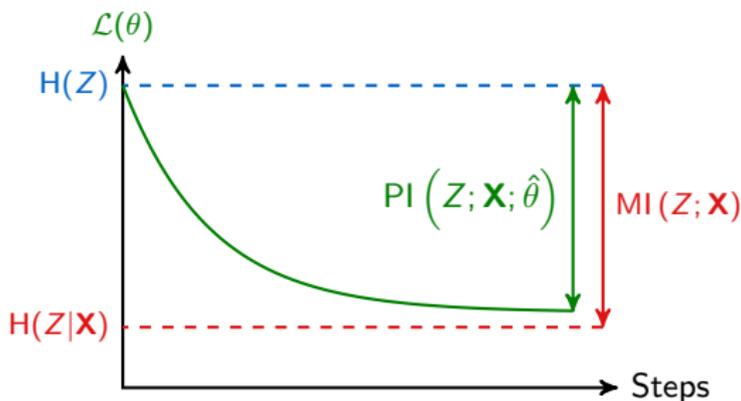
$$\text{PI}(Z; \mathbf{X}; \hat{\theta}_{N_p}) \xrightarrow[N_p \rightarrow \infty]{\mathcal{P}} \sup_{\theta} \text{PI}(Z; \mathbf{X}; \theta) \leq \text{MI}(Z; \mathbf{X})$$



## Tightness of the Lower Bound

To what extent the gap PI/MI is negligible?

Gap composed of three kinds of errors:

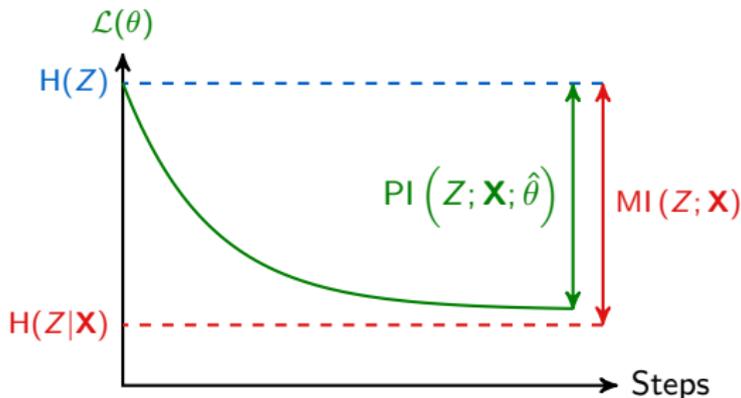


## Tightness of the Lower Bound

To what extent the gap PI/MI is negligible?

Gap composed of three kinds of errors:

- Approximation error:  $\sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) - \text{MI}(Z; \mathbf{X}) \leq 0$

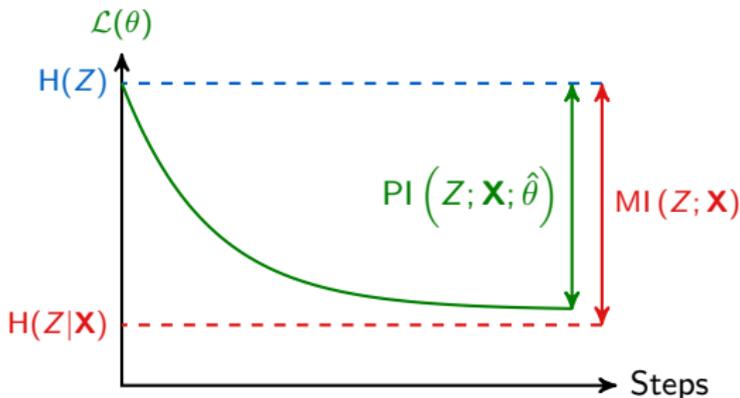


## Tightness of the Lower Bound

To what extent the gap PI/MI is negligible?

Gap composed of three kinds of errors:

- ▶ Approximation error:  $\sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) - \text{MI}(Z; \mathbf{X}) \leq 0$
- ▶ Estimation error:  $N_p < \infty \implies \sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) \rightarrow \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \hat{\theta}_{N_p})$

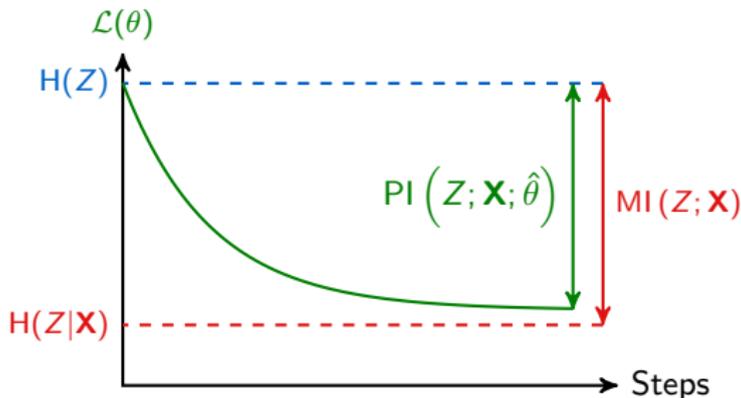


## Tightness of the Lower Bound

To what extent the gap PI/MI is negligible?

Gap composed of three kinds of errors:

- ▶ Approximation error:  $\sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) - \text{MI}(Z; \mathbf{X}) \leq 0$
- ▶ Estimation error:  $N_p < \infty \implies \sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) \rightarrow \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \hat{\theta}_{N_p})$
- ▶ Optimization error:  $\hat{\theta}_{N_p}$  unknown,  $\theta_{SGD}$  instead, by SGD



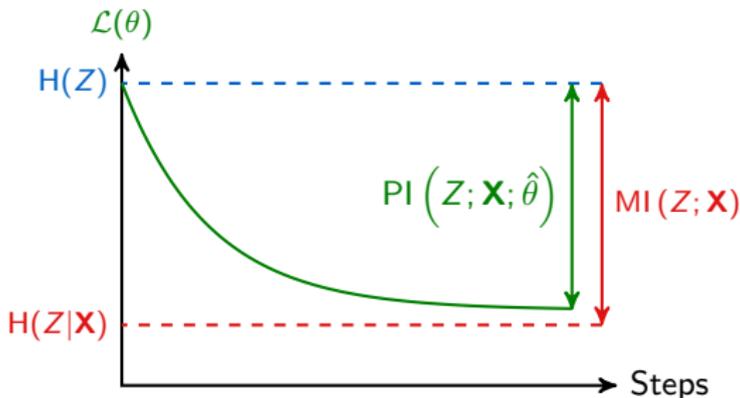
## Tightness of the Lower Bound

To what extent the gap PI/MI is negligible?

Gap composed of three kinds of errors:

- ▶ Approximation error:  $\sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) - \text{MI}(Z; \mathbf{X}) \leq 0$
- ▶ Estimation error:  $N_p < \infty \implies \sup_{\theta \in \Theta} \text{PI}(Z; \mathbf{X}; \theta) \rightarrow \widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \hat{\theta}_{N_p})$
- ▶ Optimization error:  $\hat{\theta}_{N_p}$  unknown,  $\theta_{SGD}$  instead, by SGD

$\implies$  Ideally each error must be discussed through simulations and experiments



## Outline

1. Context
2. SCA Optimization Problem versus Deep Learning Based SCA
3. NLL Minimization is PI Maximization
4. Simulation results
5. Experimental results

## Settings of the simulations

### Leakage model

- ▶ Hamming weight with additive gaussian noise ( $\sigma \in [0.01; 3.2]$ )
- ▶ Draw an *Exhaustive* dataset: estimation error negligible

## Settings of the simulations

### Leakage model

- ▶ Hamming weight with additive gaussian noise ( $\sigma \in [0.01; 3.2]$ )
- ▶ Draw an *Exhaustive* dataset: estimation error negligible

### PI/MI computation

- ▶ Computation of  $MI(\mathbf{X}; \mathbf{Z})$  with Monte-Carlo simulations

## Settings of the simulations

### Leakage model

- ▶ Hamming weight with additive gaussian noise ( $\sigma \in [0.01; 3.2]$ )
- ▶ Draw an *Exhaustive* dataset: estimation error negligible

### PI/MI computation

- ▶ Computation of  $MI(\mathbf{X}; \mathbf{Z})$  with Monte-Carlo simulations
- ▶ Training of a one layer MLP with 1,000 neurons to maximize  $PI(\mathbf{Z}; \mathbf{X}; \theta) = n - \mathcal{L}(\theta)$ , where  $n = 4$  bits

## Settings of the simulations

### Leakage model

- ▶ Hamming weight with additive gaussian noise ( $\sigma \in [0.01; 3.2]$ )
- ▶ Draw an *Exhaustive* dataset: estimation error negligible

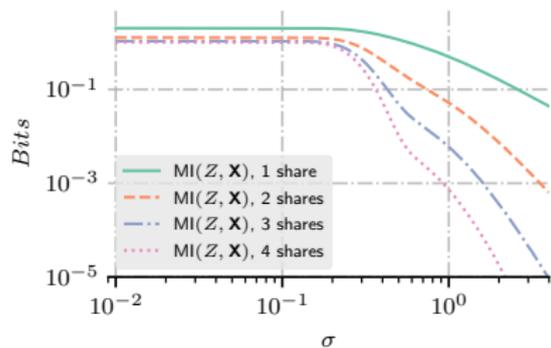
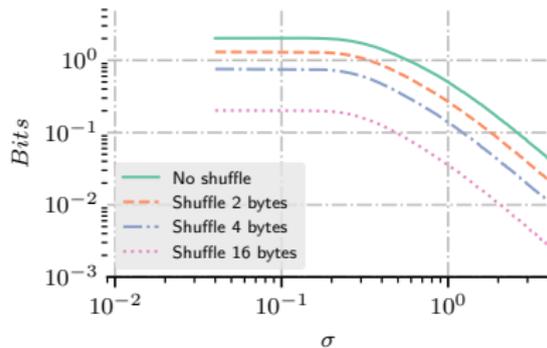
### PI/MI computation

- ▶ Computation of  $MI(\mathbf{X}; \mathbf{Z})$  with Monte-Carlo simulations
- ▶ Training of a one layer MLP with 1,000 neurons to maximize  $PI(\mathbf{Z}; \mathbf{X}; \theta) = n - \mathcal{L}(\theta)$ , where  $n = 4$  bits

### Several case studies

- ▶ Higher-order masking: sensitive variable split into  $d$  independent parts
- ▶ Shuffling: independent operations (e.g. 16 SBoxes in SubBytes) randomly shuffled

## Simulation results

**Figure:** H-O masking, w.r.t. level of noise**Figure:** Shuffling, w.r.t. level of noise

## Simulation results

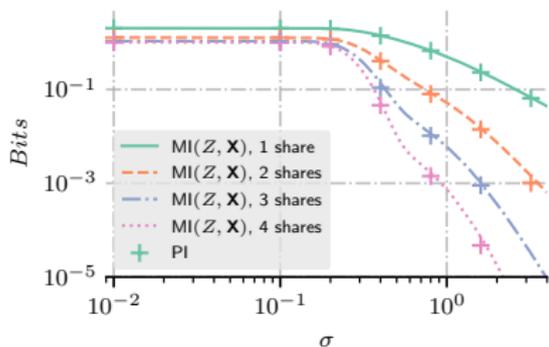


Figure: H-O masking, w.r.t. level of noise

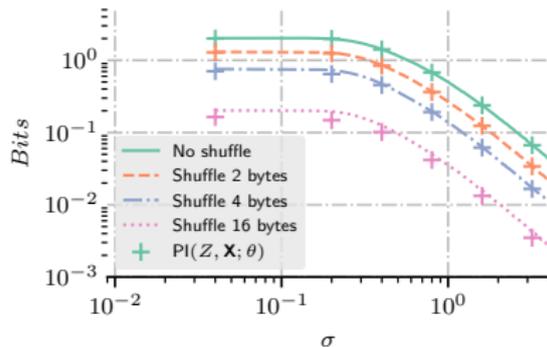


Figure: Shuffling, w.r.t. level of noise

## What to interpret

- ▶ No matter the masking order,  $PI(Z; X; \theta_{SGD}) \approx MI(Z; X)$

## Simulation results

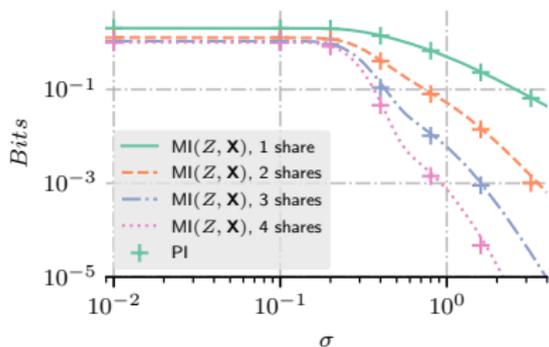


Figure: H-O masking, w.r.t. level of noise

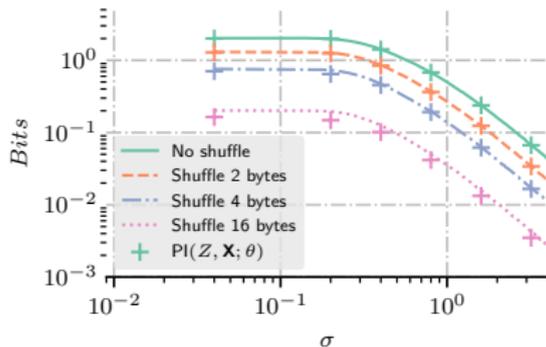


Figure: Shuffling, w.r.t. level of noise

## What to interpret

- ▶ No matter the masking order,  $PI(Z; X; \theta_{SGD}) \approx MI(Z; X)$
- ▶ For a simple MLP, the approximation error and the optimization error are negligible

## Simulation results

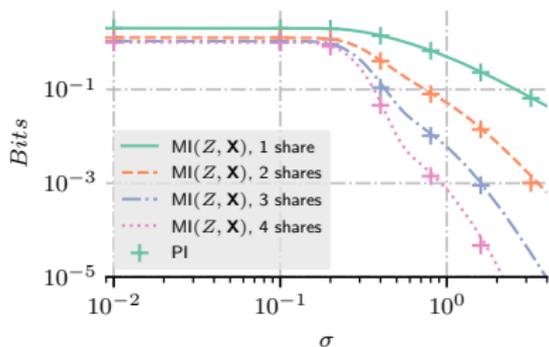


Figure: H-O masking, w.r.t. level of noise

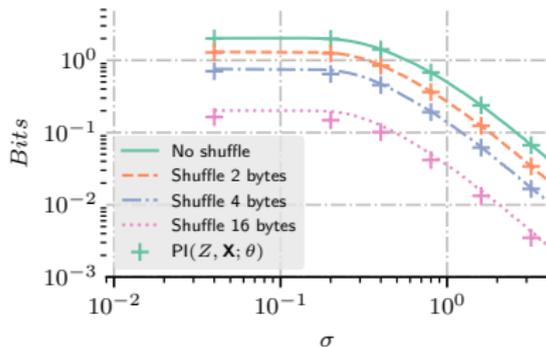


Figure: Shuffling, w.r.t. level of noise

## What to interpret

- ▶ No matter the masking order,  $PI(Z; X; \theta_{SGD}) \approx MI(Z; X)$
- ▶ For a simple MLP, the approximation error and the optimization error are negligible
- ▶ Any more *complex* model should have a negligible approximation error too
- ▶ Empirical verifications: see appendix

## Outline

1. Context
2. SCA Optimization Problem versus Deep Learning Based SCA
3. NLL Minimization is PI Maximization
4. Simulation results
5. Experimental results

## Application on Public Datasets

- ▶  $N_a(\theta)$ : number of traces obtained with key recovery.

## Application on Public Datasets

- ▶  $N_a(\theta)$ : number of traces obtained with key recovery.
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{SGD}) \approx \text{MI}(Z; \mathbf{X})$

## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{SGD}) \approx \text{MI}(Z; \mathbf{X})$

## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{SGD}) \approx \text{MI}(Z; \mathbf{X})$
- ▶ Tests on public datasets, using architectures proposed in recent papers [MDP19; Kim+19]

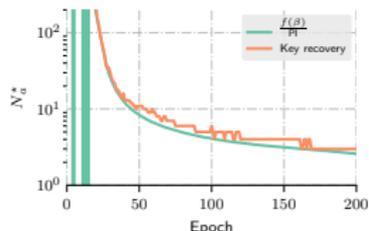
## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{SGD}) \approx \text{MI}(Z; \mathbf{X})$
- ▶ Tests on public datasets, using architectures proposed in recent papers [MDP19; Kim+19]
- ▶ Relative error  $\epsilon$  computed at final epoch

## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{SGD}) \approx \text{MI}(Z; \mathbf{X})$
- ▶ Tests on public datasets, using architectures proposed in recent papers [MDP19; Kim+19]
- ▶ Relative error  $\epsilon$  computed at final epoch

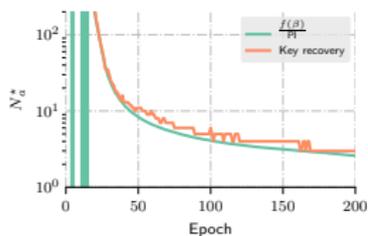
## Micro-controller protected with misalignment

Figure: AES-RD:  $\epsilon = 0.16$

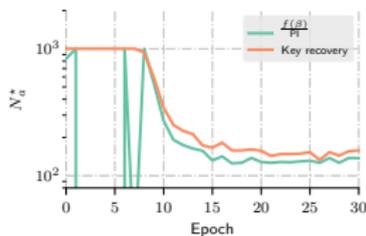
## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{\text{SGD}}) \approx \text{MI}(Z; \mathbf{X})$
- ▶ Tests on public datasets, using architectures proposed in recent papers [MDP19; Kim+19]
- ▶ Relative error  $\epsilon$  computed at final epoch

Micro-controller protected with misalignment

Figure: AES-RD:  $\epsilon = 0.16$ 

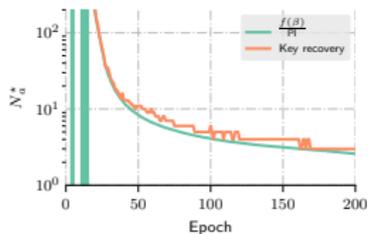
Micro-controller protected with masking

Figure: ASCAD:  $\epsilon = 0.16$

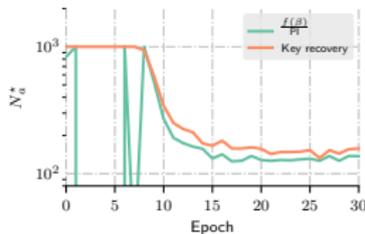
## Application on Public Datasets

- ▶  $N_a(\theta) \frac{f(\beta)}{\text{PI}(Z; \mathbf{X}; \theta)} \approx \frac{f(\beta)}{n - \mathcal{L}(\theta)}$ : number of traces obtained with key recovery?
- ▶ So far:  $N_a^* \geq \frac{f(\beta)}{\text{MI}(Z; \mathbf{X})}$  and  $\text{PI}(Z; \mathbf{X}; \theta_{\text{SGD}}) \approx \text{MI}(Z; \mathbf{X})$
- ▶ Tests on public datasets, using architectures proposed in recent papers [MDP19; Kim+19]
- ▶ Relative error  $\epsilon$  computed at final epoch

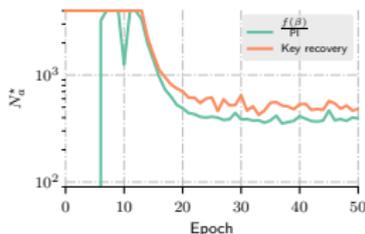
Micro-controller protected with misalignment

Figure: AES-RD:  $\epsilon = 0.16$ 

Micro-controller protected with masking

Figure: ASCAD:  $\epsilon = 0.16$ 

Implementation on FPGA (no counter-measure)

Figure: AES-HD:  $\epsilon = 0.18$

## Conclusion

Takeaway messages

## Conclusion

### Takeaway messages

1. Minimizing the **NLL loss**  $\equiv$  maximizing the **PI**  $\implies$  tight lower bound of the **MI**  $\implies$  accurate estimation of  $N_a^*$

## Conclusion

### Takeaway messages

1. Minimizing the **NLL loss**  $\equiv$  maximizing the **PI**  $\implies$  tight lower bound of the **MI**  $\implies$  accurate estimation of  $N_a^*$
2. NLL as a loss function is sound from an evaluator point of view

## Conclusion

### Takeaway messages

1. Minimizing the **NLL loss**  $\equiv$  maximizing the **PI**  $\implies$  tight lower bound of the **MI**  $\implies$  accurate estimation of  $N_a^*$
2. NLL as a loss function is sound from an evaluator point of view
3. Enables to quantitatively measure the impact of counter-measures

Thank You!

Questions?

Looking for a postdoc candidate in machine-learning-based SCA? Hire me!



## References I

- [CDP17] Eleonora Cagli, Cécile Dumas, and Emmanuel Prouff. “Convolutional Neural Networks with Data Augmentation Against Jitter-Based Countermeasures - Profiling Attacks Without Pre-processing”. In: *Cryptographic Hardware and Embedded Systems - CHES 2017 - 19th International Conference, Taipei, Taiwan, September 25-28, 2017, Proceedings*. Ed. by Wieland Fischer and Naofumi Homma. Vol. 10529. Lecture Notes in Computer Science. Springer, 2017, pp. 45–68. ISBN: 978-3-319-66786-7. DOI: 10.1007/978-3-319-66787-4\\_3. URL: [https://doi.org/10.1007/978-3-319-66787-4\\\_3](https://doi.org/10.1007/978-3-319-66787-4\_3).
- [Kim+19] Jaehun Kim et al. “Make Some Noise. Unleashing the Power of Convolutional Neural Networks for Profiled Side-channel Analysis”. In: *IACR Transactions on Cryptographic Hardware and Embedded Systems 2019.3* (2019), pp. 148–179. DOI: 10.13154/tches.v2019.i3.148-179. URL: <https://tches.iacr.org/index.php/TCHES/article/view/8292>.

## References II

- [MPP16] Housseem Maghrebi, Thibault Portigliatti, and Emmanuel Prouff. “Breaking Cryptographic Implementations Using Deep Learning Techniques”. In: *Security, Privacy, and Applied Cryptography Engineering - 6th International Conference, SPACE 2016, Hyderabad, India, December 14-18, 2016, Proceedings*. Ed. by Claude Carlet, M. Anwar Hasan, and Vishal Saraswat. Vol. 10076. Lecture Notes in Computer Science. Springer, 2016, pp. 3–26. ISBN: 978-3-319-49444-9. DOI: 10.1007/978-3-319-49445-6\\_1. URL: [https://doi.org/10.1007/978-3-319-49445-6\\\_1](https://doi.org/10.1007/978-3-319-49445-6\_1).
- [MDP19] Loïc Masure, Cécile Dumas, and Emmanuel Prouff. “Gradient Visualization for General Characterization in Profiling Attacks”. In: *Constructive Side-Channel Analysis and Secure Design - 10th International Workshop, COSADE 2019, Darmstadt, Germany, April 3-5, 2019, Proceedings*. Ed. by Ilia Polian and Marc Stöttinger. Vol. 11421. Lecture Notes in Computer Science. Springer, 2019, pp. 145–167. ISBN: 978-3-030-16349-5. DOI: 10.1007/978-3-030-16350-1\\_9. URL: [https://doi.org/10.1007/978-3-030-16350-1\\\_9](https://doi.org/10.1007/978-3-030-16350-1\_9).

## References III

- [Pic+18] Stjepan Picek et al. “The Curse of Class Imbalance and Conflicting Metrics with Machine Learning for Side-channel Evaluations”. In: *IACR Transactions on Cryptographic Hardware and Embedded Systems* 2019.1 (2018), pp. 209–237. DOI: [10.13154/tches.v2019.i1.209-237](https://doi.org/10.13154/tches.v2019.i1.209-237). URL: <https://tches.iacr.org/index.php/TCHES/article/view/7339>.
- [Vey+12] Nicolas Veyrat-Charvillon et al. “Shuffling against Side-Channel Attacks: A Comprehensive Study with Cautionary Note”. In: *Advances in Cryptology - ASIACRYPT 2012 - 18th International Conference on the Theory and Application of Cryptology and Information Security, Beijing, China, December 2-6, 2012. Proceedings*. Ed. by Xiaoyun Wang and Kazue Sako. Vol. 7658. Lecture Notes in Computer Science. Springer, 2012, pp. 740–757. ISBN: 978-3-642-34960-7. DOI: [10.1007/978-3-642-34961-4\\_44](https://doi.org/10.1007/978-3-642-34961-4_44). URL: [https://doi.org/10.1007/978-3-642-34961-4\\_44](https://doi.org/10.1007/978-3-642-34961-4_44).

## Our home dataset

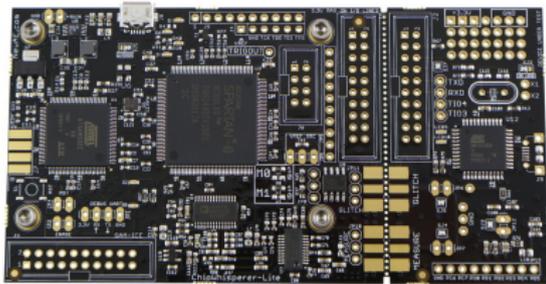


Figure: ChipWhisperer-Lite board

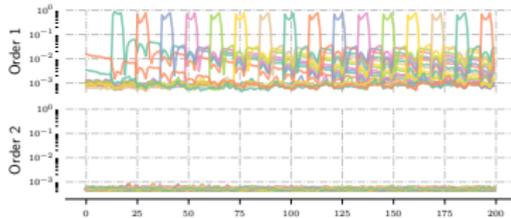


Figure: SNR at orders  $d = 1, 2$

---

### Algorithm 1 loadData

---

```
1: LD r0, X           ▷ Loads the first byte in r0
2: CLR r0             ▷ Clears the register
3: ST X, r0           ▷ Stores 0 in the plaintext array
4: LD r0, X           ▷ Do it again to clear the bus
5: CLR r0
6: ST X, r0
7: LD r0, X           ▷ One more time to be sure
8: CLR r0
9: ST X+, r0
```

---

Loads sequentially an array of 16 bytes into one register and clears it  $\implies$  no joint leakage at order  $d \geq 2$ .

500,000 traces acquired.

We only work on  $n = 4$  bits,

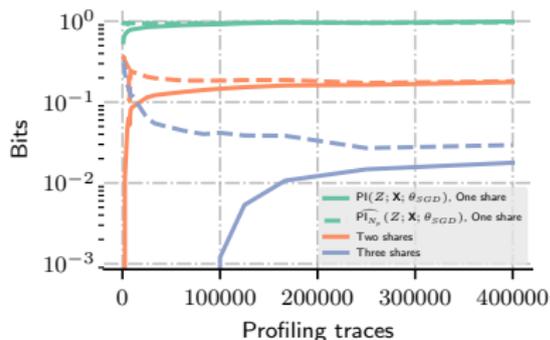
$|\mathcal{Z}| = 2^n = 16$ .

## Experiment on ChipWhisperer-Lite: masking

- ▶ Emulation of order  $d$  leakages:  
 $Z = \bigoplus_{i \in [0, d]} \text{plain}[i]$  for  
 $d \in \{0, 1, 2\}$
- ▶ Extraction of Pols according to SNR.
- ▶ Learning curve:  $\text{PI}(Z; \mathbf{X}; \theta_{SGD})$   
and  $\widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta_{SGD})$  w.r.t.  $N_p$   
 $\implies$  measures the estimation error.

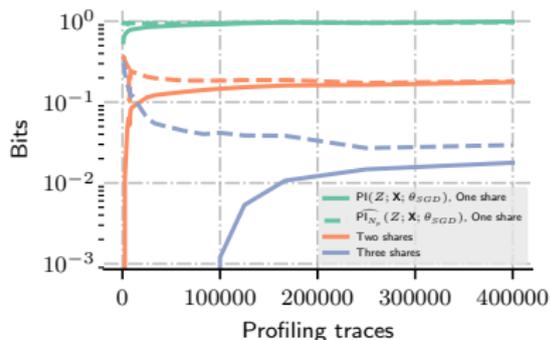
## Experiment on ChipWhisperer-Lite: masking

- ▶ Emulation of order  $d$  leakages:  
 $Z = \bigoplus_{i \in [0, d]} \text{plain}[i]$  for  
 $d \in \{0, 1, 2\}$
- ▶ Extraction of Pols according to SNR.
- ▶ Learning curve:  $\text{PI}(Z; \mathbf{X}; \theta_{SGD})$   
 and  $\widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta_{SGD})$  w.r.t.  $N_p$   
 $\implies$  measures the estimation error.



## Experiment on ChipWhisperer-Lite: masking

- ▶ Emulation of order  $d$  leakages:  
 $Z = \bigoplus_{i \in [0, d]} \text{plain}[i]$  for  
 $d \in \{0, 1, 2\}$
- ▶ Extraction of Pols according to SNR.
- ▶ Learning curve:  $\text{PI}(Z; \mathbf{X}; \theta_{SGD})$   
 and  $\widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta_{SGD})$  w.r.t.  $N_p$   
 $\implies$  measures the estimation error.

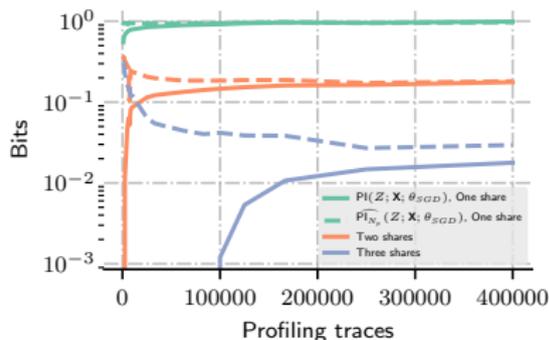


## What to interpret

- ▶  $\approx$  one decade lost for each new masking order  $\implies$  masking remains sound

## Experiment on ChipWhisperer-Lite: masking

- ▶ Emulation of order  $d$  leakages:  
 $Z = \bigoplus_{i \in [0, d]} \text{plain}[i]$  for  
 $d \in \{0, 1, 2\}$
- ▶ Extraction of Pols according to SNR.
- ▶ Learning curve:  $\text{PI}(Z; \mathbf{X}; \theta_{SGD})$   
 and  $\widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta_{SGD})$  w.r.t.  $N_p$   
 $\implies$  measures the estimation error.

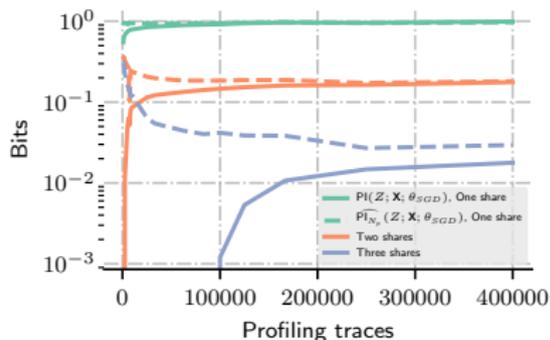


## What to interpret

- ▶  $\approx$  one decade lost for each new masking order  $\implies$  masking remains sound
- ▶ Masking has an effect on the estimation error

## Experiment on ChipWhisperer-Lite: masking

- ▶ Emulation of order  $d$  leakages:  
 $Z = \bigoplus_{i \in [0, d]} \text{plain}[i]$  for  
 $d \in \{0, 1, 2\}$
- ▶ Extraction of Pols according to SNR.
- ▶ Learning curve:  $\text{PI}(Z; \mathbf{X}; \theta_{SGD})$   
 and  $\widehat{\text{PI}}_{N_p}(Z; \mathbf{X}; \theta_{SGD})$  w.r.t.  $N_p$   
 $\implies$  measures the estimation error.



## What to interpret

- ▶  $\approx$  one decade lost for each new masking order  $\implies$  masking remains sound
- ▶ Masking has an effect on the estimation error
- ▶ For  $d = 3$ ,  $N_p < 100,000$ , no information !

## Experiment 5: shuffling

- ▶ Emulation of order  $c$  shuffling:  
 $Z = \text{plain}[i]$  where  $i$  is randomly drawn from a subset of  $c$  indices
- ▶ Complete trace:  $D = 250$

## Experiment 5: shuffling

- ▶ Emulation of order  $c$  shuffling:  
 $Z = \text{plain}[i]$  where  $i$  is randomly drawn from a subset of  $c$  indices
- ▶ Complete trace:  $D = 250$

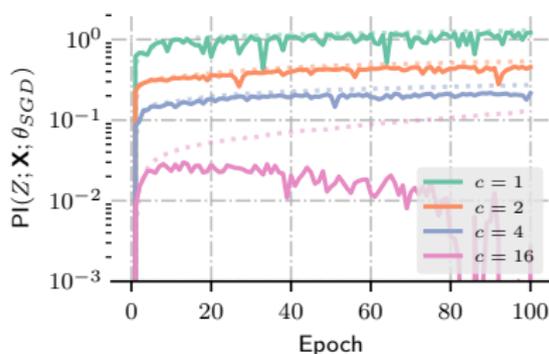


Figure: Exp. 5, shuffling

## Experiment 5: shuffling

- ▶ Emulation of order  $c$  shuffling:  
 $Z = \text{plain}[i]$  where  $i$  is randomly drawn from a subset of  $c$  indices
- ▶ Complete trace:  $D = 250$

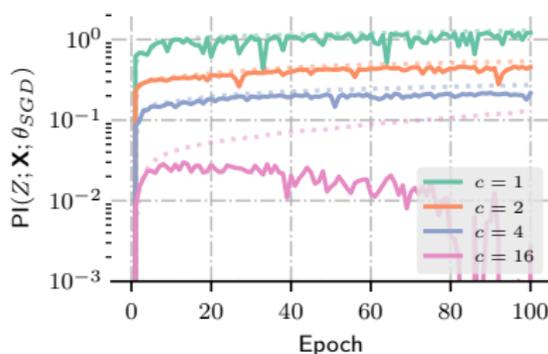


Figure: Exp. 5, shuffling

## What to interpret

- ▶ Linear decrease of PI, as expected [Vey+12]

## Experiment 5: shuffling

- ▶ Emulation of order  $c$  shuffling:  
 $Z = \text{plain}[i]$  where  $i$  is randomly drawn from a subset of  $c$  indices
- ▶ Complete trace:  $D = 250$

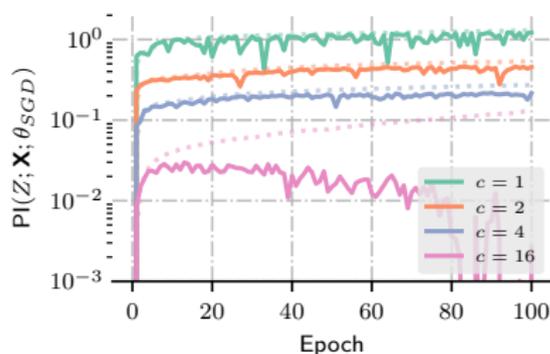


Figure: Exp. 5, shuffling

## What to interpret

- ▶ Linear decrease of PI, as expected [Vey+12]
- ▶ Clearly over-fitting: the estimation error non-negligible