

Perceived Information Revisited: New Metrics to Evaluate Success Rate of Side-Channel Attacks

Akira Ito¹, Rei Ueno², Naofumi Homma²

- 1 NTT Social Informatics Laboratory
- 2 Tohoku University

Background of this work



- DL-SCA is one of the most powerful attacks.
 - Many studies on DL-SCA have been conducted recently.
- Training an NN model requires a performance metric.



Major metrics (e.g., CE loss, acc.) are not suitable for SCA.

• Accuracy of 0% does not mean DL-SCA will fail.

However, computation cost of success rate (SR) is too high!

Contributions



- Analysis of relation between cross entropy (CE) loss function and SR
 - Explain why CE loss is not suitable to measure the performance of DL-SCA.

- Effective CE/PI (ECE/EPI), new metrics for DL-SCA
 - ECE/EPI are more useful metrics than CE/PI for SCAs.
 - EPI can enable us to estimate (the upper-bound of) SR.

Relation between NLL and MI



- Negative log likelihood (NLL) is used as loss function.

 - $NLL = -\frac{1}{m} \sum_{i=1}^{m} \log q(Z_i | X_i; \theta)$ NLL minimization is equivalent to maximum likelihood estimation.
- NLL can be regarded as approximation of CE.
 - If the number of traces m is sufficiently large, then •

NLL $\approx -\mathbb{E} \log q(Z|X;\theta) = \operatorname{CE}(q)$

Relation between mutual information (MI) and CE

 $I(Z; \mathbf{X}) \ge H(K) - CE(q) \approx H(K) - NLL$ Perceived information (PI) $J_q(Z; \mathbf{X}) = H(K) - CE(q)$ denotes how much information NN can extract.



de Chérisey et al. prove the following theorem.



Side-channel can be seen as communication channel.





de Chérisey et al. prove the following theorem.





de Chérisey et al. prove the following theorem.





de Chérisey et al. prove the following theorem.



Side-channel can be seen as communication channel.



Extension for DL-SCAs



■ Intuitively, we expect the following inequality holds:

 $\underline{\xi(\operatorname{SR}_m(q))} \le \underline{mJ_q(Z; \boldsymbol{X})} = m(H(K) - \operatorname{CE}(q))$

How much entropy does attacker need when using **model** *q* and *m* traces?

Amount of information **model** q can extract with m traces

- If this holds, we can estimate SR by using PI (i.e., CE)
 - > Masure et al. experimentally showed that this inequality would hold.

Unfortunately, this does not hold.

Theorem (probability distribution conversion which retains SR)

Let q be a model. Define a following conversion of q with an inverse temperature $\beta > 0$: $q_{\beta}(z \mid \boldsymbol{x}; \theta) = \frac{q(z \mid \boldsymbol{x}; \theta)^{\beta}}{\sum_{z'} q(z' \mid \boldsymbol{x}; \theta)^{\beta}}$

For any $\beta > 0$, the success rate using q is equivalent to that using q_{β} .

Results of conversion using β





■ NLL (CE) value and distribution shape change with β .

But, SR/GE does not change with β .

• There is counterexample q_{β} of following inequality:

 $\xi(\operatorname{SR}_m(q)) \le m J_q(Z; \boldsymbol{X}) = m(H(K) - \operatorname{CE}(q))$

Effective CE/PI (ECE/EPI)



SR is invariant, but CE/PI varies with the value of β .

• CE/PI are not appropriate metrics for DL-SCA.

Proposed metrics: ECE and EPI (effective PI)

$$\operatorname{CE}^*(q) \coloneqq \inf_{\beta \in (0,\infty)} \operatorname{CE}(q_\beta),$$

$$J_q^*(Z; \boldsymbol{X}) \coloneqq \sup_{\beta \in (0, \infty)} J_{q_\beta}(Z; \boldsymbol{X}) = H(Z) - CE^*(q)$$

- Basic idea: remove the uncertainty of CE/PI in terms of SR
- Conject following inequality holds using EPI.

Conjecture

DL-SCAs on masked software/hardware implementations





Processing time of each method



Processing time per one epoch [s]

	Empirical SR evaluation	Proposed method	Ratio
ASCAD	14.1	0.0378	373
ті	145	0.531	273

- SR is evaluated by bootstrapping.
 - > Use 100 bootstrap samples to estimate SR value.

Proposed method is several hundreds faster than empirical evaluation.

Concluding remarks



- Analysis of relation between CE loss and SR
 - Conversion changes CE loss but not SR
 - CE/PI has uncertainty in terms of SR
- Effective CE/PI (ECE/EPI), new metrics for DL-SCA
 - Can easily estimate the attack performance (e.g., SR and GE)

- Future work
 - Formal proof of our conjecture (inequality of SR and EPI)



Settings of experiments



	Training	Test
ASCAD	50,000	10,000
TI	4,000,000	4,000,000

Model comparison



- Compare four pretrained models for ASCAD dataset
 - MLP and CNN models proposed in original ASCAD paper
 - CNN models proposed by Zaid et al. and Wouters et al.

Lack of bins means # of required traces is greater than 10,000.



Copyright NTT CORPORATION

How to calculate ECE/EPI

 $CE(q_{\beta})$ has the following properties:

- $\operatorname{CE}(q_{\beta}) \to n \text{ as } \beta \to 0$
- $\operatorname{CE}(q_{\beta}) \to \infty \text{ as } \beta \to \infty$
- $CE(q_{\beta})$ is a strictly convex function of β .

■ Newton method can find the minimum value of $CE(q_\beta)$.

• The local minimum of $CE(q_{\beta})$ is its global minimum.





How to use NN for key recovery

Negative log likelihood (NLL) is used as a score of each key

$$\text{NLL}^{(k)} = -\frac{1}{m} \sum_{i=1}^{m} \log q(S(k \oplus T_i) | \boldsymbol{X}_i; \theta)$$

- NLL is inversely proportional to ٠ the product of probability.
- Attack Procedure:
 - 1. Calculate NLL for each key using m traces
 - 2. Get k whose the minimum NLL value among all candidates



Inference using NNs



■ NN is used to estimate intermediate value from a trace.



- In profiling phase, NN trains plausible probability distribution.
- In attack phase, trained NN estimates secret information.