

Side-Channel Attacks Against the Human Brain: the PIN Code Case Study

J. Lange, **C. Massart**, A. Mouraux and F.-X. Standaert

UCL Crypto Group, Belgium.

Cosade 2017, Paris, France



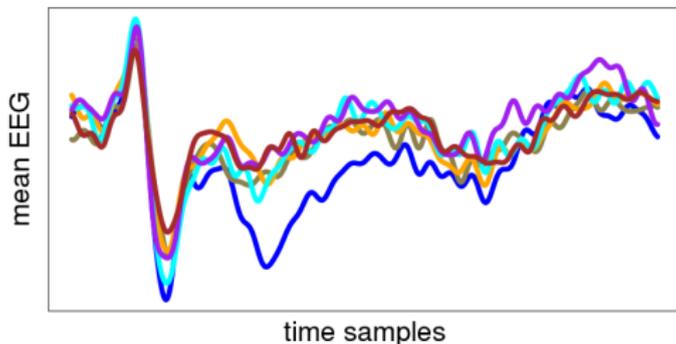
Introduction

- ▶ Side-channel attacks allow PIN code recoveries
 - ▶ e.g., Le Bouder et al., *A Template Attack against Verify PIN Algorithms*, SECRIPT 2016
- ▶ Can we apply them to BCIs & EEG signals?



Motivation

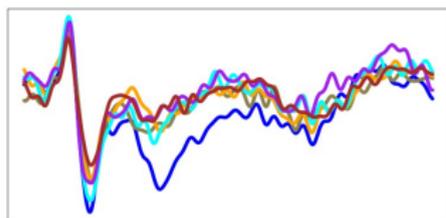
- ▶ Feasibility shown by Martinovic et al. (USENIX 2012)
 - ▶ i.e., there is exploitable information in EEG signals



- ▶ BCIs more and more commercialized (e.g., for gaming)



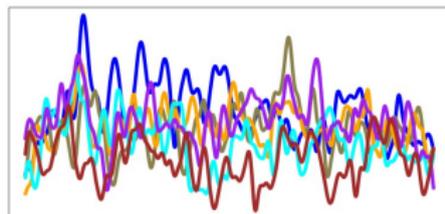
Challenge: low (& irregular) SNR



Range of the
Standard Deviation

\approx

Differences in the
mean traces



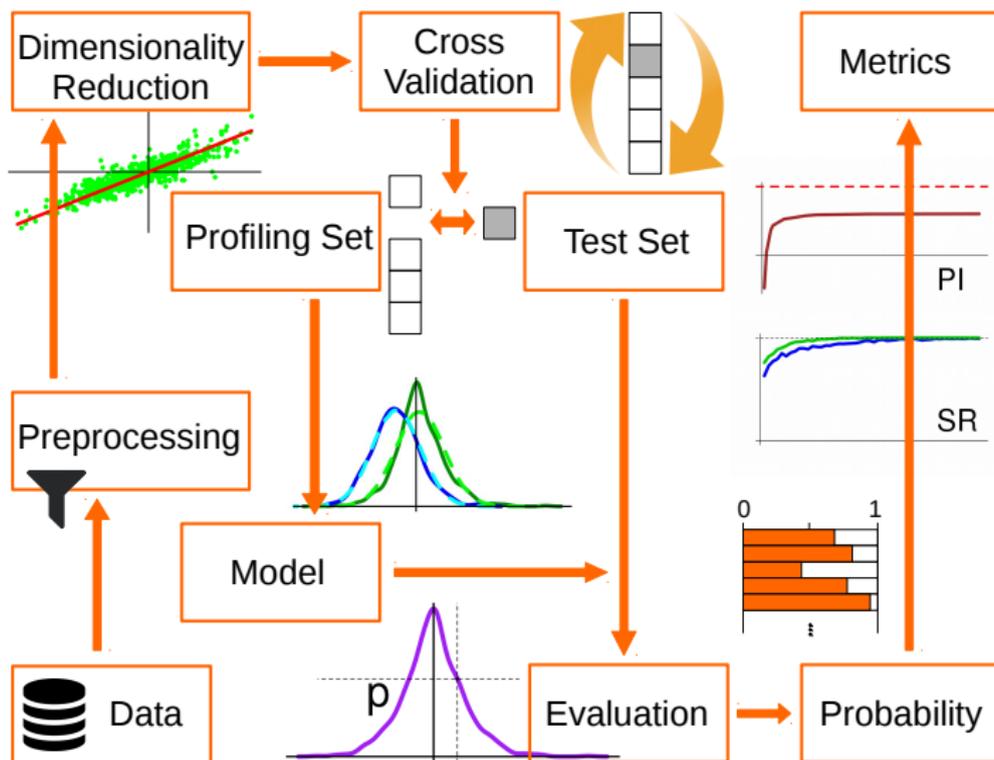
Main questions

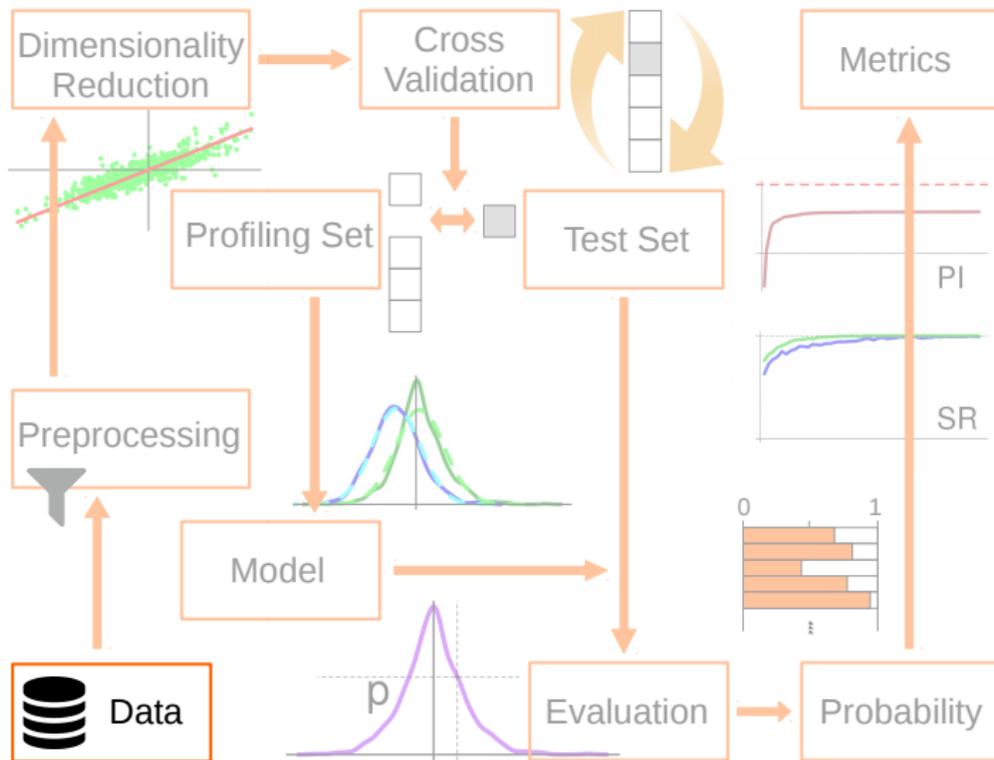
- ▶ Can we extract PINs exactly (or only partially)?
- ▶ Can we extract them with sufficient confidence?
- ▶ How do supervised (aka profiled) and unsupervised (aka non-profiled) attacks compare?
- ▶ *How similar/different are different subjects?*
- ▶ *What are the consequences for security & privacy?*

Note: results can be viewed as positive or negative!

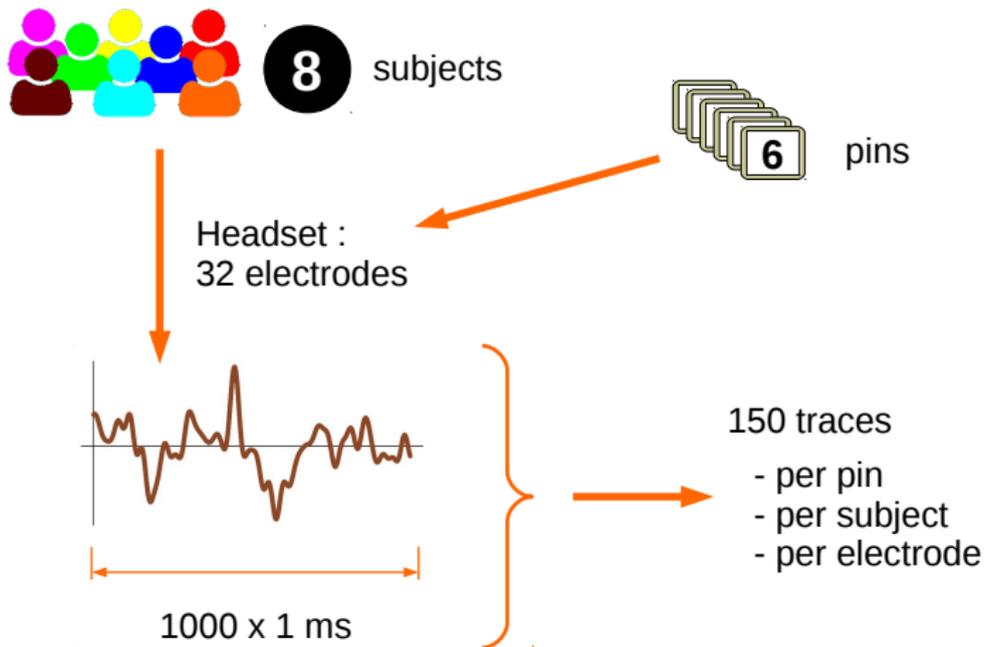
Related works: semantic associations and incongruities

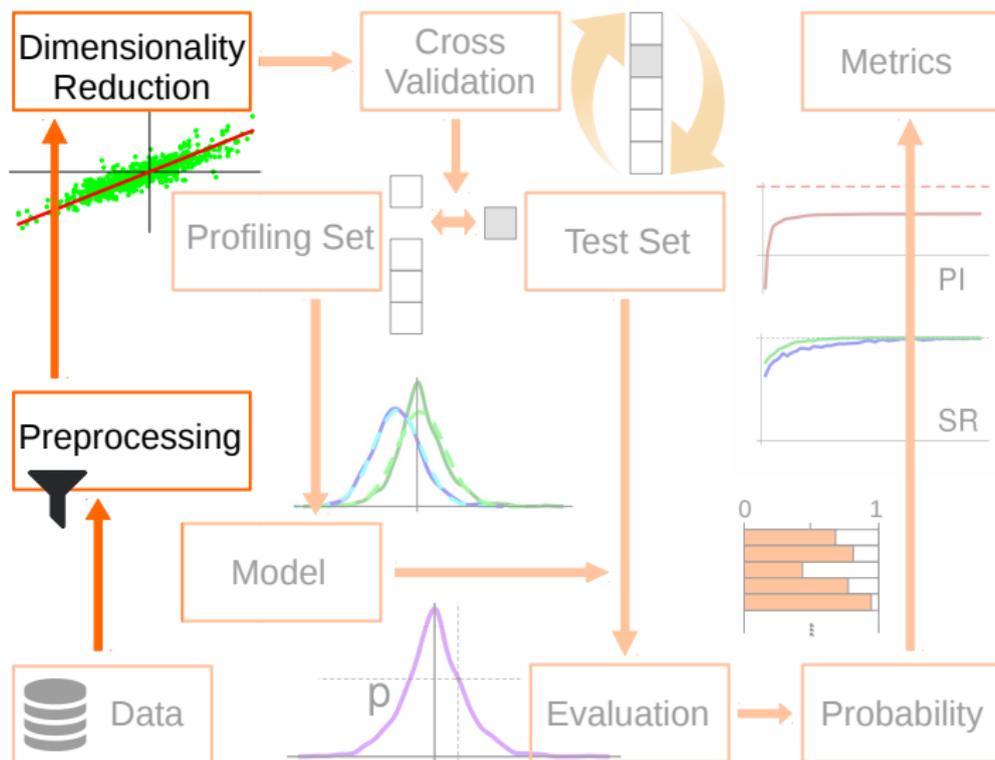






The data

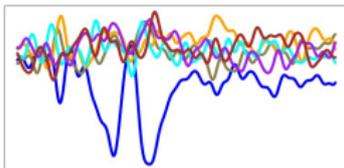




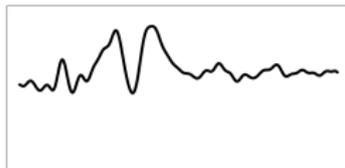
Dimensionality reduction

Principal Component Analysis (PCA)

Mean traces

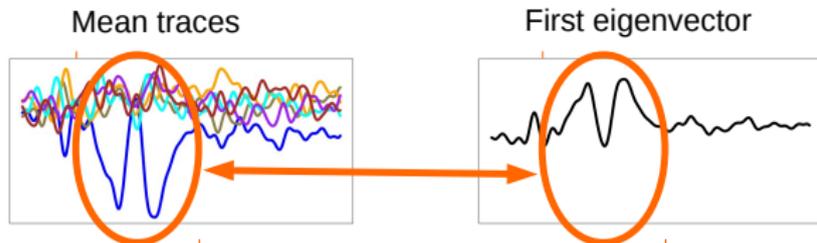


First eigenvector



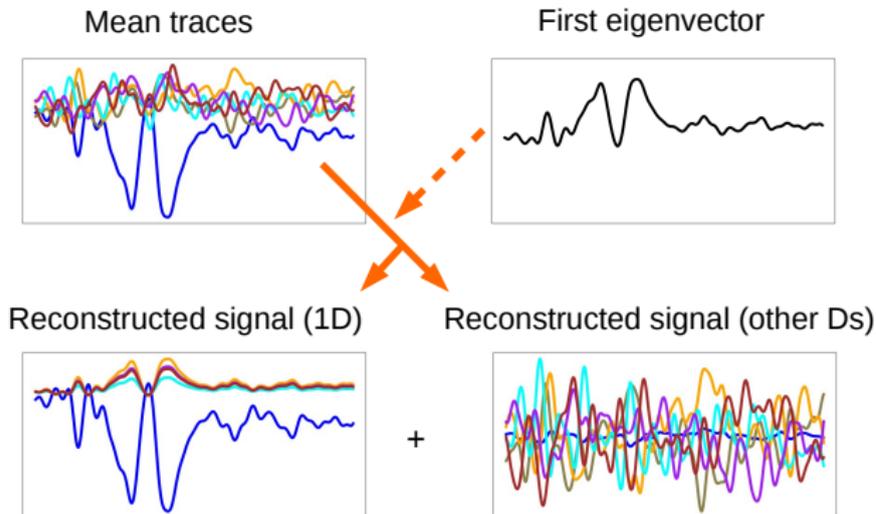
Dimensionality reduction

Principal Component Analysis (PCA)



Dimensionality reduction

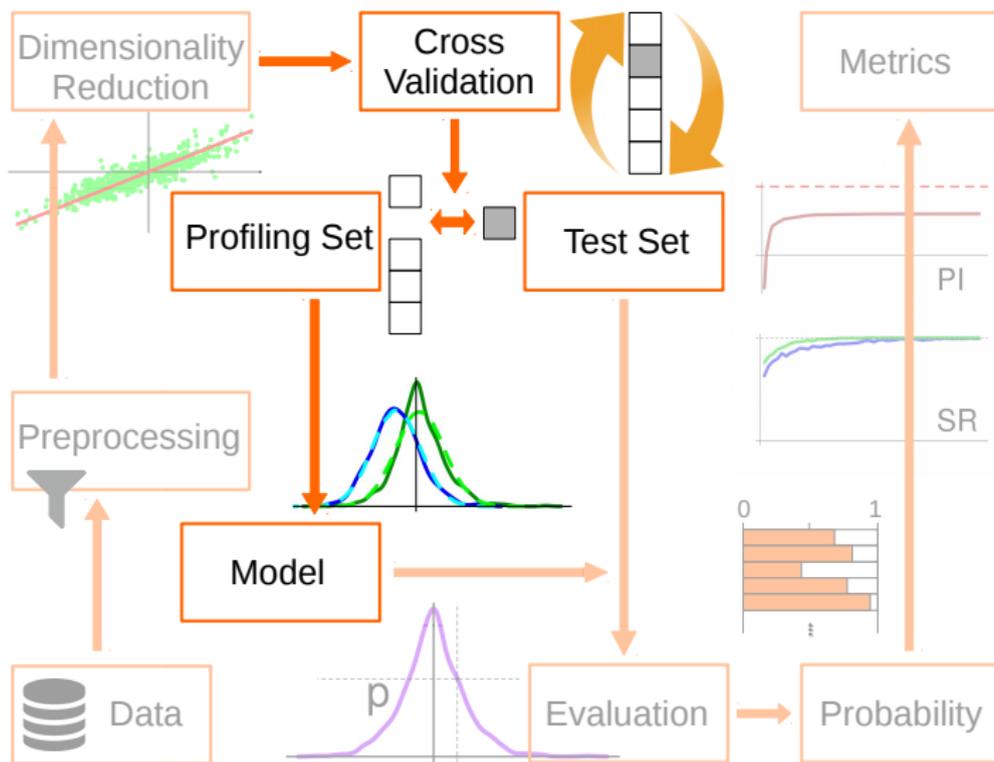
Principal Component Analysis (PCA)



Principal Component Analysis

- + Average PCA: one dimension is sufficient
 - + Projects 1000 dimensions to a single one
 - Estimating means becomes expensive with many PINs
- ⇒ Raw PCA also studied in the paper
(Requires more dimensions and outliers management)





Profiling/modeling

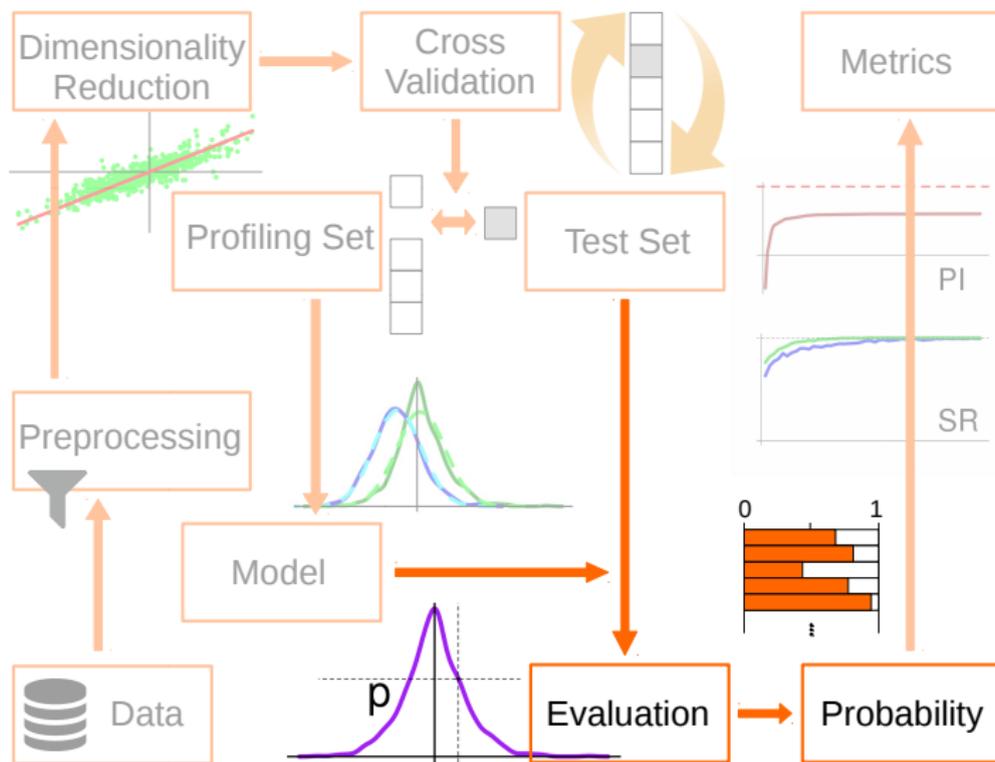
- ▶ **Gaussian** estimation with mean $\hat{\mu}$ and variance $\hat{\sigma}$:

$$\hat{f}_g(x) = \frac{1}{\sqrt{2\pi\hat{\sigma}}} \exp\left(-\frac{1}{2} \left(\frac{x - \hat{\mu}}{\hat{\sigma}}\right)^2\right)$$

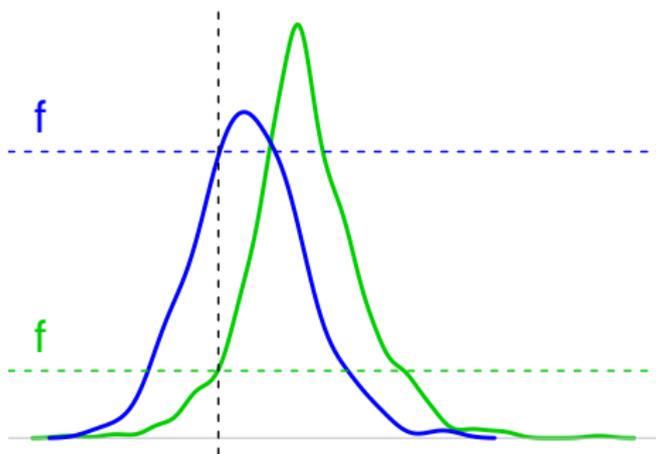
- ▶ **Kernel** density estimation with bandwidth parameter h and samples x_1, \dots, x_n :

$$\hat{f}_k(x) = \frac{1}{nh} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x - x_i}{h}\right)^2\right)$$





Model evaluation



Probability generation

- ▶ From the estimated PDFs

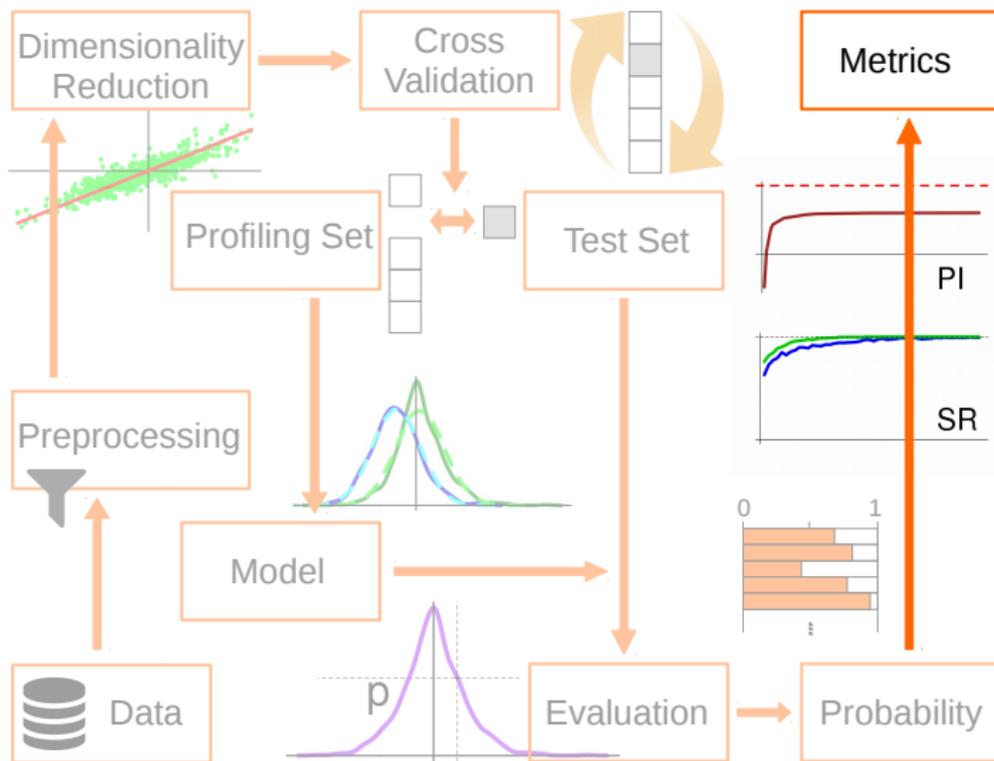
$$\hat{\mathbf{f}}_{model}[\mathbf{o} \mid \mathbf{p} = \text{correct PIN}] = f,$$

$$\hat{\mathbf{f}}_{model}[\mathbf{o} \mid \mathbf{p} = \text{incorrect PIN}] = f.$$

- ▶ Produce probabilities thanks to Bayes

$$\hat{\mathbf{Pr}}_{model}[\mathbf{p} \mid \mathbf{o}] = \frac{\hat{\mathbf{f}}_{model}[\mathbf{o} \mid \mathbf{p}] \cdot \mathbf{Pr}[\mathbf{p}]}{\sum_{\mathbf{p}^*} \hat{\mathbf{f}}_{model}[\mathbf{o} \mid \mathbf{p}^*] \cdot \mathbf{Pr}[\mathbf{p}^*]}$$



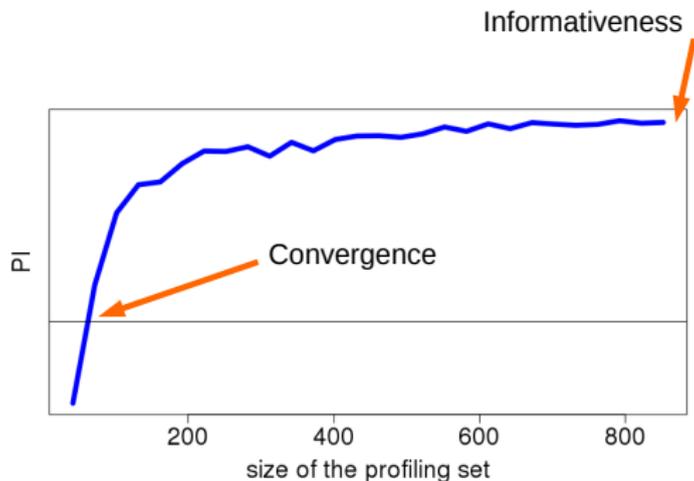


Metrics

- ▶ **Perceived Information (PI):** amount of information extracted from the observations (given a model)
- ▶ **Success Rate:** probability of correct classification (estimated for correct and incorrect PIN values)
- ▶ **Average rank:** average position of the correct PIN value in the sorted list of 6 possible ones



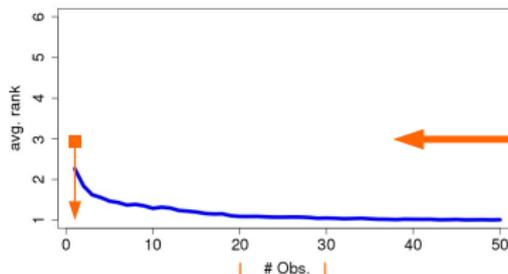
Profiled attacks: PI



- ▶ Convergence reached after \approx 200 to 400 traces



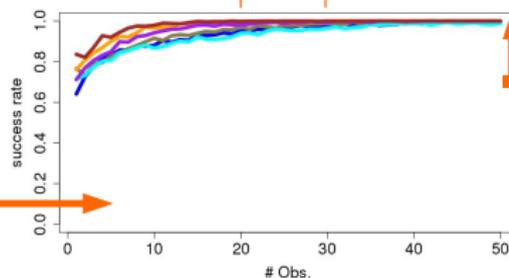
Profiled attacks: SR & avg. rank



Rank 3 to 1

Need about 20-30 traces

Success Rate from
65-85% to 100%



Profiled attacks: summary

- ▶ PIN recovery for most subjects (7 out of 8)
 - ▶ Failure due to another “distinguishable” event
 - ▶ Seems inherent to the investigated setup

⇒ We minimized false negatives (to allow enumeration)

- ▶ Answers our first questions: PINs can be extracted from EEG signals partially and with good confidence
 - ▶ \exists scenarii where this can be damaging
- ▶ Profiling more expensive than online attack
 - ▶ Given a good model is available

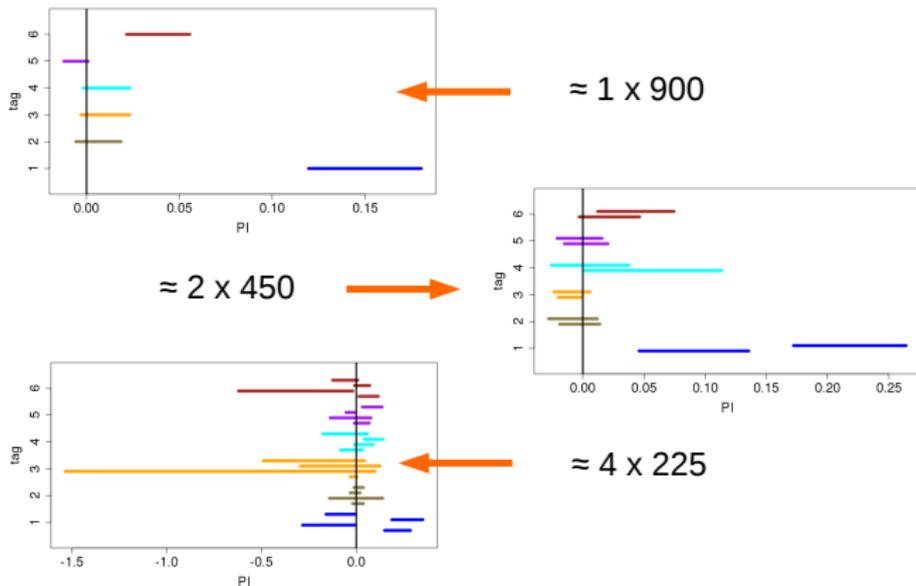


Unprofiled attacks

- ▶ Consider each PIN to be the correct one
- ▶ Estimate the PI on-the-fly for each case
 - ▶ And compute confidence intervals
- ▶ Correct PIN is expected to have the highest PI



Unprofiled attacks: results



- ▶ Nicely correlated with the profiling cost (slide 21)



Other results

- ▶ **Portability:** attack one subject with a model built from others' data: less successful (5 out of 8)
- ▶ **Privacy:** target the subjects' identities instead of their PIN: positive results obtained for all users



Conclusions

- ▶ Information available and exploitable with confidence
- ▶ Yet not sufficient for full (4-digit) PIN recovery
- ▶ Mostly because of signal instability / subjects' focus
- ▶ Biggest risk here: reduction of the guessing entropy



Conclusions

- ▶ Information available and exploitable with confidence
- ▶ Yet not sufficient for full (4-digit) PIN recovery
- ▶ Mostly because of signal instability / subjects' focus
- ▶ Biggest risk here: reduction of the guessing entropy

More generally...

- ▶ Targets of smaller cardinality would be more worrying
- ▶ Privacy is also more worrying (unbounded data)
- ▶ Motivation for MPC, FHE, ...
- ▶ Much more research needed



Thanks !

