



I Have Got to Warn You, It Is a Learning Robot: Using Deep Learning Attribution Methods for Fault Injection Attacks

Karim M. Abdellatif





LEDGER

Hardware Wallet Manufacturer



Donjon
Ledger's Security Research Team



- **Fault injection:** Perturbing the chip during sensitive operations:
 - Power and clock glitches
 - Electromagnetic fault injection (EMFI)
 - Body biasing injection (BBI)
 - Laser fault injection (LFI)
- **Side-channel:** Investing leakages such as EM, power, or time to perform:
 - Simple power analysis (SPA)
 - Differential power analysis (DPA)
 - Profiling attacks

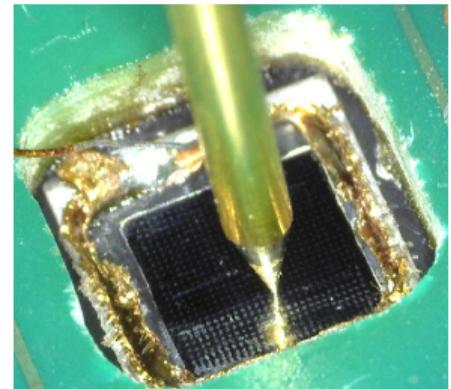




Motivation

- Working on black-box fault injection evaluations takes a lot of time.
- A lot of parameters should brute-forced:
 - Example: BBI or laser fault injection require tuning the following parameters: pulse power, pulse width, **vulnerable timing moments**, and XY point.
- Identifying vulnerable timing moments is one of the big challenges, especially under the case of countermeasures that require injecting multiple faults.

Having reverse engineering tools would be very useful in such evaluations.



(BBI attack ¹)

¹Donjon, "Breaking A Recent SoC's Hardware AES Accelerator Using Body Biasing Injection", HW.io 2022.



Outline

Deep Learning in Hardware Security

Deep Learning Attribution Methods

Practical Challenge: DS28C36 from Analog Devices

Applying DL Attribution Methods into Fault Injection

Tooling

Conclusion

DEEP LEARNING IN HARDWARE SECURITY



- DL-based SCAs ²
 - Several devices for learning and test
 - Better efficiency in case of countermeasures ³
- DL-based leakage detection ⁴
 - It uses DL attribution methods to detect POIs.
 - Better than classical statistical techniques in case of countermeasures

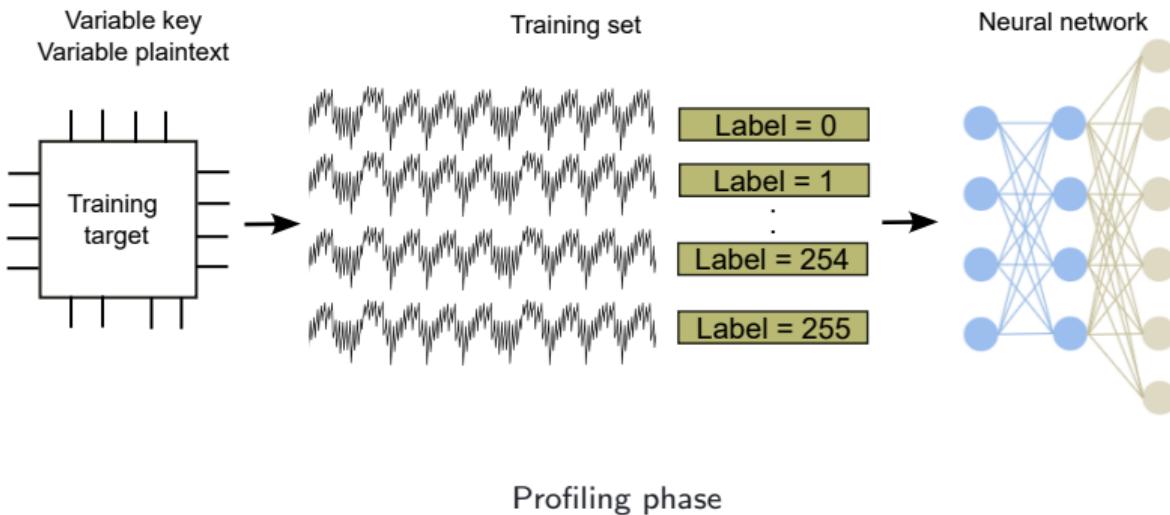
²H. Maghrebi, T. Portigliatti, and E. Prouff. "Breaking cryptographic implementations using deep learning techniques", SPACE 2016.

³E. Cagli, C. Dumas, and E. Prouff "Convolutional neural networks with data augmentation against jitter-based countermeasures: Profiling attacks without pre-processing", CHES 2017

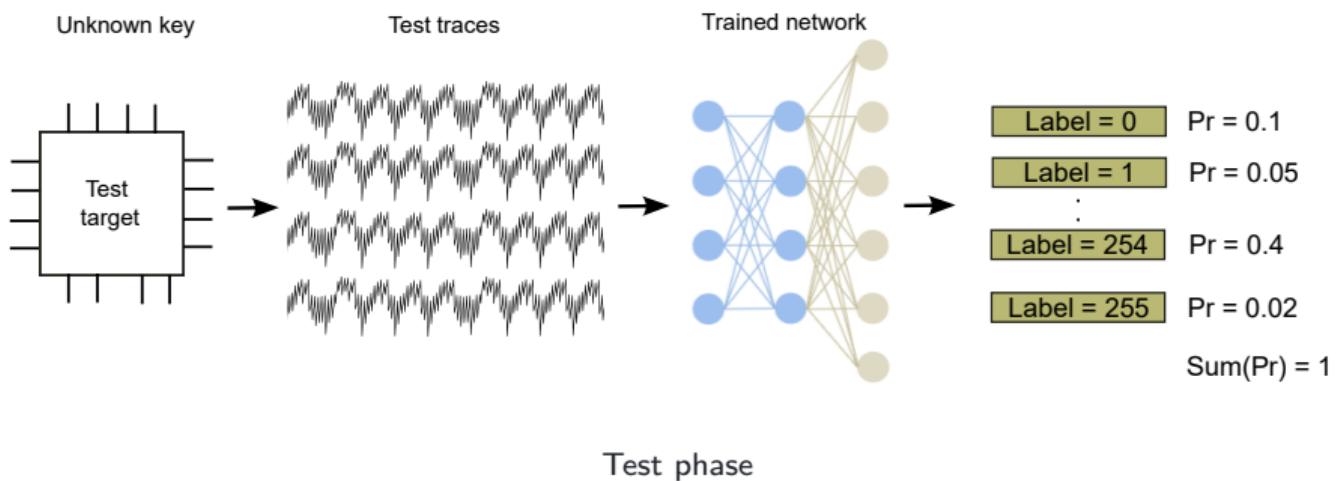
⁴L. Masure et al, Gradient Visualization for General Characterization in Profiling Attacks, IACR.



DL-based SCAs



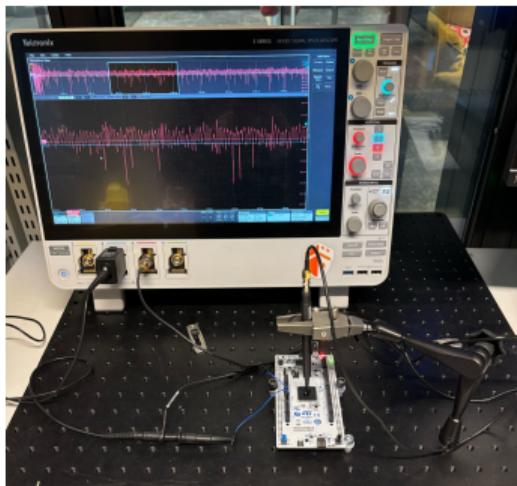
DL-based SCAs



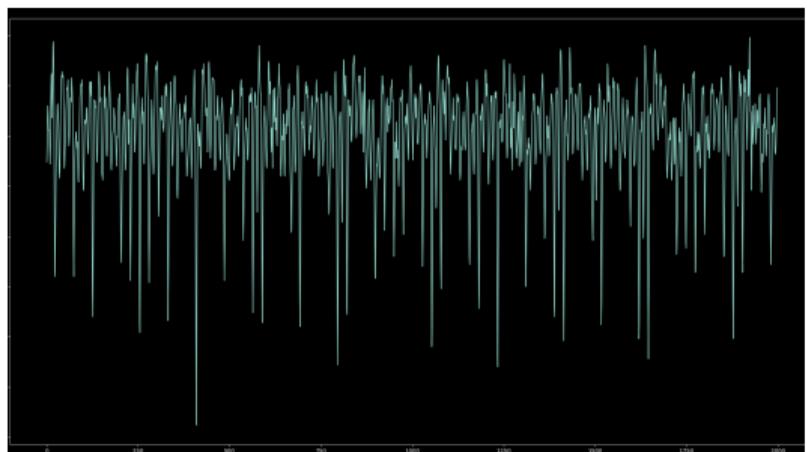


Practical example

Running AES-128 (first round) on a 32-bit MCU.



EM setup for STM32U5 - Donjon



EM signal



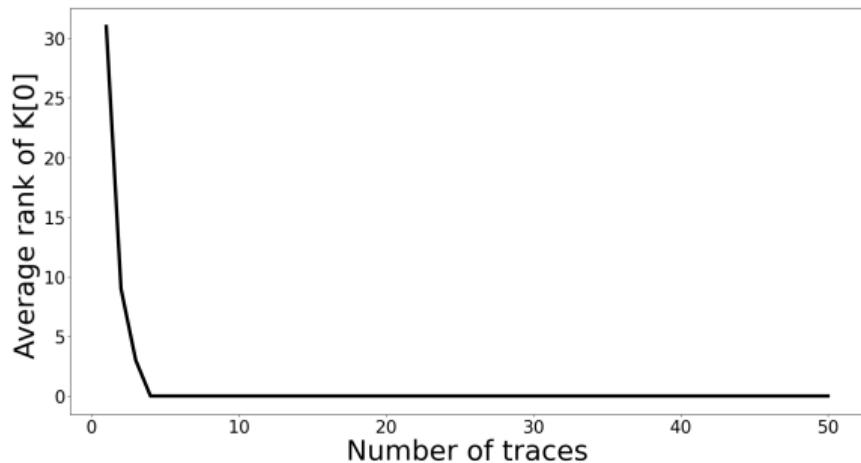
MLP-based example

```
1 def mlp_model(sample_len, range_outer_layer):
2     model = Sequential()
3     model.add(Dense(20, input_dim=sample_len, activation=tf.nn.relu))
4     model.add(Dense(10, activation=tf.nn.relu))
5     model.add(Dense(range_outer_layer, activation=tf.nn.softmax))
6     model.compile(
7         optimizer="adam",
8         loss="categorical_crossentropy",
9         metrics=["accuracy"],
10    )
11    return model
```



Few traces to attack unknown key

- 500K traces for profiling
- 1K traces for test
- Labels on Sbox output



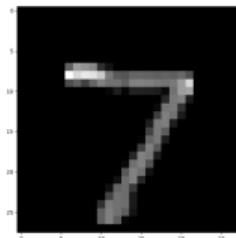
DEEP LEARNING ATTRIBUTION METHODS



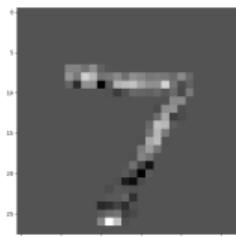
Attribution methods

- Such methods are used to interpret and understand the decisions made by deep neural networks.
- Identify which input features (e.g., pixels in an image) are most influential in the model's predictions.

Gradient-Based Methods, Activation-Based Methods, ...



Input challenge (predict 7)

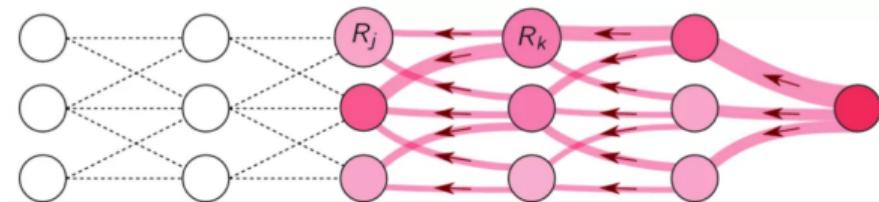


Attribution result



Activation-Based Methods: Layer-wise Relevance Propagation (LRP)

It lies in tracing back the contributions of input nodes to the final prediction.



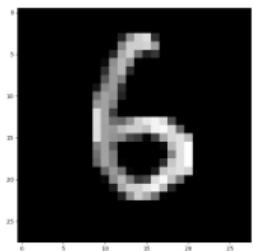
$$R_j = \sum_k \frac{z_{jk}}{\sum_j z_{jk}} R_k \quad (1)$$

Illustration of the LRP procedure ⁵

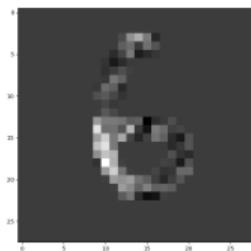
where j and k denote neurons in consecutive layers, and $z_{jk} = a_j w_{jk}$ is the activation of the neuron j multiplied by the weight between neuron j and neuron k .

⁵S. Bach et al. 'On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation.' PloS one 10.7 (2015).

MNIST example



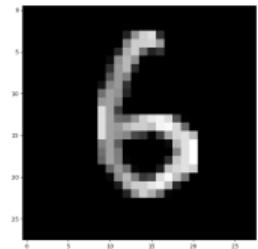
Input challenge



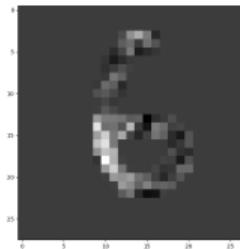
LRP



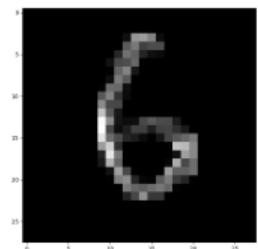
Other methods



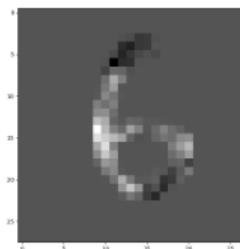
Input challenge



LRP



Taylor



Input



Application into side-channel

Previous work⁶ exists in the side-channel domain.

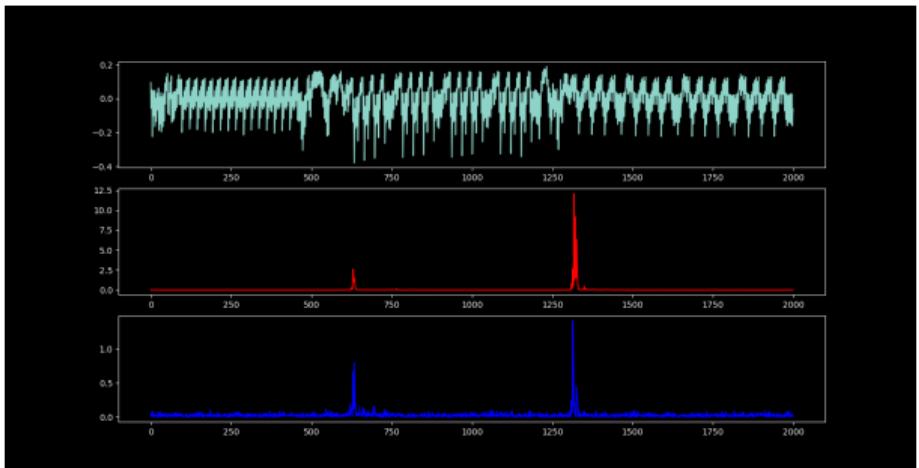
Advantages: detecting leakage points in case of countermeasures (ex: masking and jitter) unlike SNR or T-test.

⁶B. Hettwer et al., " neural network attribution methods for leakage analysis and symmetric key recovery", SAC-2019.



Practical example - AES

- ① Collecting power traces from an AES
- ② Profiling on $S\text{box}[0]$
- ③ Reverse-engineering using attribution methods
- ④ Comparing it with SNR/NICV

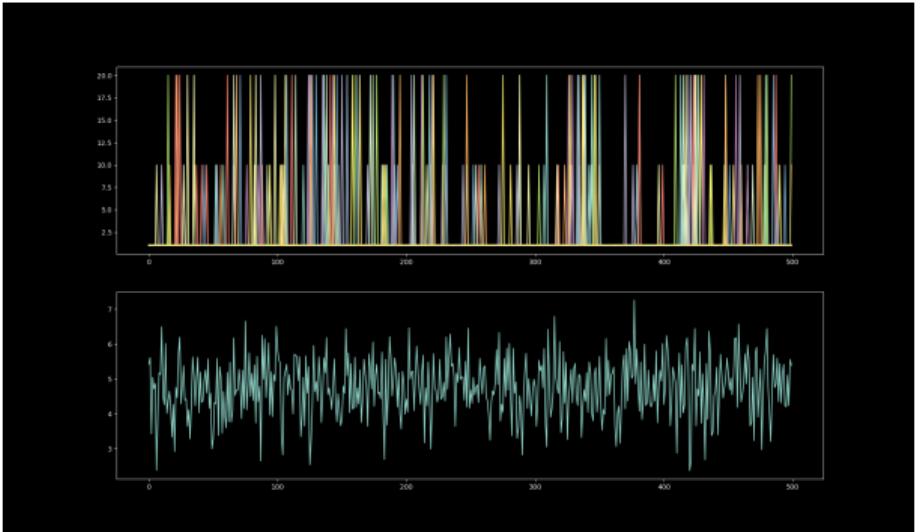


Difference between **SNR** and **LRP**



Another example

- ① Two different values with jitter
- ② T-test fails!

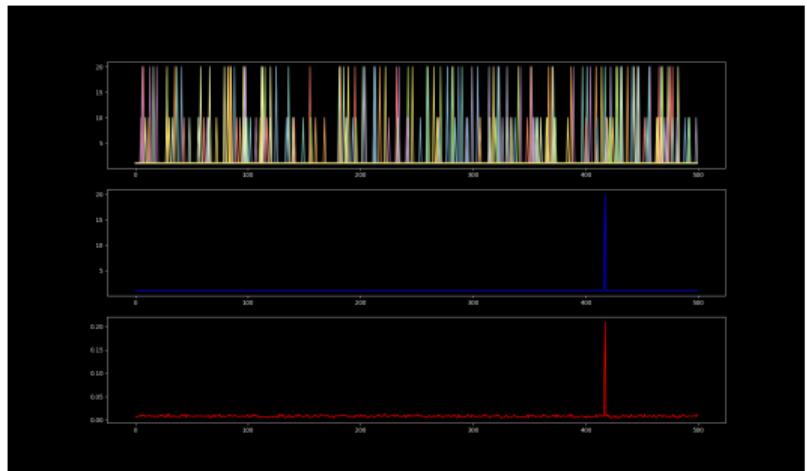


Upper: desynchronized traces, below: T-test



Using attribution methods

- ① Two different values are randomized.
- ② Profiling on the two labels of them
- ③ Timing is very well detected.
- ④ Advantages:
 - Decision scalability (one trace)
 - It can defeat countermeasures.



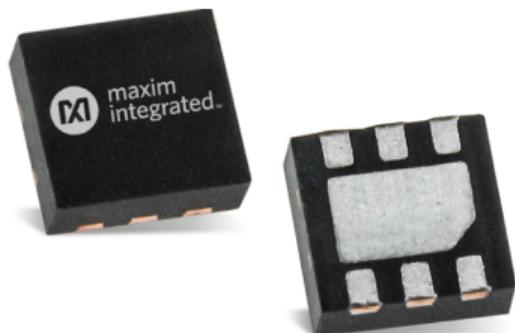
Sample trace and LRP

PRACTICAL CHALLENGE: DS28C36 FROM ANALOG DEVICES

Security features ⁷



- ECC-256 computation engine
- FIPS 180 SHA-256 computation engine
- TRNG with NIST SP 800-90B compliant entropy source with function to read out
- 17-Bit one-time settable, non-volatile decrement-only counter with authenticated read
- **8Kbit of EEPROM for user data, keys, and certificates**
- The full data sheet is not available and this required some reverse to find the available commands and their parameters.



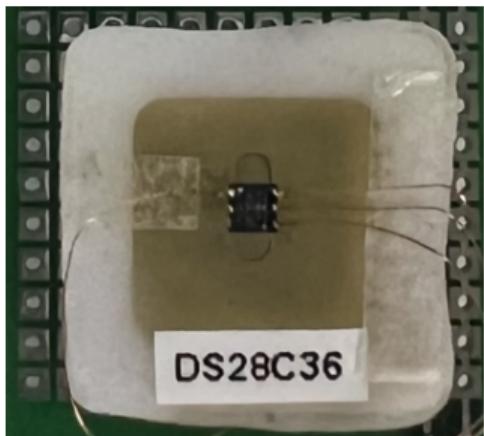
⁷<https://www.analog.com/media/en/technical-documentation/data-sheets/DS28C36.pdf>



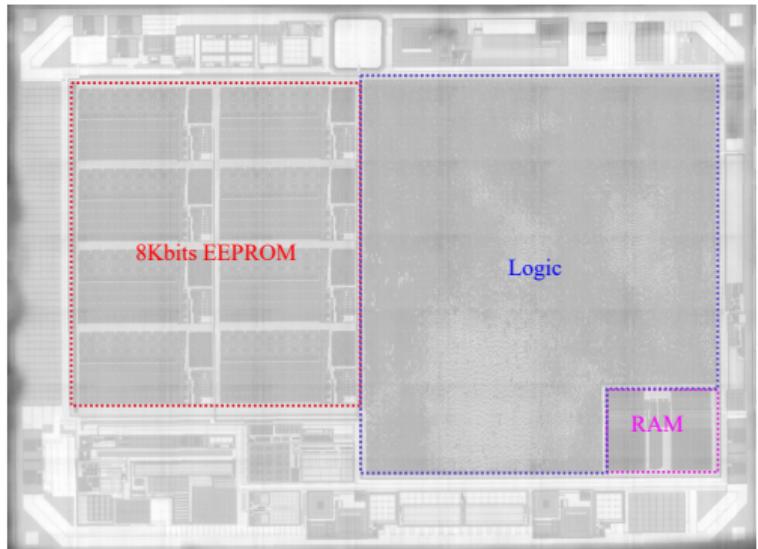
Page	Description
0 to 15	User pages
16 to 21	Public keys (x and y)
22 to 24	Private keys
25 to 26	Secret pages
27	Counter
28 to 29	Random
30 to 31	RAM buffer



Sample preparation



Decapped chip

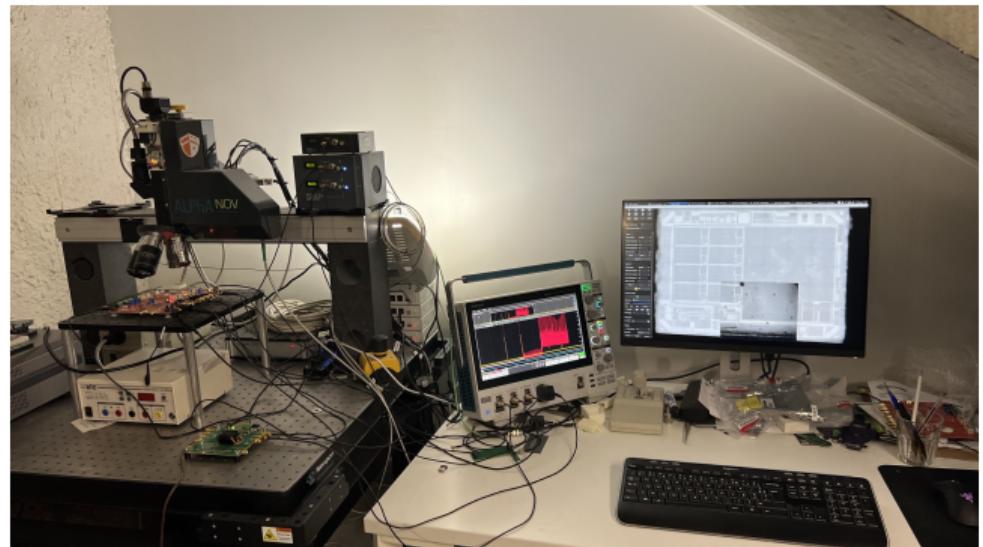


Infrared backside image



Setup

- An infrared pulsed laser source and a microscope for focusing
- A Scaffold⁸ board
- A Tektronix MSO44 oscilloscope
- DUT: DS28C36



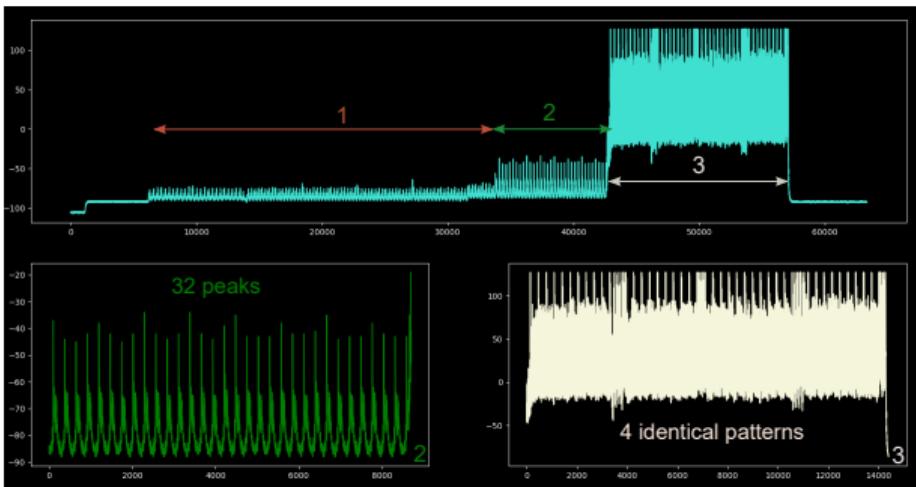
Setup

⁸O. Heriveaux. Scaffold. <https://github.com/Ledger-Donjon/scaffold>



Read page command

```
1 write_data(page_number, data)
2 read_page(page_number)
3 save_power_trace()
```

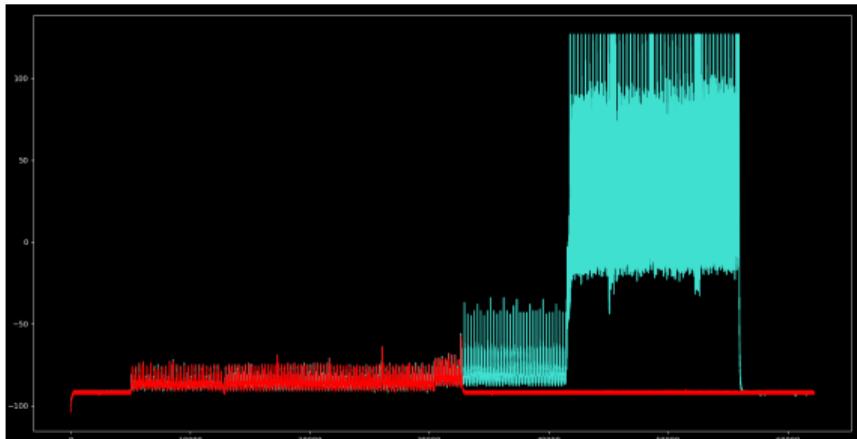


Power consumption in case of unprotected page



Unprotected and protected page

```
1 write_data(page_number, data)
2 read_page(page_number)
3 save_power_trace()
4 lock_page(page_number)
5 read_page(page_number)
6 save_power_trace()
```



Protected and unprotected



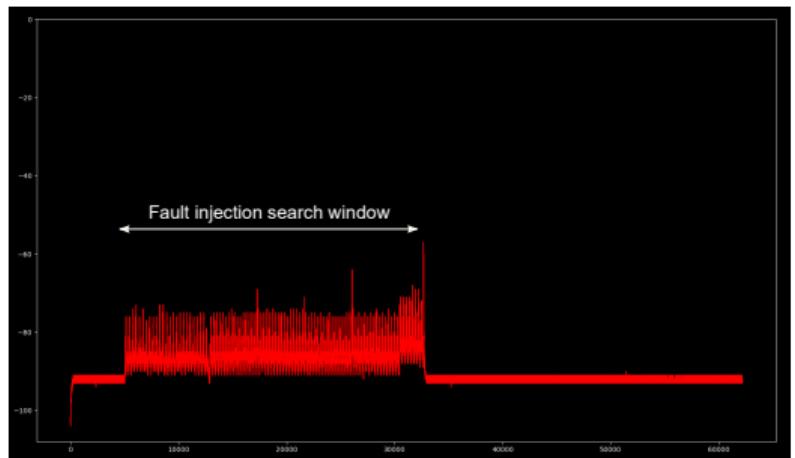
Attack scenario

Step 1:

```
1 write_data(page_number, data)
2 lock_page(page_number)
```

Step 2:

```
1 while True:
2     prepare_fault() # single pulse
3     chip_restart()
4     read_page(page_number)
5     save_log()
6     move_laser()
```



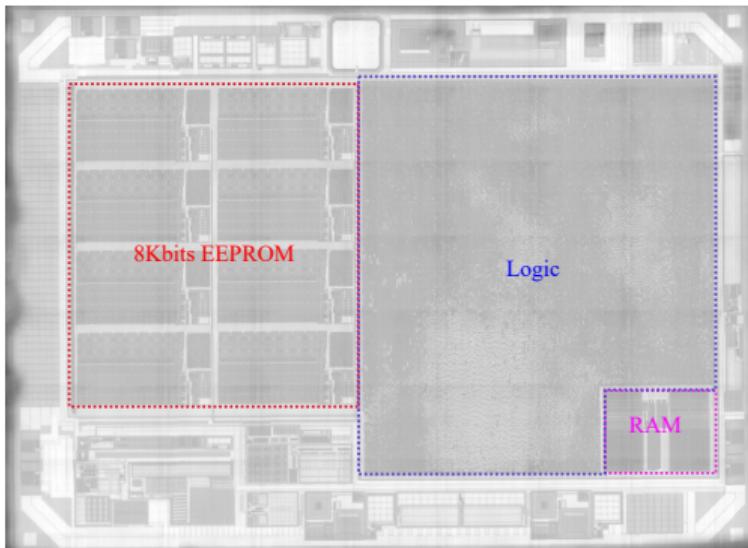
Power consumption when the page is locked

Investigation



Page configuration (bit or bits) can be:

- Stored in the EEPROM
- Stored in eFuses
- Manipulated in the logic
- Temporarily stored in RAM



IR image

Results

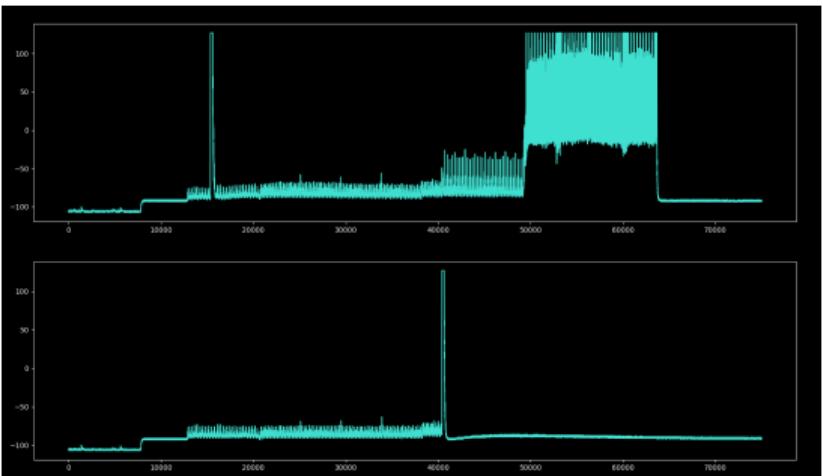


Number	Chip response	Note
0	2155fffffffffffff fffff	Locked
1	fffff fffff	Timeout
2	NACK	I2C communication error
3	21aab8289516978a7b25eb1d8a317f6c6a 71718b4d47de4754ac32a1d1c5adb7d324	Public key slot
4	21aa208cf9a7dc7fcdb5437775fea79aa 2c95f5795ed2bfe883082a2ada0585694f	Needs to be investigated

Investigation



- Correct read page trace for response 3
- In case of response 4, no EEPROM read
- It seems to be a RAM, or RNG content



Difference between response 3 (upper) and response 4 (below)



- The chip seems to be protected against single fault attacks? (black box evaluation)
- **Reverse-engineering the Read Page command is the only way to understand clearly.**
- **How can we do that? deep learning?**

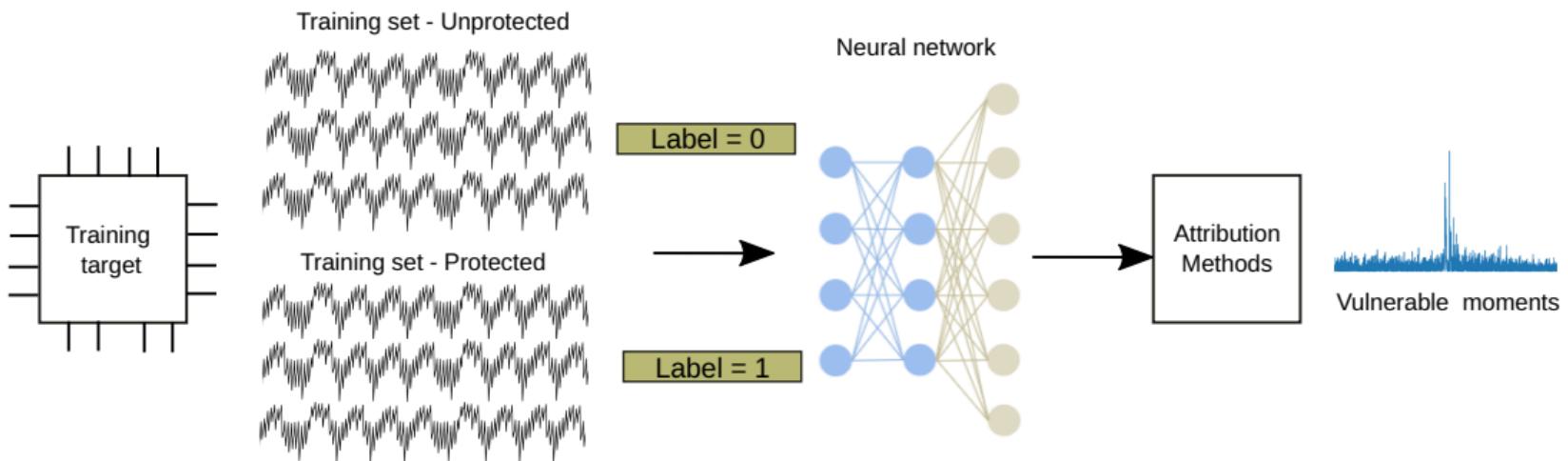
APPLYING DL ATTRIBUTION METHODS INTO FAULT INJECTION



- We will apply the DL attribution methods, which are used in SCAs to detect sensitive operations, into fault injection (FI).
- The main purpose is to detect when sensitive bits are processed.
- More precisely, we will try to locate on the power consumption trace, the manipulation of the page protection bit/bits.
 - The first set is collected when the page is unlocked (50K traces).
 - The second set is collected when the **same page** is locked (50K traces).



Methodology





Methodology

```
80     leakages = np.concatenate((leakages_0, leakages_1), axis=0) → Combining two datasets
81     metadata = np.concatenate((labels_0, labels_1), axis=0)
82     x_train = normalization((leakages), feature_range=(-1, 1)) → Performance improvement
83     GUESS_RANGE = 2 # protected and unprotected → 0 and 1
84     model = model_mlp(x_train.shape[1], GUESS_RANGE)
85     # Train the model
86     profile_engine = Profile(model, leakage_model=leakage_model)
87     EPOCHS = 5
88     profile_engine.data_augmentation(aug_mixup)
89     profile_engine.train(
90         x_train=x_train, → Learning phase
91         metadata=metadata,
92         guess_range=GUESS_RANGE,
93         epochs=EPOCHS,
94         batch_size=10,
95         validation_split=0.1,
96         data_augmentation=False,
97     )
```



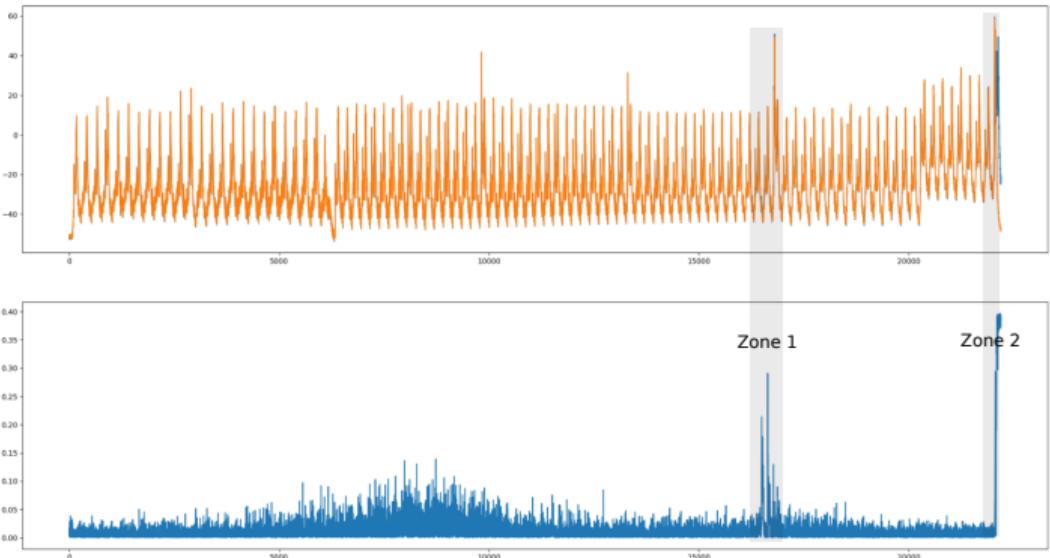
Methodology

```
115     model_wo_sm = innvestigate.model_wo_softmax(model)
116     gradient_analyzer = innvestigate.analyzer.LRP(model_wo_sm) → LRP
117     trace_sample = x_test[0]
118     trace = trace_sample.reshape(1, x_test.shape[1])
119     vis_trace = gradient_analyzer.analyze(trace)[0]
```



Result

- Two zones
- Protected against single fault attacks

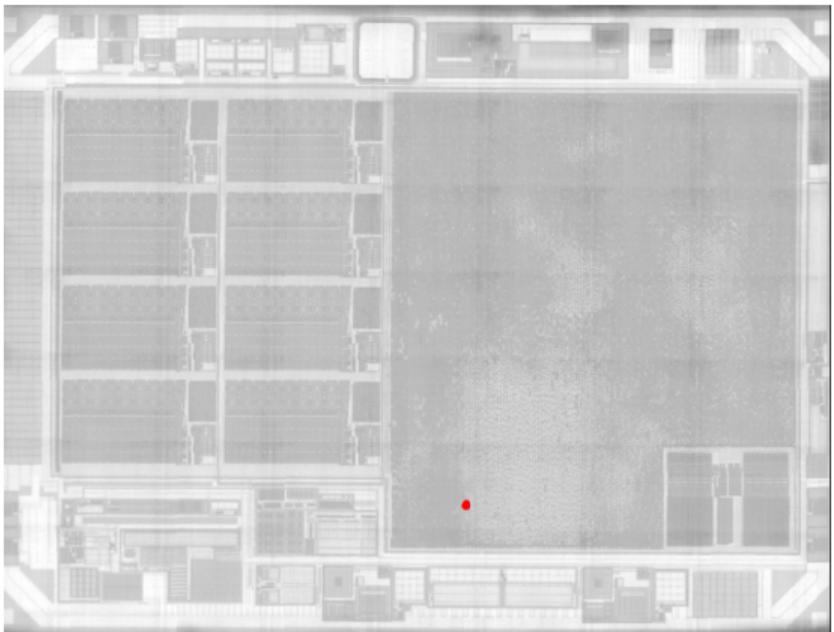


LRP result



Successful faults

- Scanning the chip with double pulses
- Scanning the logic area
- Successful fault

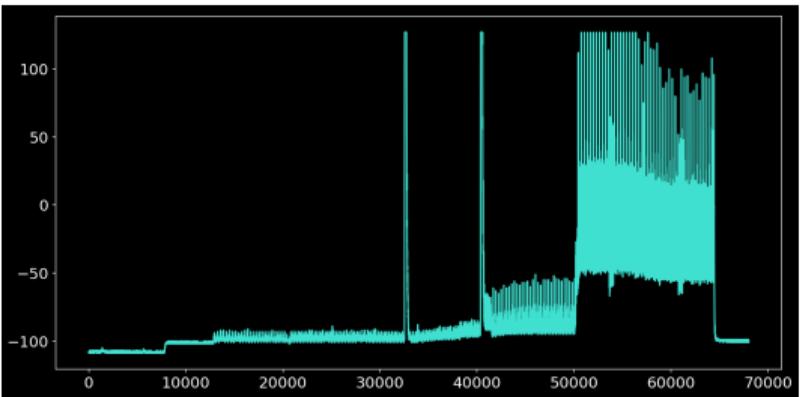


Vulnerable spot



Double fault attack

- Fixing the laser beam on the correct location
- Success rate close to 99%



Power trace in case of a successful fault

This confirms the efficiency of DL attribution methods.



- The presented attack is applicable on all the user pages.
- It isn't applicable on permanent-protected pages used for P256 curve private keys.
 - The chip passed a fixed unidentified value for these pages.

TOOLING



- Latest side-channel attacks using deep learning
- Leakage detection using deep learning attribution methods



Scadl: Open source tool - Donjon

Clone and investigate!

CONCLUSION



- DL attribution methods involved in this work, can be used when performing fault injection attacks in black box context.
- Manufacturers **must** consider such technique for testing countermeasures in addition to leakage detection techniques to detect vulnerable timing moments.
- Using double verification against fault injection attacks is not efficient enough if it is used alone.
- Manufacturers must at minimum combine it with strong hardware and/or software jitter as an additional countermeasure.

THANK YOU. QUESTIONS?



Karim M. Abdellatif, PhD
e-mail: karim.abdellatif@ledger.fr