# Google

# A hacker guide to deep-learning based AES side channel attacks

**Elie Bursztein**
Google, @elie

**Jean-Michel Picod**
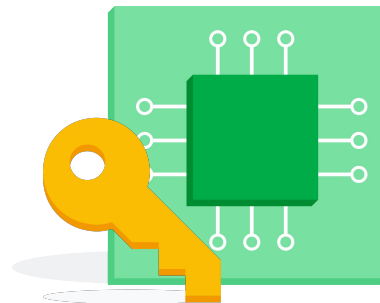Google, @jmichel_p

with the help of **many** Googlers and external collaborators

Security and Privacy Group

Side channel attacks are one of the **most efficient ways to attack secure hardware**

Google

Security and Privacy Group

A side-channel attack was used to recover the Trezor bitcoin wallet private key

https://jochen-hoenicke.de/crypto/trezor-power-analysis/

Security and Privacy Group

Side-channels attacks requires a lot of **domain expertise** and are **implementation specific**
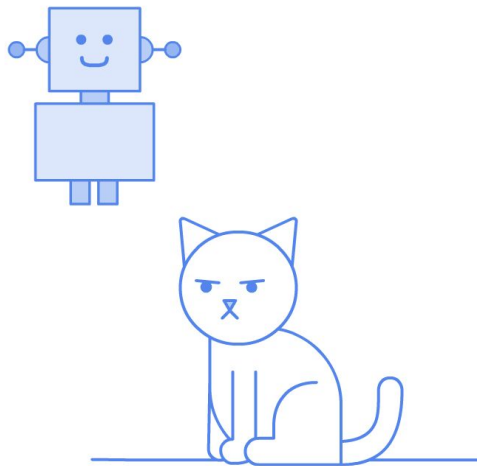
Google

Is there a better and more generic way to perform side-channels attacks?

Google

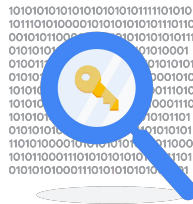Deep-learning is posed to revolutionize hardware side-channel cryptanalysis

Google

Security and Privacy Group
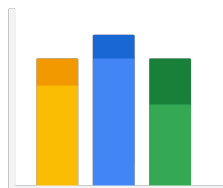
# AI? Really?

Google

Security and Privacy Group

# Template attack on steroids

No trace processing

Direct attack point targeting

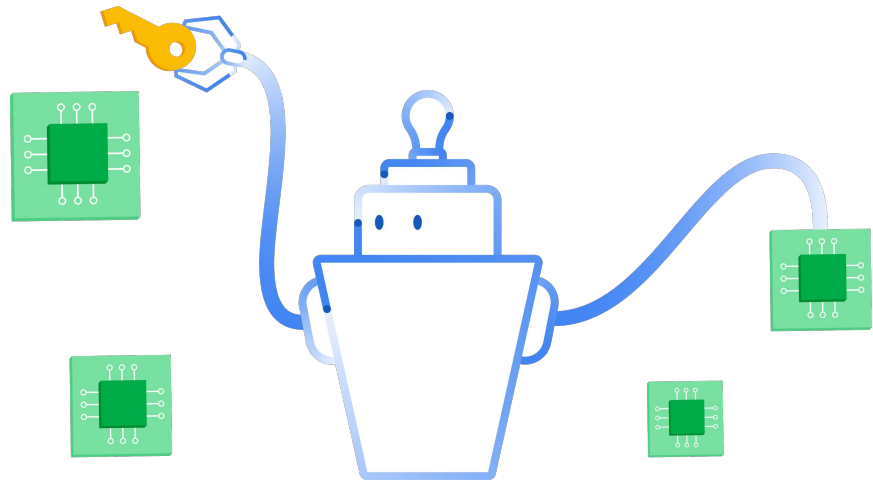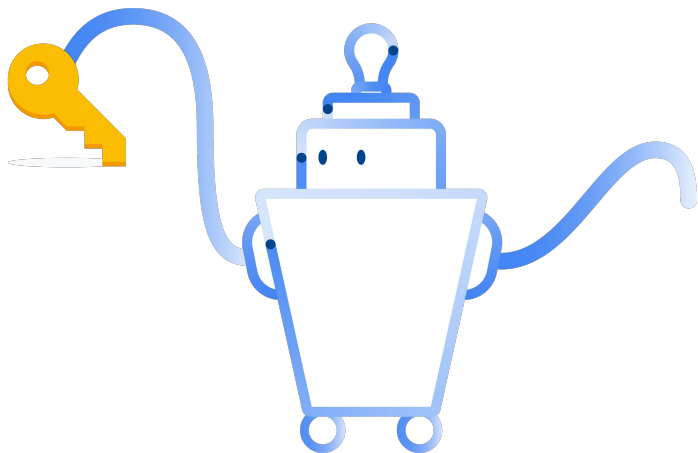Efficient probabilistic attack

Better and intuitive success metrics

Security and Privacy Group

# Attacks are going to be better over time as **deep learning scales with data and computing**



Deep learning

Other AI Algorithms

1980s & 1990s

2018+

More compute & data

Performance

Scale (data, model, computing)

Google

Security and Privacy Group

How to use deep-learning to recover AES keys baked in hardware **in practice**

Google

Security and Privacy Group

Side
Channel
Attacks
Assisted with
Machine
Learning

Google

Security and Privacy Group

Talk is based on some of the results of a joint research project with many collaborators on **hardening hardware cryptography**

Security and Privacy Group

# Code and slides
https://elie.net/scaaml

Google

Security and Privacy Group

# Disclaimer

This talk purposely focuses on showcasing how to a get SCAAML attack working end-to-end rather than discussing state of the art attacks.

Google

Security and Privacy Group

# Agenda

What are side-channels?

What is deep-learning?

Hacker's guide to
AES SCAAML attacks

What's next

Google

Security and Privacy Group

What are
side-channels?

Google

Security and Privacy Group

A side-channel attack is **an indirect measurement** of a **computation result** via an **auxiliary mechanism**
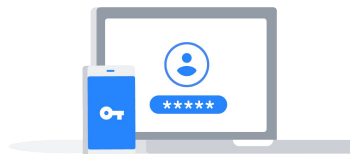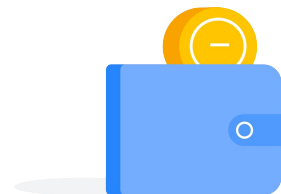
Google

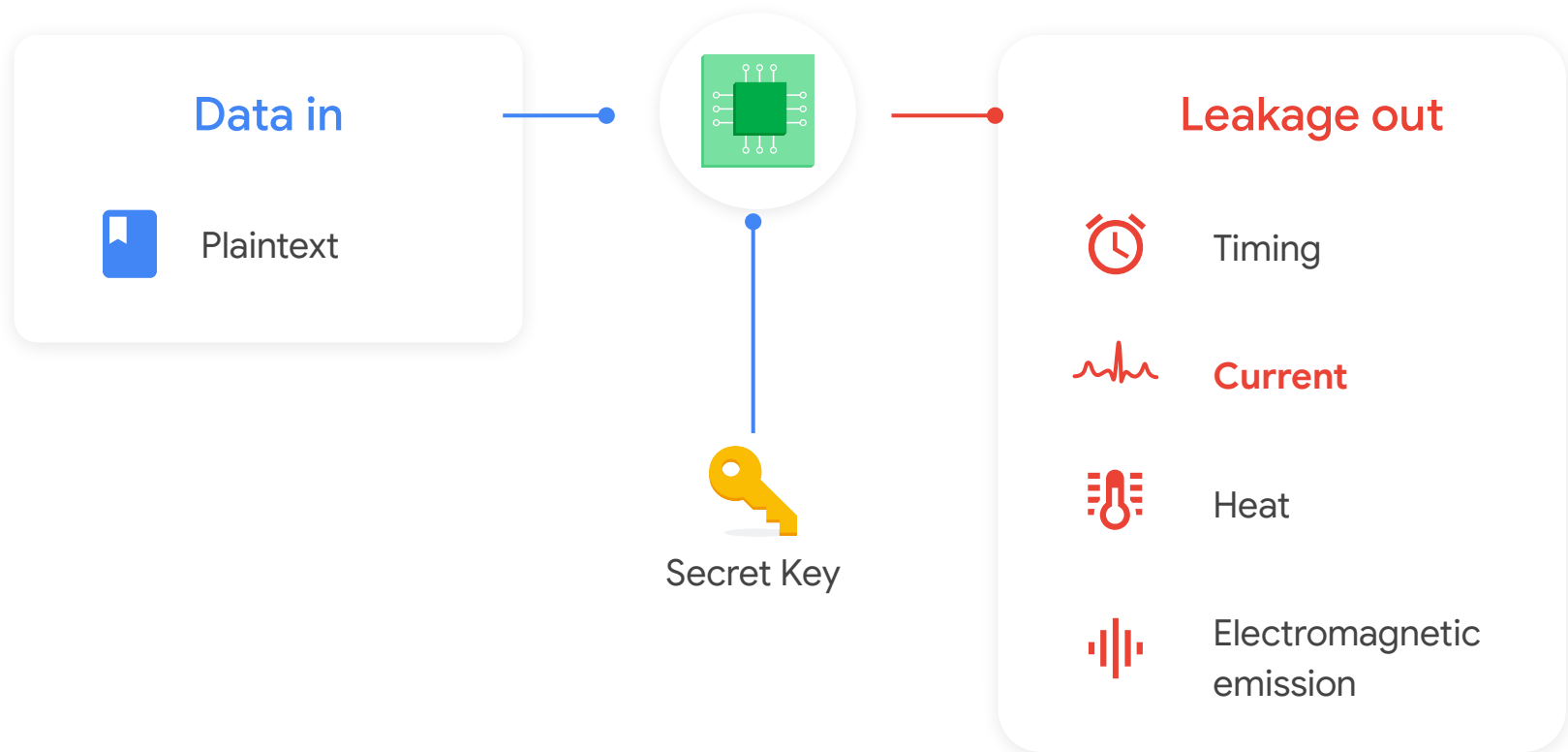Security and Privacy Group

# SCA real-world applications
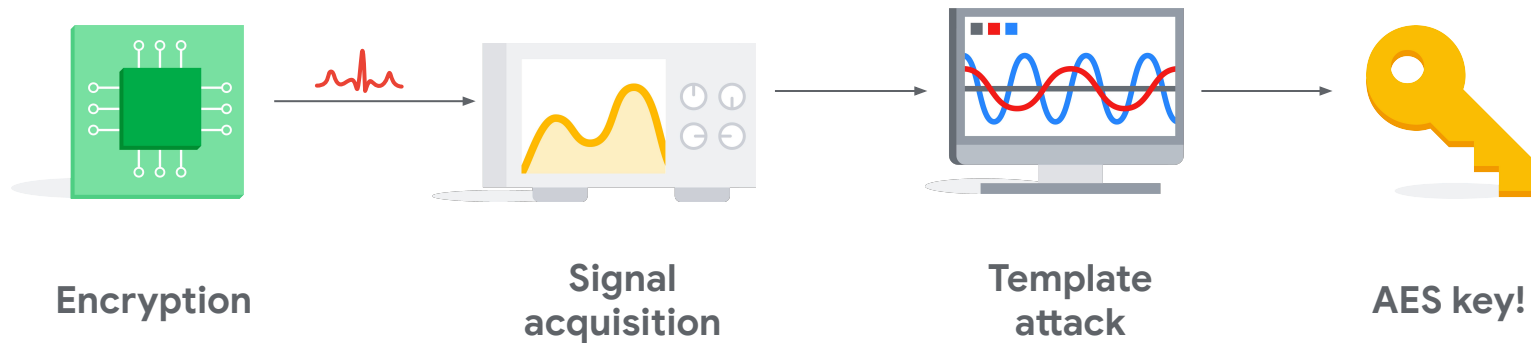
**Recover Encryption keys**

Perform blind SQL injections

Steal passwords & pins

Extract cryptowallet private keys

Google

Security and Privacy Group
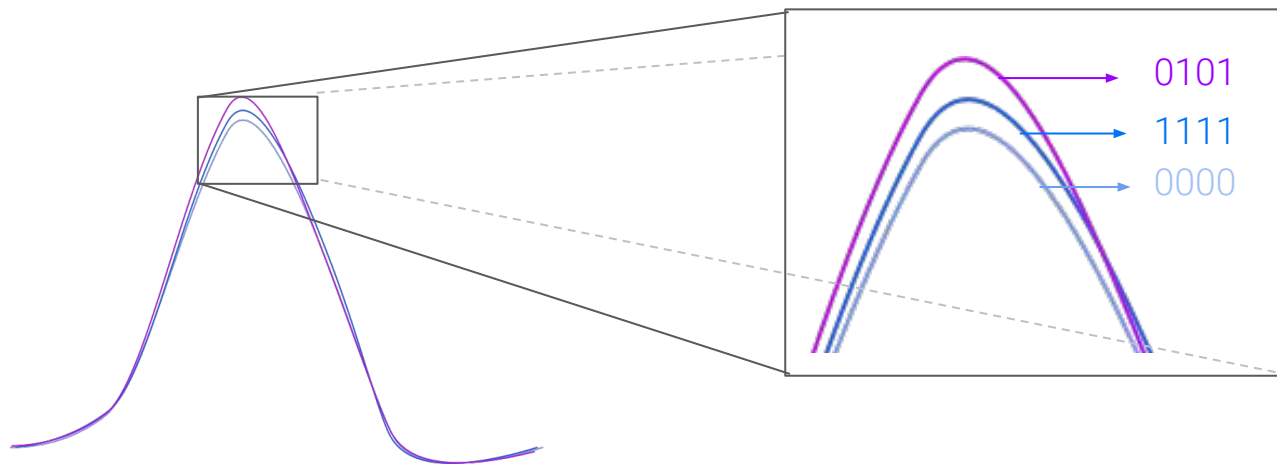
# Data in

📖 Plaintext

Secret Key

# Leakage out

⏰ Timing

〰️ **Current**

🌡️ Heat

▮▮▮ Electromagnetic emission

Google

Security and Privacy Group

# SCA in a nutshell



Encryption → Signal acquisition → Template attack → AES key!

Google

Security and Privacy Group

# Crypto computation side-effects are measurable



0101
1111
0000

Google

Security and Privacy Group

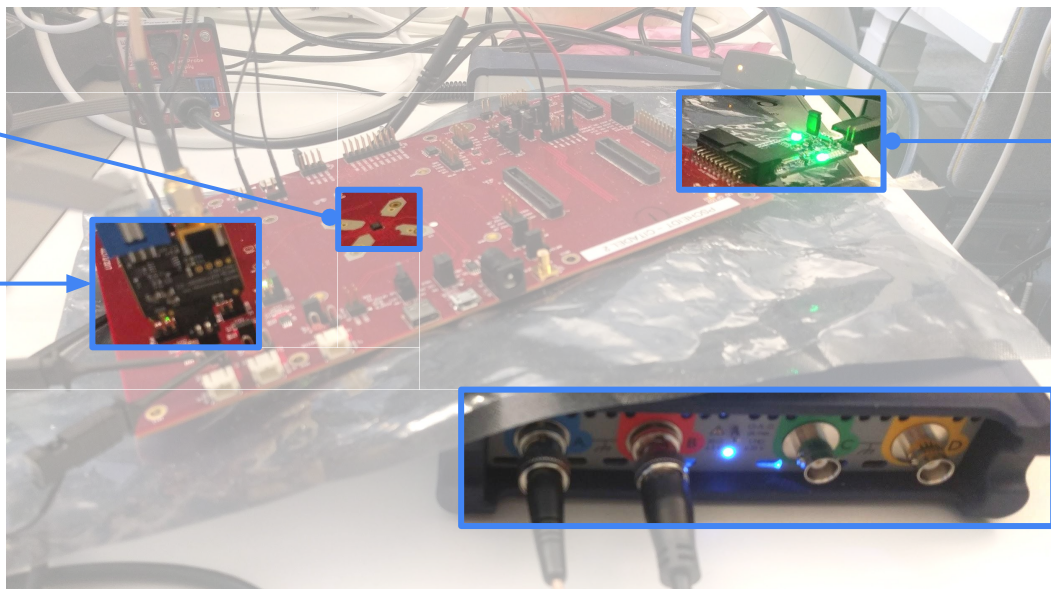# Lightly protected AES power trace

# DIY hardware setup from early days

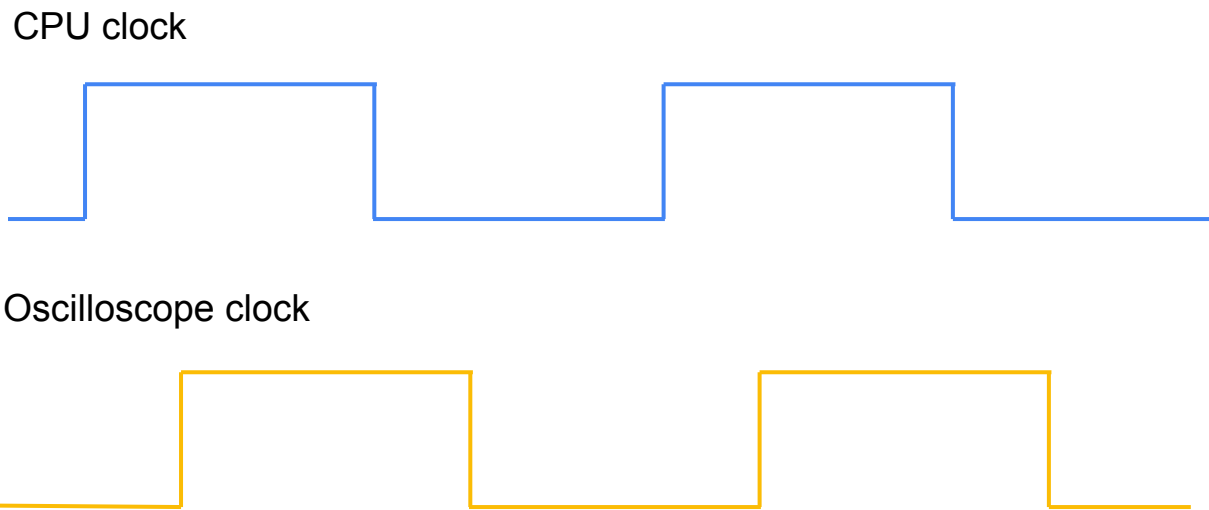

Target chip

Probe

Communication interface

Oscilloscope

Google

Security and Privacy Group

# NewAE Chipwhisperer is an easy and cheap all-in-one entry to side-channel attacks

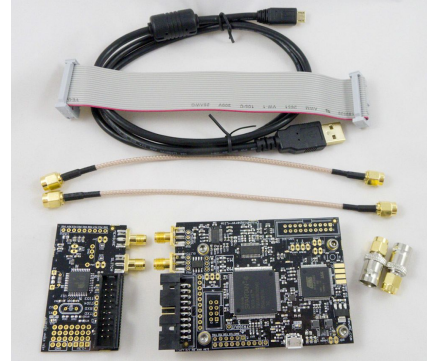https://newae.com/tools/chipwhisperer/



Recommendation based on usage

Google

Security and Privacy Group

For many chips a higher sampling rate is needed due to their clock speed so you need a faster oscilloscope

Security and Privacy Group

**Asynchronous capture** used for blackbox attacks like SCAAML needs **at least 4x the CPU clock speed**

Google

Security and Privacy Group

# NewAE Chipwhisperer Pro + Picoscope 6000 is what we use for our SCA research

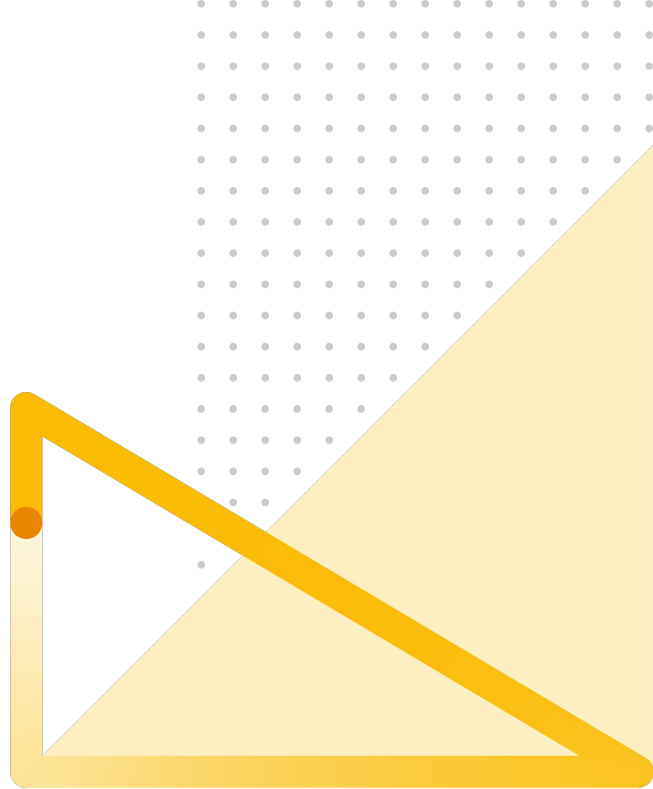This is not an ad :) it is a recommendation
based on what we use

Google

Security and Privacy Group
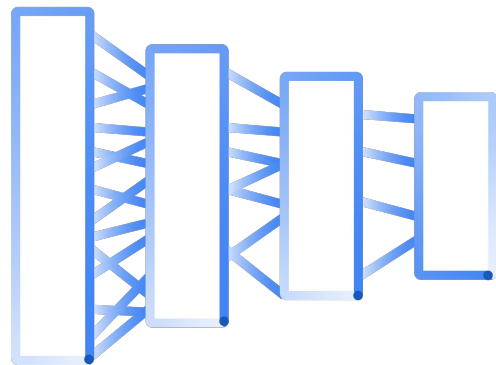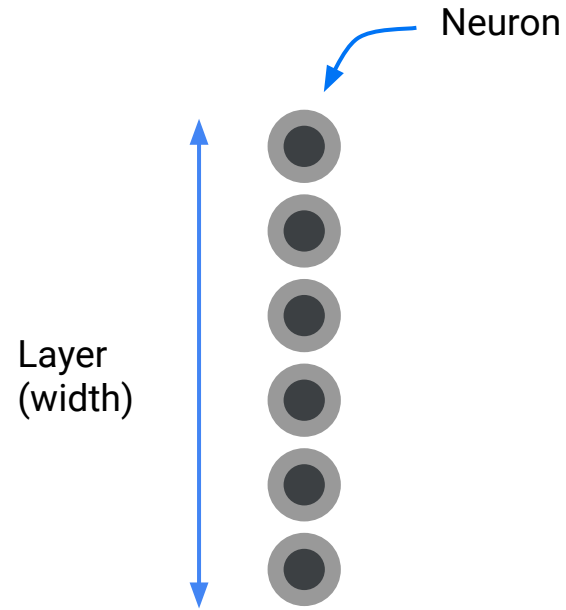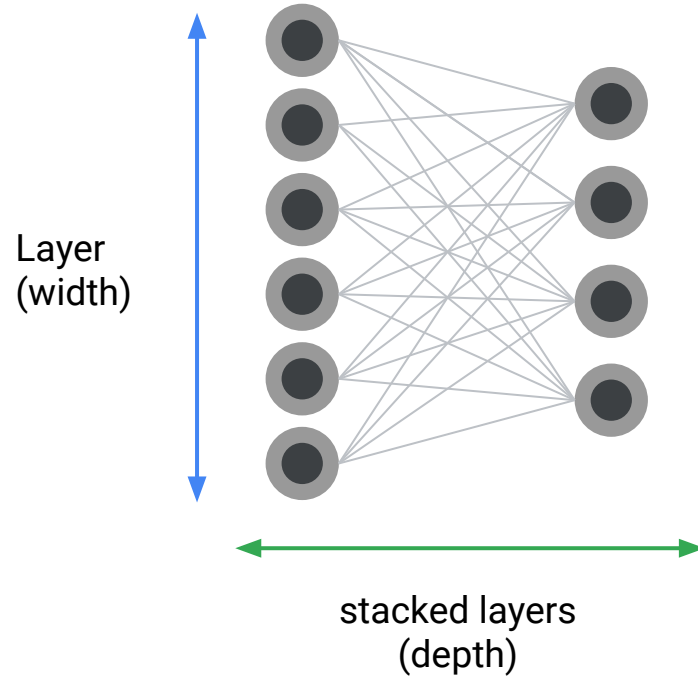
# What is
# deep-learning?

Google

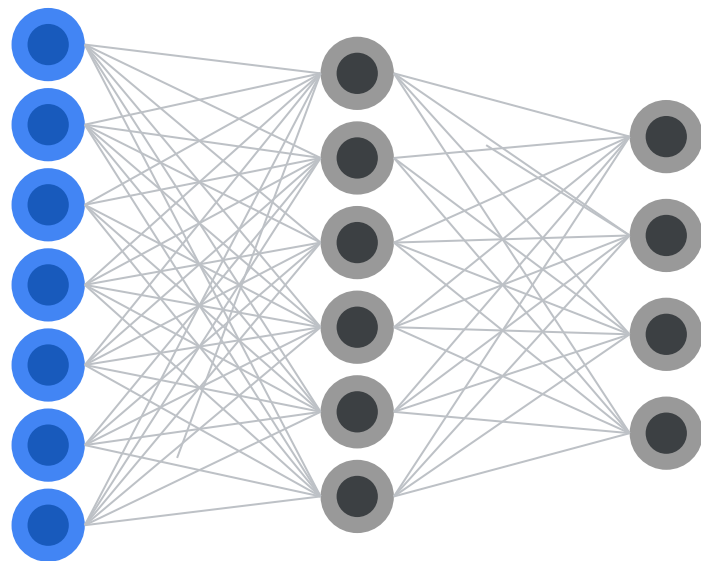At its core deep-learning is basically a neural network with many layers



Google

Security and Privacy Group

Neuron

Layer
(width)

Google

Security and Privacy Group

Layer
(width)

stacked layers
(depth)

Google

Security and Privacy Group

Input layer    hidden layers

Google    Security and Privacy Group

Input layer      hidden layers      Output layer

Dog

Cat

Google

Security and Privacy Group

Input layer          hidden layers          Output layer

Google          Security and Privacy Group

Input layer · hidden layers · Output layer → Dog, Cat

Google · Security and Privacy Group

Input layer        hidden layers        Output layer

Dog

Cat

Google

Security and Privacy Group

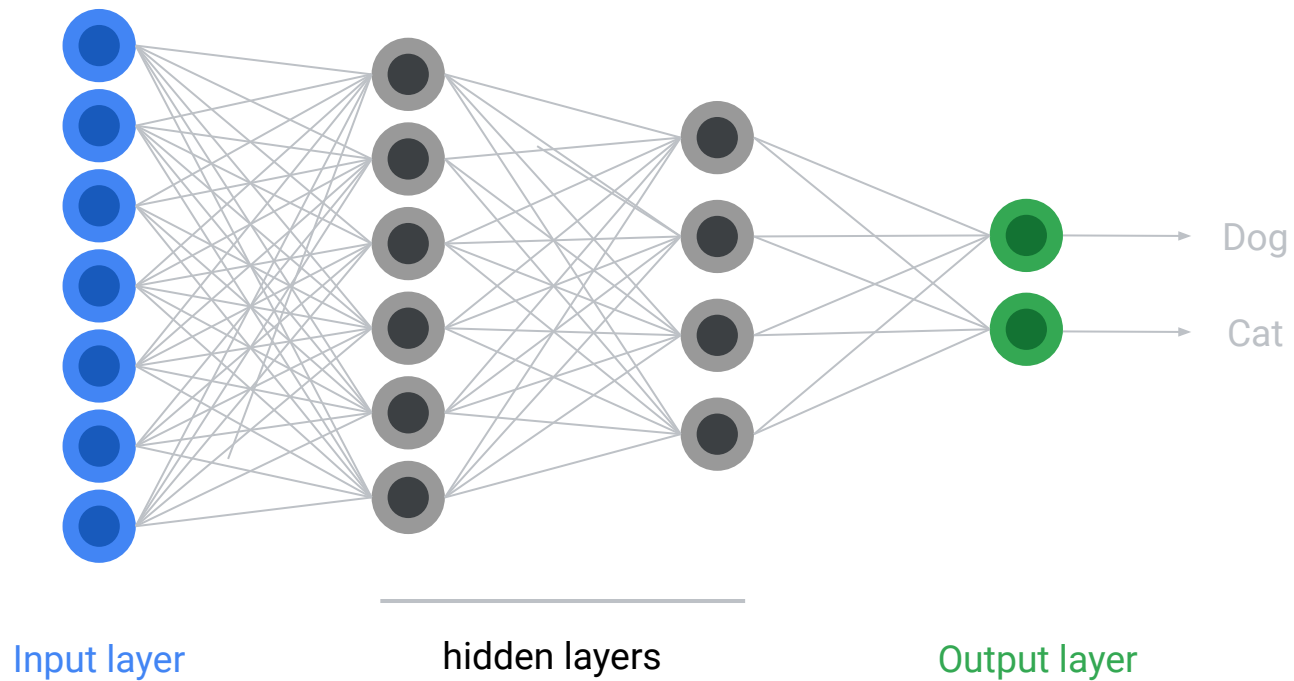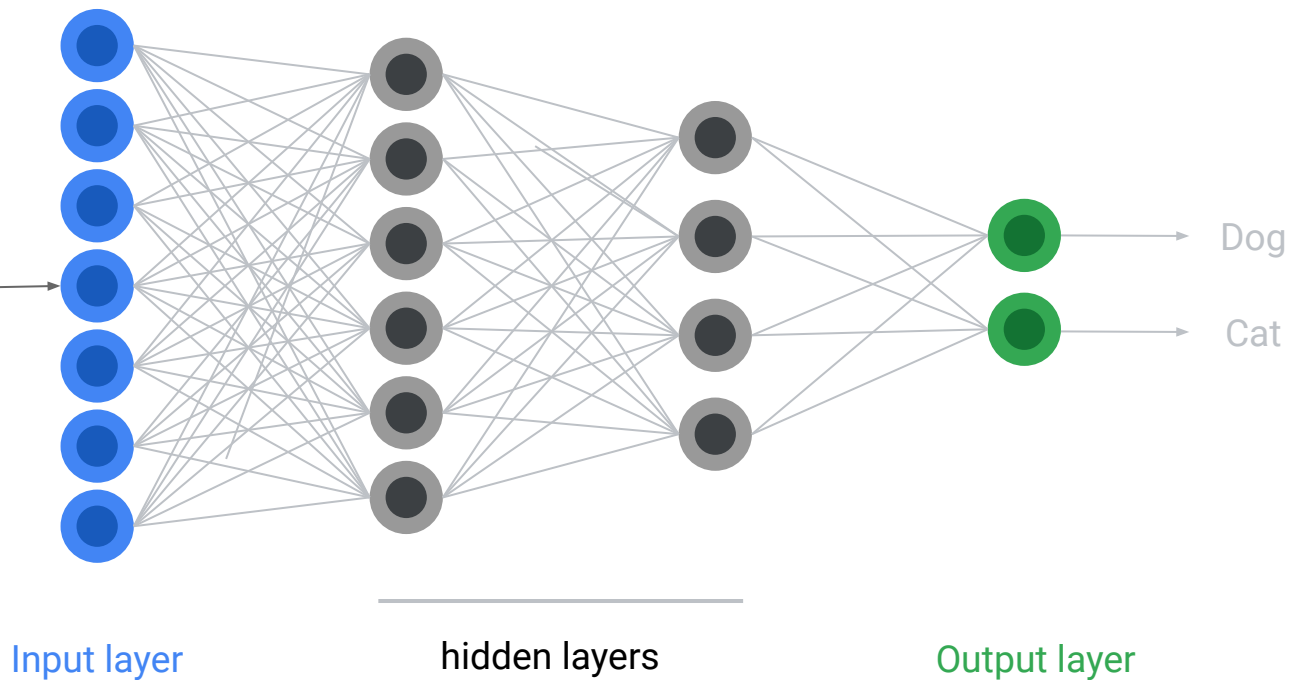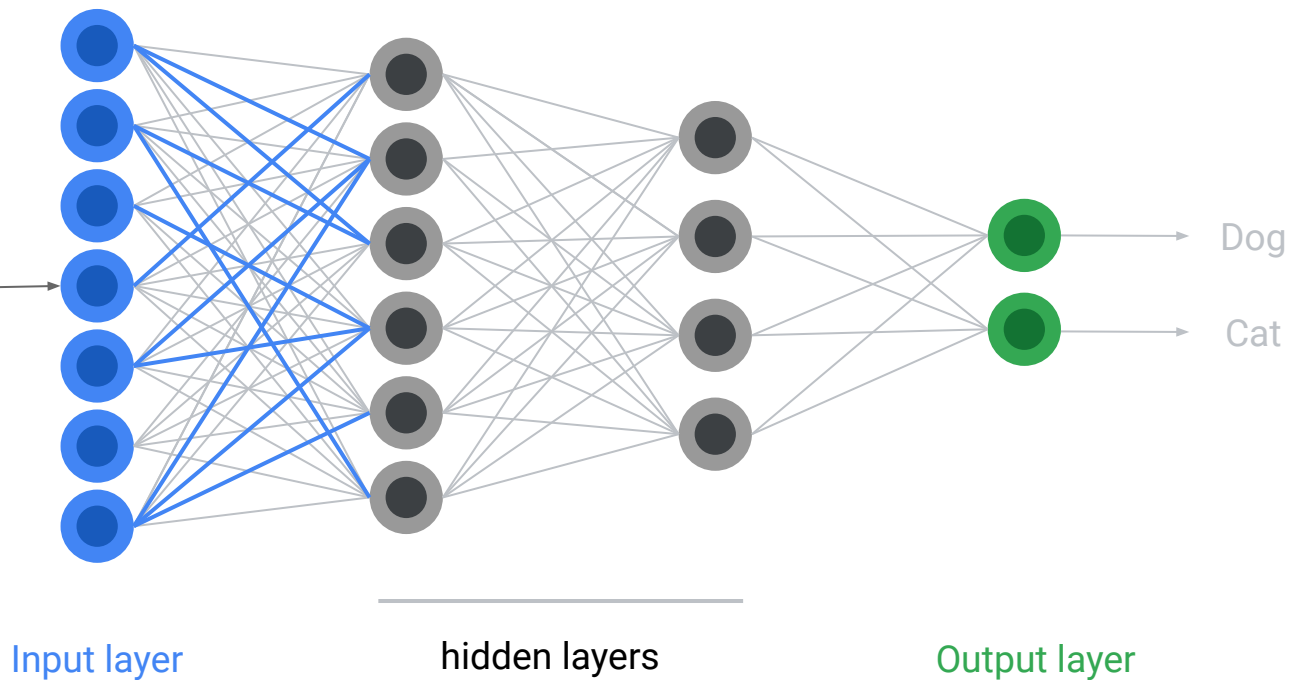Input layer      hidden layers      Output layer

Dog

Cat

Google

Security and Privacy Group

Input layer

hidden layers

Output layer

Dog

Cat

Google

Security and Privacy Group
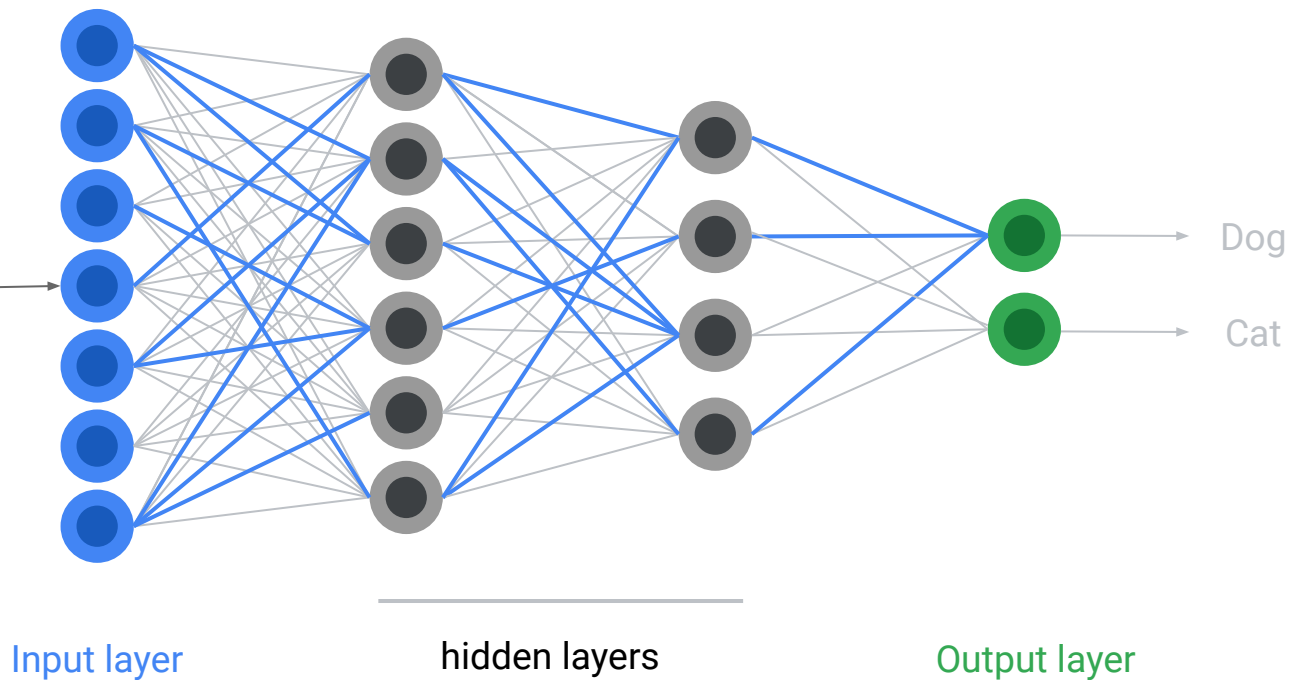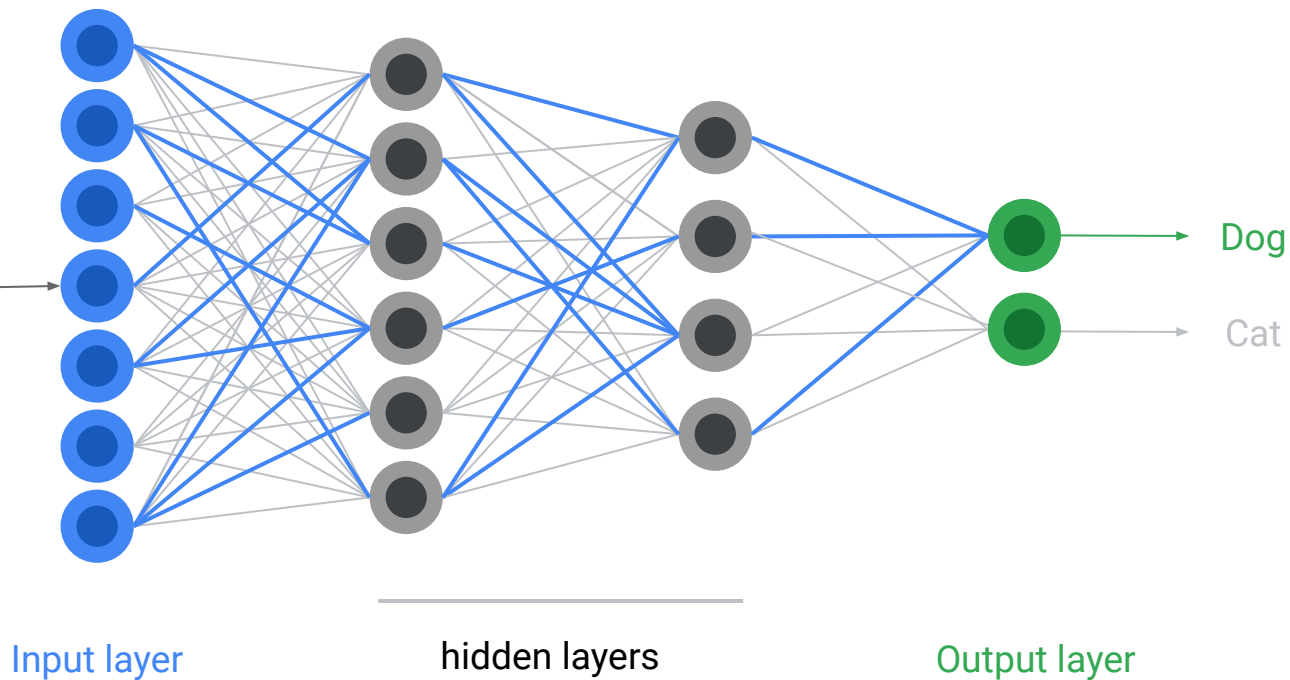
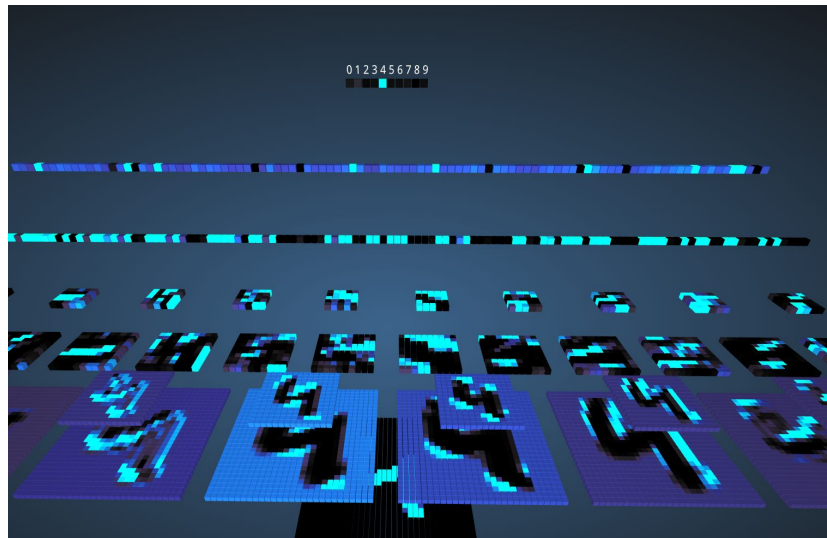# Different use-cases use different types of layers and network architectures



Handwritten digit recognition visualization from
http://scs.ryerson.ca/~aharley/vis/conv/

Google

Security and Privacy Group

What do I need to train deep learning models?

Google

Security and Privacy Group

# Tensorflow to write and train your model

Google

Security and Privacy Group

TPU v2

You need a **hardware accelerator (GPU or TPU)** as training on CPU is impossibly slow

Google

Security and Privacy Group

Demo code is available
on Colab: a hosted
Python notebook with
Tensorflow and free
GPU/TPU time

https://colab.research.google.com/

Google

Security and Privacy Group

# Hacker guide to SCAAML attacks

Google

Security and Privacy Group

Performing a
SCAAML attack
**step by step**

Google

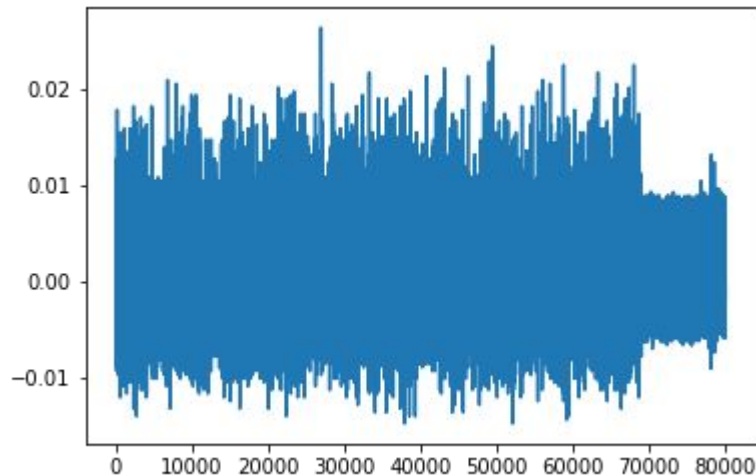Security and Privacy Group

**Goal**: train a model that can recover the AES keys from the STM32F415 TinyAES implementation using as few power traces as possible

Google

Security and Privacy Group

# SCAAML game plan



**Encryption**

**Signal acquisition (cw + pico)**

**Predictions using DNN**

**Combine DNN predictions**

**AES key!**

Google

Security and Privacy Group

Dataset is composed of **50000 raw power traces with 80000 points per trace**, without any processing or cutting, that were connected asynchronously



Sample trace from the TinyAES dataset used in this talk
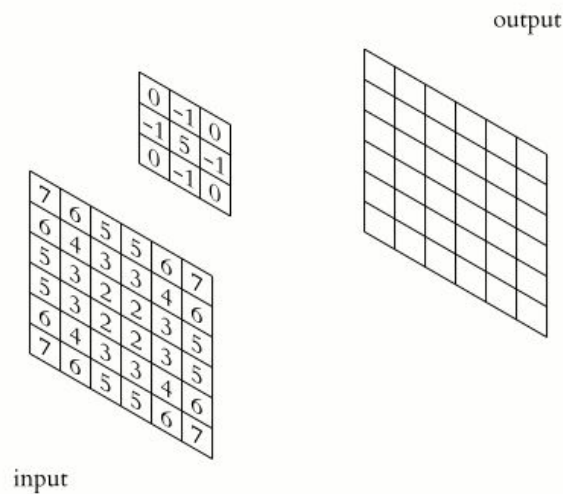
Google

Security and Privacy Group

What model
architecture to use?

Google

Security and Privacy Group

We are going to use a **convolutional network architecture**



input

output

Google

Security and Privacy Group

**Constants**

```python
dropout_rate = 0.3
filters = 32
kernel_size = 5
num_convolutions = 5
pool_size = 4
```

**Input**

```python
inputs = layers.Input(shape=(trace_size, 1))
x = inputs
```

**Pooling**

```python
x = layers.MaxPooling1D(pool_size)(x)
```

**Convolutions**

```python
for _ in range(num_convolutions):
    x = layers.SeparableConv1D(filters, kernel_size)(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    filters *= 2
```

**Pooling**

```python
x = layers.GlobalMaxPool1D()(x)
```

**Denses**

```python
x = layers.Dropout(dropout_rate)(x) # better with it
x = layers.Dense(256, activation='relu')(x)
x = layers.BatchNormalization()(x)  # helps
x = layers.Dropout(dropout_rate)(x)
```

**softmax**

```python
outputs = layers.Dense(256, activation='softmax')(x)
```
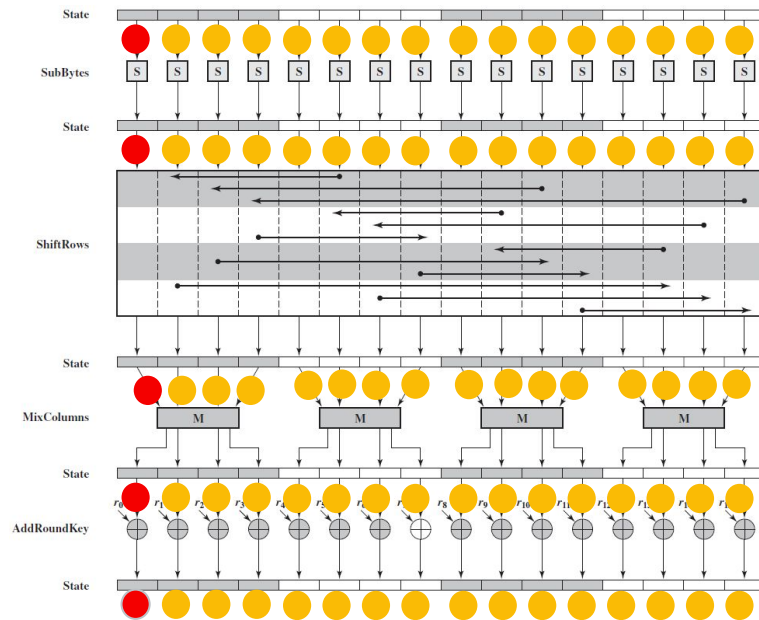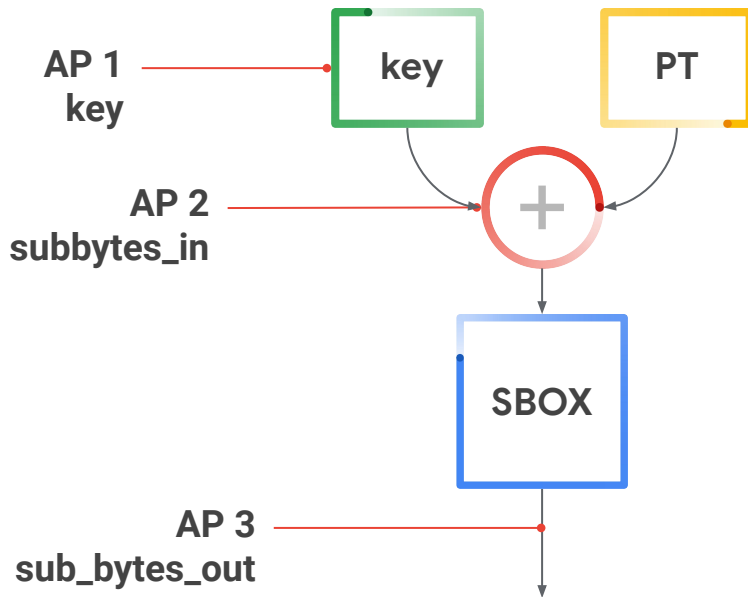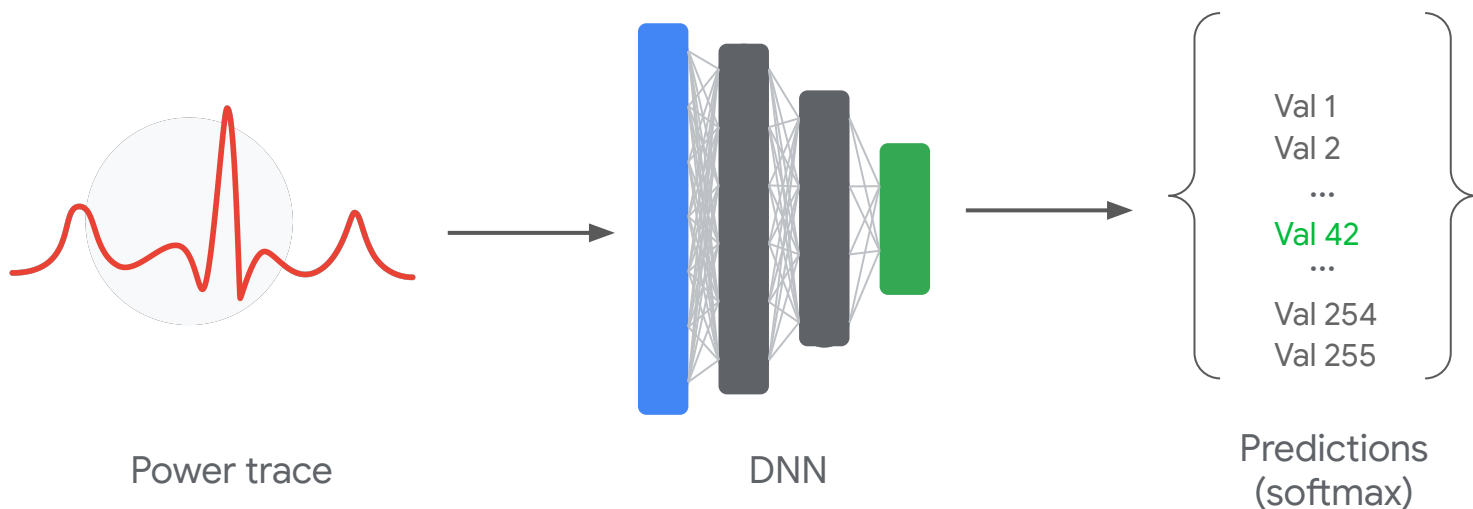
Security and Privacy Group

Google

What the model should predict?

Google

Security and Privacy Group

# AES
## attack points



Google

Security and Privacy Group

Target any of the **initial three AES attacks points** as they are **easily invertible**



AP 1
key

AP 2
subbytes_in

AP 3
sub_bytes_out

key

PT

+

SBOX

Google

Security and Privacy Group
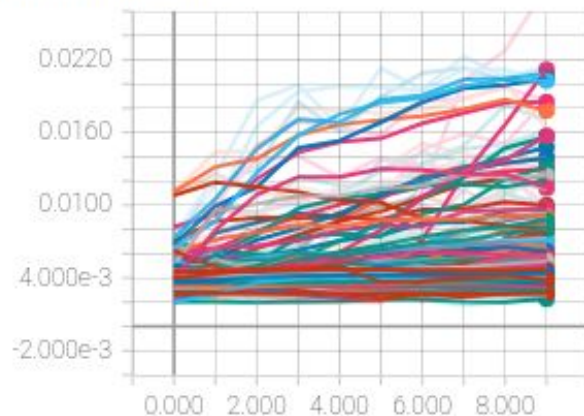
Power trace → DNN → Predictions (softmax)

Val 1
Val 2
...
Val 42
...
Val 254
Val 255

**predict a single byte at the time**
**256 predictions per model:** one for each attack point potential value

Google

Security and Privacy Group
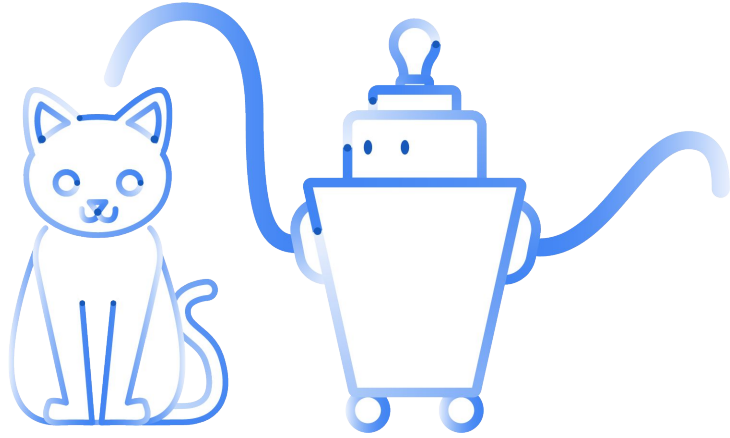
Security and Privacy Group

# Learning crypto is hard ... most models won't converge



val_categorical_accuracy

Google

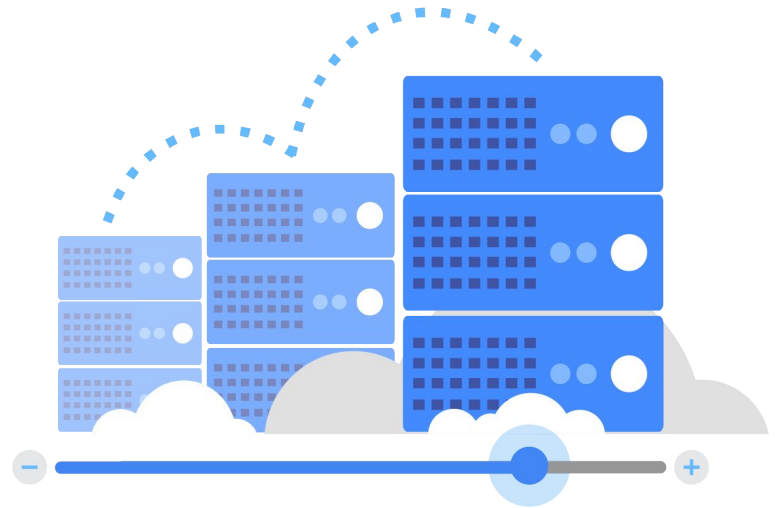Security and Privacy Group

How do I find a
model that work?

Google

Security and Privacy Group

SCAAML models **are hard to find by hand** so instead it is best to **use hyper-tuning to find the right model automatically**

Trained 1000+ to find the right one using **Keras Tuner** and **Kubernetes** on **Google Cloud**

Google

Hypertuning found **very effective models** however **none of them are simple**

Google

Security and Privacy Group

Input

Pooling

Convolutions

Convolutions with skip-connection

Pooling

Residual blocks

Pooling

Denses

softmax

Google

```python
x = inputs

x = layers.MaxPooling1D(pool_size)(x)  # helps

x = layers.Conv1D(16, kernel_size, strides=strides, padding='same',
activation='relu')(x)
x = layers.BatchNormalization()(x)

x = layers.Conv1D(filters, kernel_size, strides=strides, padding='same',
activation='relu')(x)
x = layers.BatchNormalization()(x)

for idx in range(num_convolutions):
    filters *= 2
    residual = layers.Conv1D(filters, 1, strides=strides, padding='same')(x)
    x = layers.SeparableConv1D(filters, kernel_size, padding='same')(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.Conv1D(filters, kernel_size, padding='same')(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.MaxPooling1D(kernel_size, strides=strides, padding='same')(x)
    x = layers.add([x, residual], name='sortcut_%s' % (idx))

for idx in range(nun_residuals):
    residual = x
    x = layers.Conv1D(filters, kernel_size, padding='same')(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.Conv1D(filters, kernel_size, padding='same')(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.Conv1D(filters, kernel_size, padding='same')(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation('relu')(x)
    x = layers.add([x, residual], name='residual_%s' % (idx))

x = layers.GlobalMaxPool1D()(x)

x = layers.Dense(256, activation='relu')(x)
x = layers.BatchNormalization()(x)  # helps
outputs = layers.Dense(256, activation='softmax')(x)
```
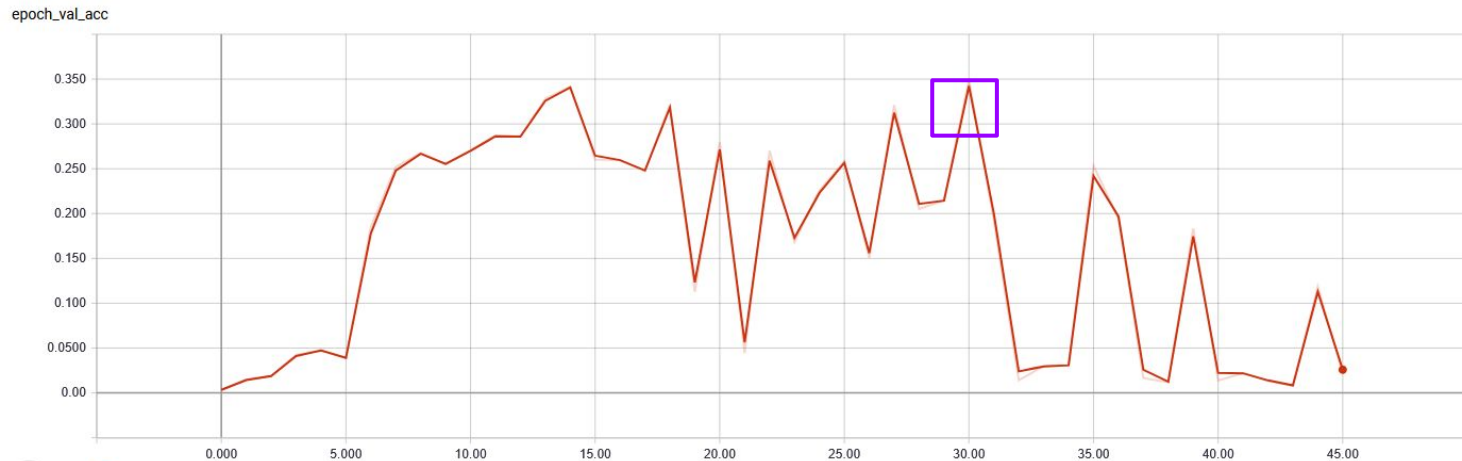
Google

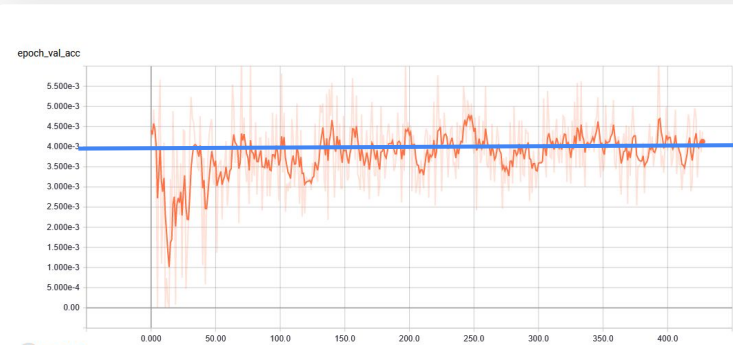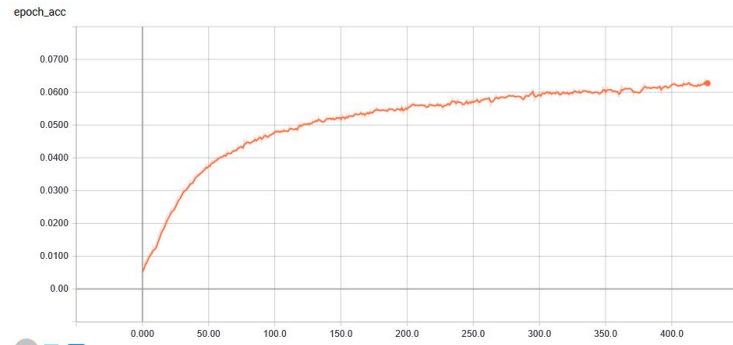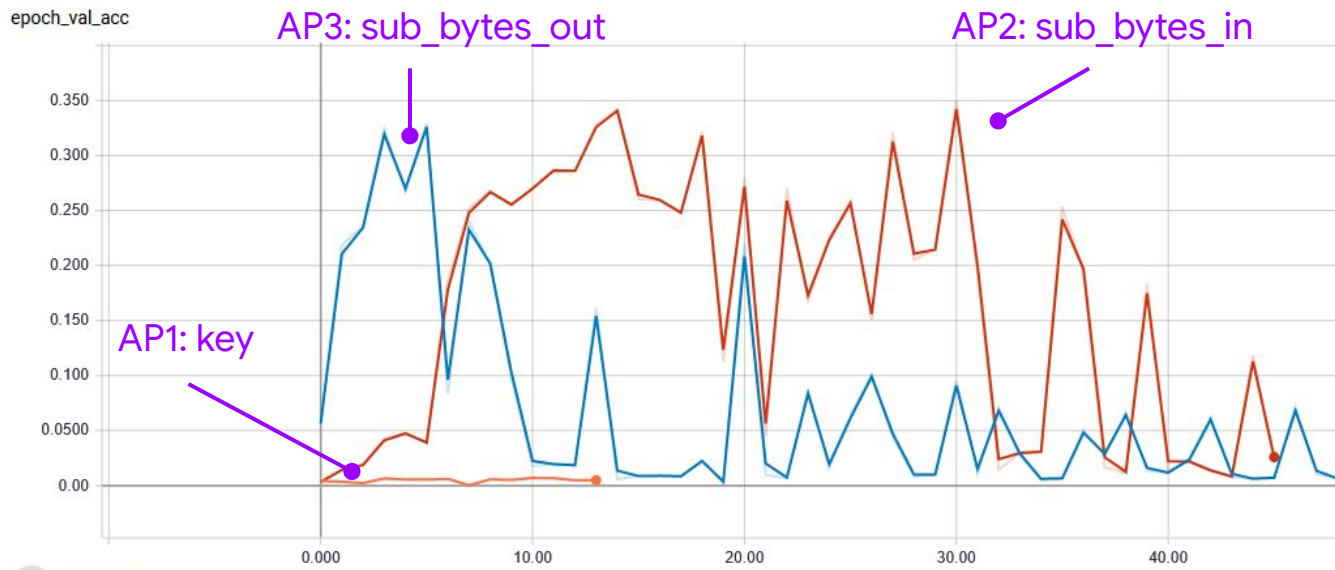Security and Privacy Group

our model reached 34.94% validation accuracy
before collapsing

Google

Security and Privacy Group

Data augmentation
can help but if badly
configured it prevents
the model from
converging



Google

Security and Privacy Group

**Choosing the right attack point matters** to get the best performance. The **best attack point varies from architecture to architecture**
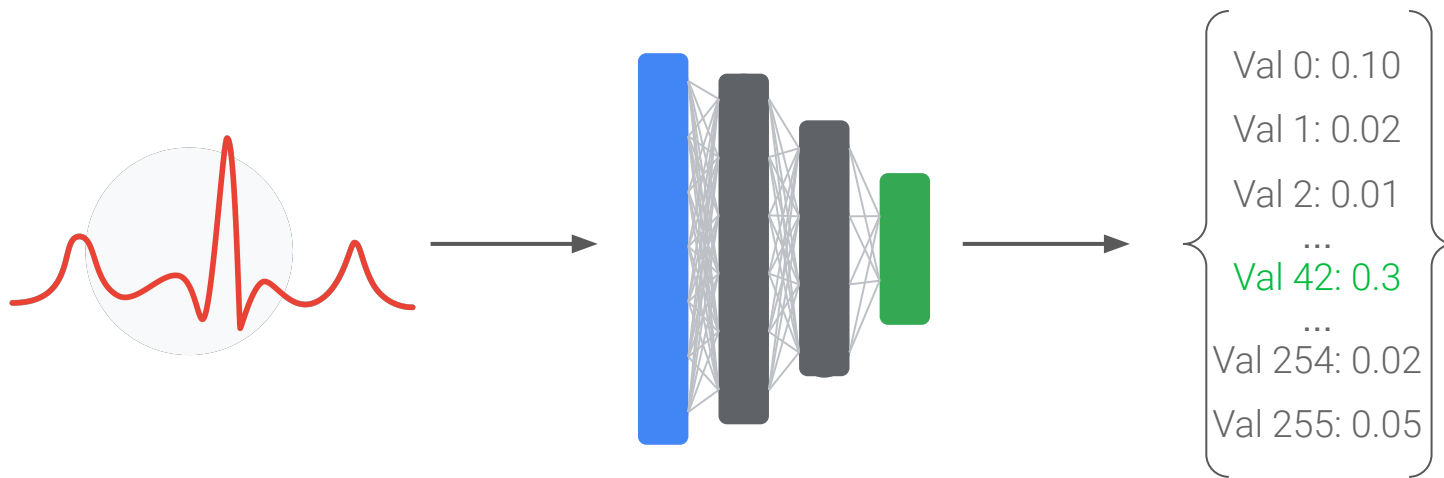
Google

Security and Privacy Group

How do I recover the key?

Google

Security and Privacy Group

Leverage all model predictions on many traces to carry out probabilistic attacks
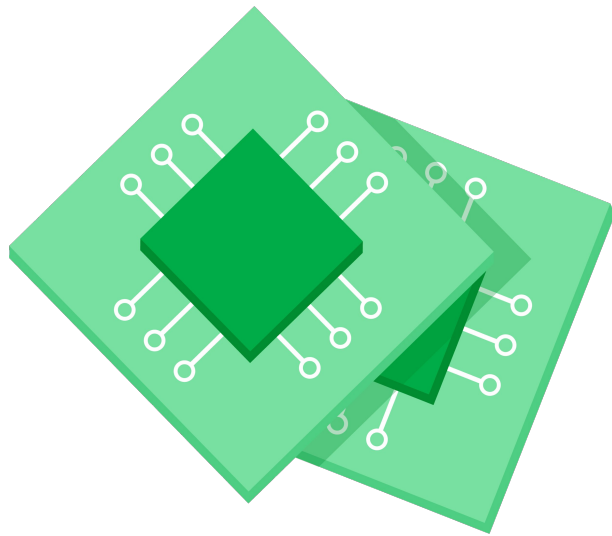
Google

# Probabilistic attack: **single trace**



Val 0: 0.10
Val 1: 0.02
Val 2: 0.01
...
Val 42: 0.3
...
Val 254: 0.02
Val 255: 0.05

Google

Security and Privacy Group

# Probabilistic attack: **summing traces\***

{
Val 0: 0.10
Val 1: 0.02
Val 2: 0.01
...
Val 42: 0.3
...
Val 254: 0.02
Val 255: 0.05
}

+ ... +

{
Val 0: 0.08
Val 1: 0.04
Val 2: 0.05
...
Val 42: 0.12
...
Val 254: 0.03
Val 255: 0.10
}

→

{
Val 0: 4.4
Val 1: 5.3
Val 2: 3.2
...
Val 42: 21.4
...
Val 254: 2.9
Val 255: 4.2
}

Google

*sum uses log10 + ε

Security and Privacy Group

# Does it work across chips?

Google

Security and Privacy Group

Use **a different chip** to create the **holdout dataset** used to **evaluate attack effectiveness**

How to evaluate attack effectiveness?

Google

Security and Privacy Group

# Success metrics

| Metric | Description | Baseline |
|--------|-------------|----------|
| Top 1 | Number of bytes correctly predicted | 0.004% (1/256) |
| Top 5 | Number of times correct byte is in top5 | 0.02% (5/256) |
| Mean rank | Average rank of the correct byte | 128 |
| Max rank | Maximum rank of the correct byte | 256 |

Google

Security and Privacy Group

Holdout dataset is composed of **100 keys** with **300 power traces** for each key that use a different plaintext
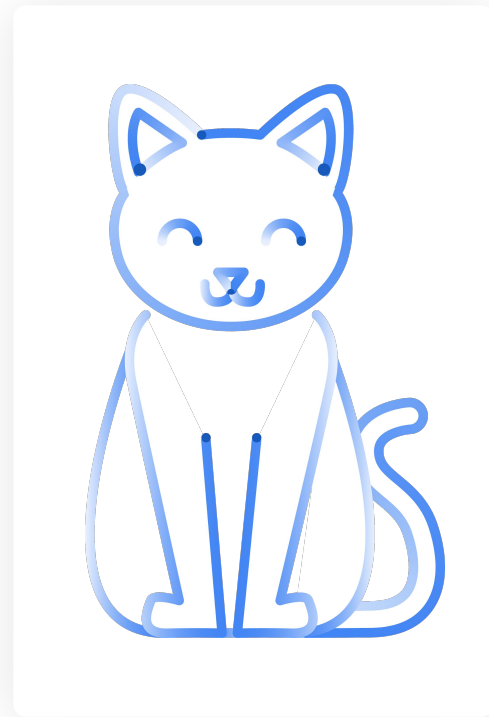
Security and Privacy Group

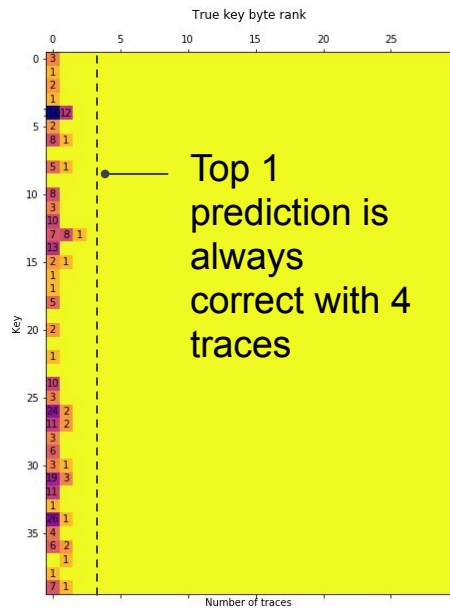Success! We recovered 100% of the keys!

Google

Security and Privacy Group

# Results: **perfect score!**

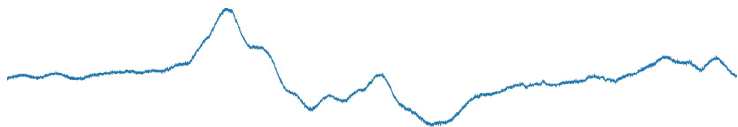| Metric | Baseline | Results |
|---|---|---|
| Top 1 | 0.004% (1/256) | **100%** |
| Top 5 | 0.02% (5/256) | **100%** |
| Mean rank | 128 | **0** |
| Max rank | 256 | **0** |

Google

Security and Privacy Group

Despite having "only a 30% accuracy" our model allows to **recover automatically 100% of the bytes with at most 4 traces (81% with a single trace!) on a different chip**

Google

Security and Privacy Group

# How about protected implementations?

Google

Security and Privacy Group

Unprotected power trace

Protected power trace

Hardened implementation needs significantly more advanced techniques, computation and data

Google

Security and Privacy Group

# What's next?

Google

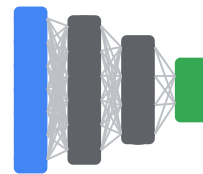Security and Privacy Group

# Testbed key numbers

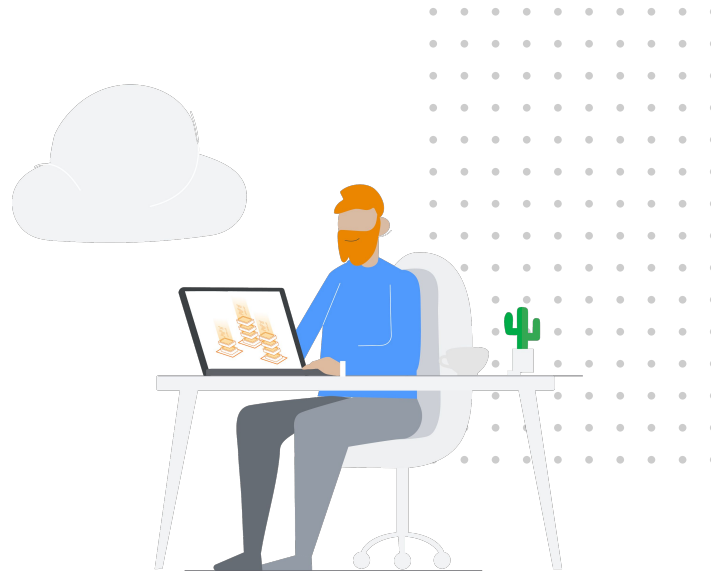**6 AES implementation**

**9M+ power traces**

**330GB storage**

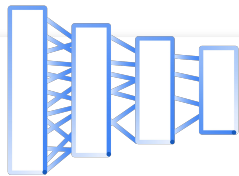**5000+ models trained**

Hope the initial draft of our paper will be public in a few weeks with our results
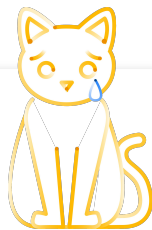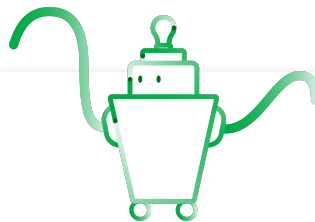
Google

Security and Privacy Group

# Takeaways



Deep-learning is the future of hardware SCA



Training model for SCA is hard



Automation is key to success



It's just the beginning

Google

Security and Privacy Group

SCAAML allow to focus on crypto algorithms design and analysis by automatically leveraging computing and AI improvements to assess their security

Google

Security and Privacy Group

Keep up with our progress on deep-learning side-channel attacks: https://elie.net/scaaml

Google

Security and Privacy Group