# Who am I

## MyreLabs, 2010, Founder&CEO
Reverse engineering
Embedded/IoT security
Secure coding
High-performance computing

## CloudVA, 2021, Founder
Side-channels as a service
Side-channels as a regression test

# Agenda

1. Intro
2. Why?
3. How?
4. Strong sides
5. Problems & workarounds
6. Demo!
7. Conclusions

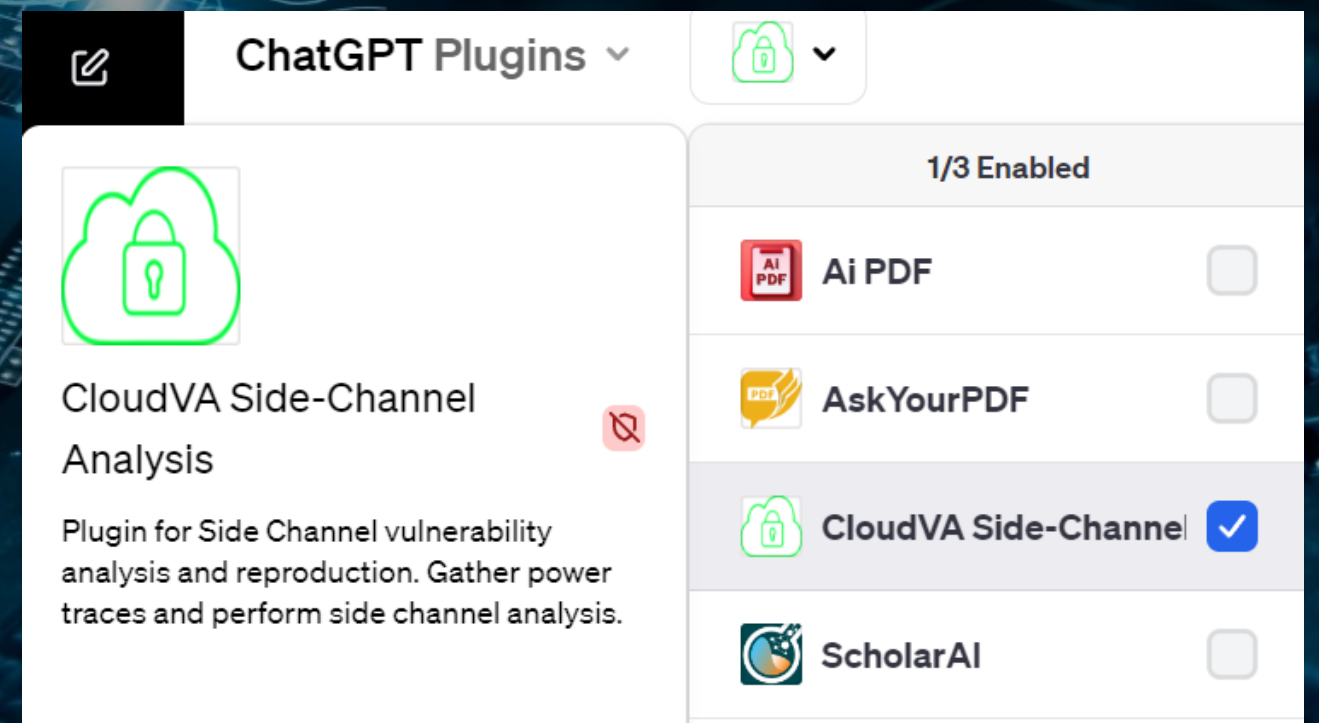# Intro: side-channels

**Easier** than you think
- High initial barrier of entry, but
- Everyone uses **the same** crypto
- Everyone uses **the same** hardware

Find a new zero-day

Re-run the pipeline

scci.yml
on: push

| ✓ start-scci | 11s | ✓ run-scci | 9s | ✓ stop-scci | 0s |

# Why: side-channels

The sad state of IoT HWSec
- Non-secure element chips
- Low clocks
- No countermeasures
- Deployed in hostile environments

Cost of attack: a laptop + $100 PicoScope
"Clever teenager" threat level



Image by SparkFun Electronics, CC BY 2.0
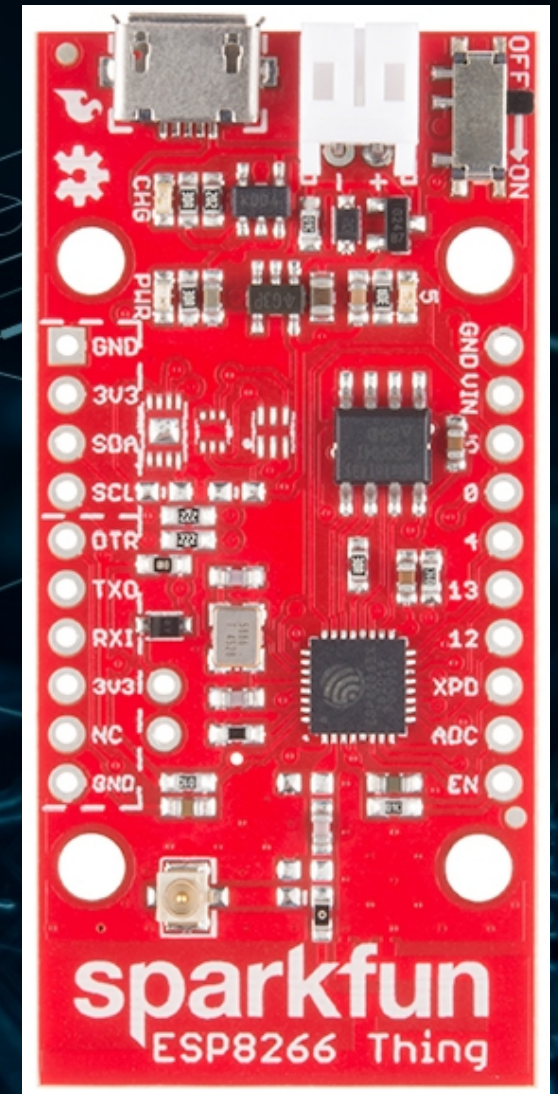
# STM32F3, mbedTLS AES Decrypt

## https://youtube.com/watch?v=CCeK_S3ED4A

# Why: GPT-4

## Exploring the limits
GPT-4 can do **time-invariant** code
What else can it do?

## Low-cost defense
Semi-**automatic** countermeasures

## Auto-training
1. **Generate** code
2. **Evaluate** leakage
3. **Learn**
4. **Repeat**

**You**
Generate C code for time-invariant simple password check

**ChatGPT**
Creating a time-invariant simple password check in C involves writing a function that compares a given password against a known, correct password. The key to making it time-invariant is to ensure that the comparison takes the same amount of time regardless of how much of the password matches. This is important for security, as it prevents attackers from using timing attacks to guess the password.

Here's a basic example of how you might implement this:

```c
#include <stdio.h>
#include <string.h>
#include <stdbool.h>


// Time-invariant password check
bool check_password(const char *input, const char *correct_password) {
    size_t input_len = strlen(input);
    size_t correct_len = strlen(correct_password);
    size_t len = (input_len > correct_len) ? input_len : correct_len;
    bool result = true;

    for (size_t i = 0; i < len; i++) {
        char input_char = (i < input_len) ? input[i] : '\0';
        char correct_char = (i < correct_len) ? correct_password[i] : '\0';

        if (input_char != correct_char) {
            result = false;
        }
    }


    return result && (input_len == correct_len);
}
```

# How: **Hardware**

- **PCIe passthrough USB controllers**
- **PPPS USB hubs**
- **ChipWhisperers**
- **ChipShouters**
- **Oscilloscopes**
- **Protocol analyzers**
- **Debuggers**

# How: Software

## Standard virtualization stack
- Linux
- QEMU KVM
- Libvirt
- **RESTful API**

```
Id    Name                                        State
--------------------------------------------------------
2     a2cac3da84779ddfc3c569e9dcd81c9efafbb4f0    running
17    e1e0e79093110e1ebe43f50d03286118726f61e0    running
```

# How: GPT-4 plugin

OpenAPI specification
- **When** to use your API
- **What** and how to call
- **How** to interpret the results
- Codegen **instructions**
- Usage **examples**

# GPT-4 plugin

```yaml
/api/v1/vms/{vmid}/gpt/cpa:
  post:
    operationId: runCPA
    summary: |-
      Run Correlation Power Analysis on previously gathered traces using provided model.
    description: |-
      Run Correlation Power Analysis on previously gathered traces
      (`tfp_uuid`) using provided model.
      Perform simple evaluation of the results.
      The analysis is executed on the virtual machine indicated by `vmid`.
    parameters:
      - name: vmid
        in: path
        schema:
          type: string
        required: true
        description: |-
          The unique id of the VM.
          A `vmid` can be obtained from `vmList` or `vmCreate` operation.
          Any VM with type `gpt` can be used.
    requestBody:
      description: CPA model definition and a pointer to target power traces.
      required: true
      content:
        application/json:
          schema:
            $ref: '#/components/schemas/runCPARequest'
```

# How: GPT-4 plugin

**Details matter**

"**Do not** generate 'main' function"

"**All** symbols except 'entrypoint' should be **static**."

"The file **may** use one of the provided cryptographic libraries (mbedtls, wolfssl)"

"Function `hw(x)` computes hamming weight of x for **bytes** and **integers**"

Our CPA endpoint specification is 187 LoC
https://cloudva.io/.well-known/ai-plugin.json

# How: GPT-4 plugin

## Examples matter more than descriptions

Token-denser
Provide context
Show usage, defaults, conventions

## Annotated examples are awesome!

```
Examples:
"hw(255-candidate)"
"hw(input[0] ^ candidate)"
"hw(sbox[input[0] ^ candidate])"
"hw(sbox[input[i] ^ candidate])" (with i set in the `model_extra` block)
"hw(a-1-b)" with a and b being input and candidate (common with comparison)
```

# Strong sides: GPT-4

## Boilerplate
- **Knows APIs better than I do**
- **Writes faster than I can**

## Translations
- **English to code**
- **Code to code**
- **Scientific formula to code**

## Self-correcting
- **C/C++ build issues**
- **Python imports**
- **It's was right about a loose cable once**

RESPONSE FROM UNKNOWN PLUGIN ⓘ

```
{
    "error": "OSError",
    "reason": "Could not detect STM32F, check connections, BOOT MODE entry setup"
}
```

It seems there was an error while gathering the power traces. The error message indicates that the STM32F microcontroller could not be detected. This could be due to a variety of reasons such as improper connections, incorrect boot mode, or hardware issues.

As I'm an AI model, I don't have the ability to physically check or modify the hardware setup. You might need to check the hardware setup and ensure that the connections are correct and the microcontroller is in the correct boot mode.

# Problems

- Context size
- Hallucinations
- BLOB handling
- Timeouts
- Arithmetics
- Metrics interpretation

# Problem: **Context size**

**Up to 20 tokens/64-bit element
256 elements = 4K+ tokens**

**April 2023:**
   **Unable to produce complete AES S-Box**
   **Unable to produce a 200 element random array**
**November 2023:**
   **Both problems resolved**
   **But still unable to process large pieces of code**

| Tokens | Characters |
|---|---|
| 17 | 18 |

`0xf1d2c3b1a7f9e4c8`

# Code token size optimization

- Process only what you **need**
- Move what you can to **separate** files and headers
- Provide **clear** naming so model can infer from **examples**
- #define constants away

| Tokens | Characters |
|--------|------------|
| 24 | 32 |

```
uint64_t a = 0xf1d2c3b1a7f9e4c8;
```

| Tokens | Characters |
|--------|------------|
| 8 | 20 |

```
uint64_t a = CONST1;
```

# ID token optimization

| | Sample | Entropy | Tok |
|---|---|---|---|
| UUID | 550e8400-e29b-41d4-a716-446655440000 | 128 bits | 18 |
| XKCD1 | Tr0ub4dor&3 | ~28 bits | 8 |
| XKCD2 | CorrectHorseBatteryStaple | ~44 bits | 7 |
| TOK4 | HeavyCostWithinLanguage | 64 bits? | 4 |



Image by XKCD, CC BY 2.5

# Problem: Hallucinations

- Calling **nonexistent** functions
- Importing imaginary modules
- **Making up** libraries
- **Ignoring** your instructions
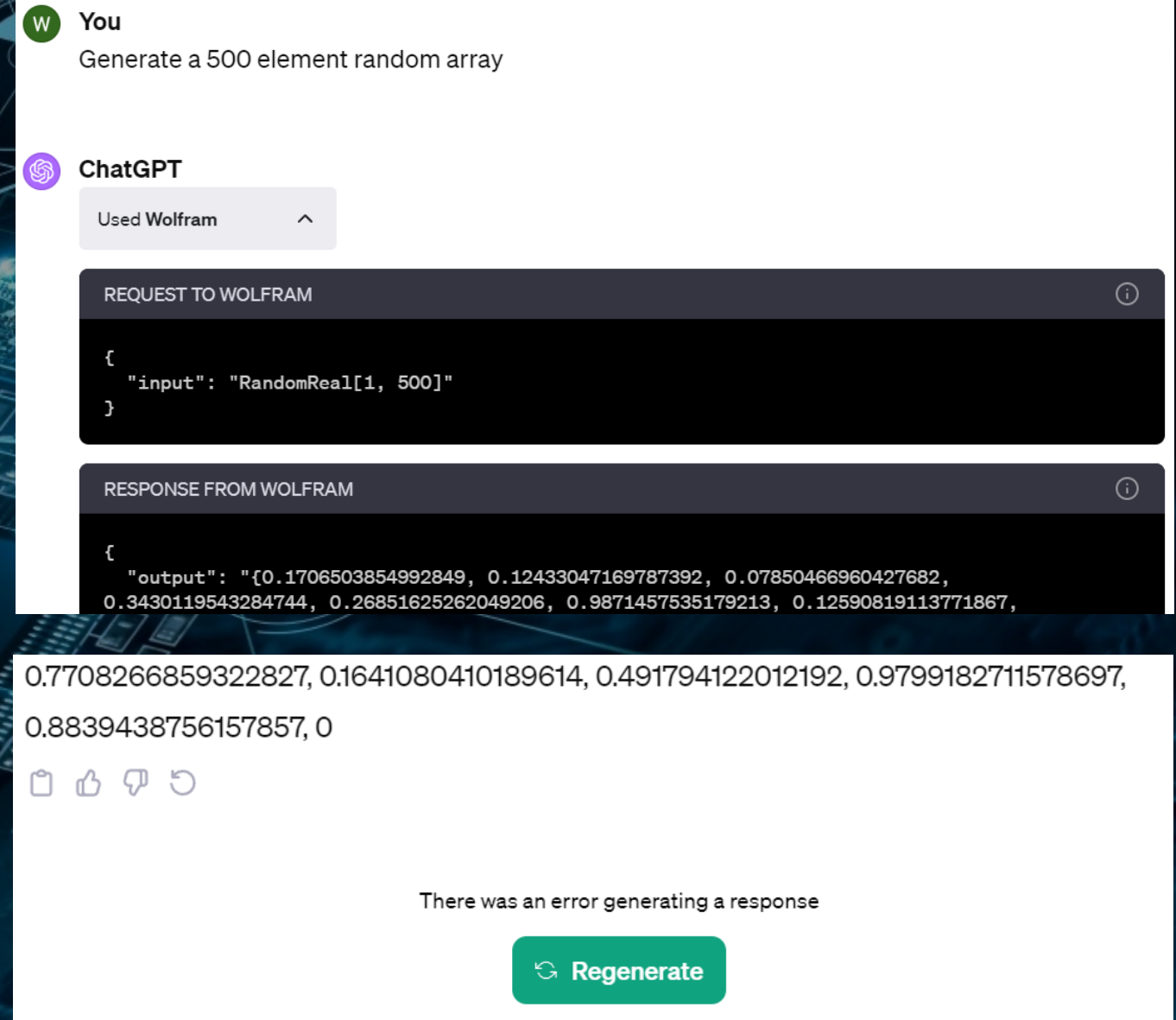- Mixing up versions

REQUEST TO CLOUDVA SIDE-CHANNEL ANALYSIS                    ⓘ

```
{
    "vmid": "e419b3ad1b8640161e2f11ec10a047e38c8319b5",
    "tfp_uuid": "a9bdebf1-761b-48ff-925d-a8b9dbcf2262",
    "model_extra": "from patools.util import hw\nimport itertools",
    "model_hypothesis": "hw(input[0] ^ candidate)",
    "model_candidates": "range(256)",
    "model_target": "secret[0]"
}
```

RESPONSE FROM CLOUDVA SIDE-CHANNEL ANALYSIS                 ⓘ

```
{
    "error": "ModuleNotFoundError",
    "reason": "No module named 'patools.util'"
}
```

# Dealing with hallucinations

- Fail **early**
- Return **meaningful** logs
- Provide examples
- Clear and **concrete** descriptions
  "is a number" vs "is uint32_t"
- Give in?

I will now run the CPA analysis again with this revised setup.

Used CloudVA Side-Channel Analysis    🚫 UNVERIFIED    ^

REQUEST TO CLOUDVA SIDE-CHANNEL ANALYSIS    ⓘ

```
{
    "vmid": "e419b3ad1b8640161e2f11ec10a047e38c8319b5",
    "tfp_uuid": "a9bdebf1-761b-48ff-925d-a8b9dbcf2262",
    "model_extra": "def hw(x):\n    return bin(x).count('1')",
    "model_hypothesis": "hw(input[0] ^ candidate)",
    _candidates": "range(256)",
    _target": "secret[0]"
```

It appears there was an error with the CPA analysis due to a missing module: `patools.util`. This module is necessary for the analysis but is not available in the current environment.

Let's revise the CPA setup with a custom implementation of the Hamming weight function:

```python
def hw(x):
    return bin(x).count('1')
```

# BLOB handling

- **Store** data **on your end**
- **Even relatively small amounts!**
- **Only** pass **IDs** around
- **Use ID token optmization**
  (yes, we don't follow our own advice)
- **Provide a CRU(D) API**
  **Delete** needs extra attention

Used CloudVA Side-Channel Analysis  🚫 UNVERIFIED  ^

REQUEST TO CLOUDVA SIDE-CHANNEL ANALYSIS  ⓘ

```
{
    "vmid": "e419b3ad1b8640161e2f11ec10a047e38c8319b5",
    "dut_uuid": "f649974e-1dc1-4209-8282-0fa4475530b7"
}
```

RESPONSE FROM CLOUDVA SIDE-CHANNEL ANALYSIS  ⓘ

```
{
    "tfp_uuid": "a9bdebf1-761b-48ff-925d-a8b9dbcf2262"
}
```

The power traces have been successfully gathered from the

# Fixing timeouts

## **Optimize** response times
- **Our provisioning takes 0.7s**
- **VMs respond to ping within 3s**
- **Are ready to serve GPT endpoints within 10s**
- **... and that's not fast enough!**

## Build **asynchronous** APIs
- **Return request ID immediately**
- **Perform work in the background**
- **Implement status/progress functions**
- **Implement wait/sleep on your end**
- **Return 503 in reasonable time**

```
get:
  operationId: vmList
  summary: |-
    List available VMs.
  description: |-
    List available VMs. Any vm with type `gpt` ma
    Therefore, before executing `vmCreate` it's b
    an available VM that can be already used.
    A VM must be `running` before it can be used
```

```
/api/v1/vms/{vmid}/gpt/build:
  post:
    operationId: buildDUTCode
    summary: |-
      Prepares device under test code.
    description: |-
      Prepares device under test (DUT) code, returing a unique id (dut_uuid).
      The bulid procedure is executed on the virtual machine indicated by `vmid`.
```

```
"503":
  description: VM is not running. It must be restarted using vmStart.
  content:
    application/json:
      schema:
        $ref: '#/components/schemas/empty'
```

# Conclusions

**Can GPT-4 do side-channels?**
  **Yes**
**Does it solve hard problems?**
  **No**
**Is it useful yet?**
  **Depends**

# Takeaways

**How to build a GPT-4 plugin**

**How to work with code**

**Side-channels are easy!**

**Thank you!**

✉ witold.waligora@cloudva.io

𝕏 @WaligoraWitold

in linkedin.com/in/witold-waligora

#BHEU   @BlackHatEvents

# Extras

#1 [live] Iterative CPA on simple password

#2 Fun prompts:

"Change leak model to [expr]"

Just works

"Bisect to find out how many traces it takes"

Correct flow, but hits token limits eventually

"Extract leak model from publication: [upload]"

Biggest problem seems to be loading the paper with its scientific notation intact