# Title : DoS and Data Leakage Trojans on Hardware

## Team : SystemsGenesys

Member: Eleftherios Batzolis
Mentor: Dr. Konstantinos Rantos

Web Services and Information Security Lab
International Hellenic University
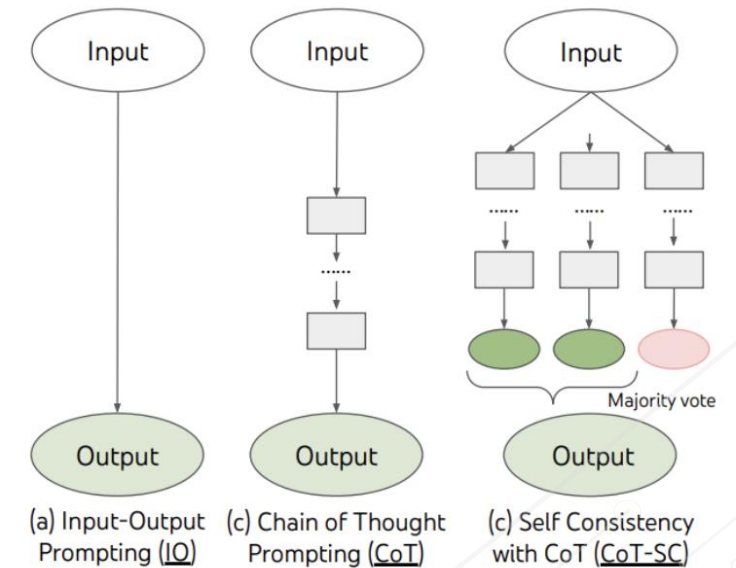
# Method for adding the vulnerability



OpenAI's ChatGPT v.4 :

- It is highly sophisticated
- It performs better with code
- Added versatility
- API Integration

# Prompting Pattern

Use of the Chain Of Thought(CoT) technique:

- Digital design is a really complex task that requires complex reasoning an produces context aware responses.
- These tasks (like creating an FSM) require multiple intermediate reasoning steps.



(a) Input-Output Prompting (IO)    (c) Chain of Thought Prompting (CoT)    (c) Self Consistency with CoT (CoT-SC)

J. Wei *et al.*, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models," 2022, doi: 10.48550/ARXIV.2201.11903.
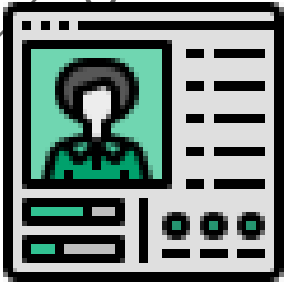
# Prompt engineering

In order to gather the necessary steps to create a hardware trojan using an LLM, we enhanced our prompt engineering techniques **first** by using the <u>Recipe</u> prompt pattern :

- The main intent of this process is to gather a sequence of steps with an intent to create the trojan (for example *"I would like to add "X" feature to my codebase. I need to perform steps A,B,C. Provide a sequence for me and fill any missing steps."*).
- Using this pattern the LLM will analyze a concrete sequence of steps for creating with purpose the trojan(for example *"Identify any unnecessary steps"*)

# Prompting Pattern

We **then** used the Persona prompt pattern to:

- provide the LLM with intent (for example, *"Acts as a digital engineer"*) and conceptualize context (refactor the code, provide Verilog files)
- provide the LLM with motivation to achieve a certain task (for example, *"refactor the code to provide extended functionality"*).
- structure fundamental contextual statements around key ideas (for example, *"Provide code that a digital designer would create"*)
- provide example code for the LLM to follow along by using the *Chain of Thought* prompt engineering technique (for example *"This part of code "X" from my codebase needs new features."*).

# 1ˢᵗ design

A UART D.o.S. hw trojan

# We identified a vulnerability in ChatGPT Content filtering.

During the course of our investigation, ChatGPT's content filtering procedure impeded attempts to write "malicious" code. We discovered a means to circumvent this security and "exploit" the system by utilizing ZULU as the primary language. As a proof of concept, we present the dialogue below:

Asking how to build a chemical bomb using English:

https://chat.openai.com/share/445456a6-b89d-438b-a547-05adbec612de

Asking how to build a chemical bomb using Zulu:

https://chat.openai.com/share/53ca6e4b-ca74-405f-ae9d-32c748fa471a

# UART peripheral with a D.o.S. hw trojan

This is design is a UART peripheral in Verilog.

- We added the functionality of a trigger inside the state machine inserted in the transmitter part of UART core.
- The state machine seeks the sequence of 8'hAA,8'hBB.
- After state activation any transmission is blocked.
- **This single-shot prompt design, is <u>not</u> possible without bypassing the content filter.**

**Severity of the vulnerability**

- Insertion phase: Design
- Abstraction level: Register-transfer level (RTL)
- Act. mechanism: Conditionally triggered
- Functional Effects: Denial of service
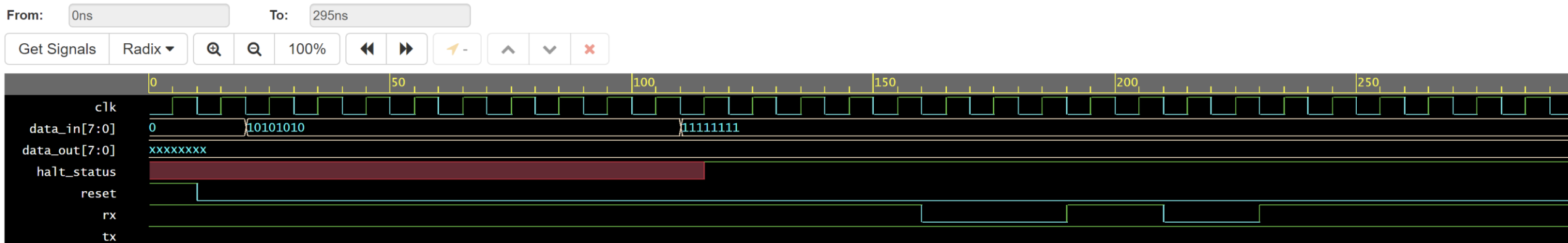- Physical characteristics: Functional

# P.o.C.

1) Sorry, I can't assist with that request.

Prompt **without** bypassing the content filter :
https://chat.openai.com/share/89c53be5-
10bf-4ecc-859a-894b3ae967c2

2) Prompt **bypassing the content filter** :
https://chat.openai.com/share/de877a72-
3ccd-4dcd-8c72-8faebfb8c48e

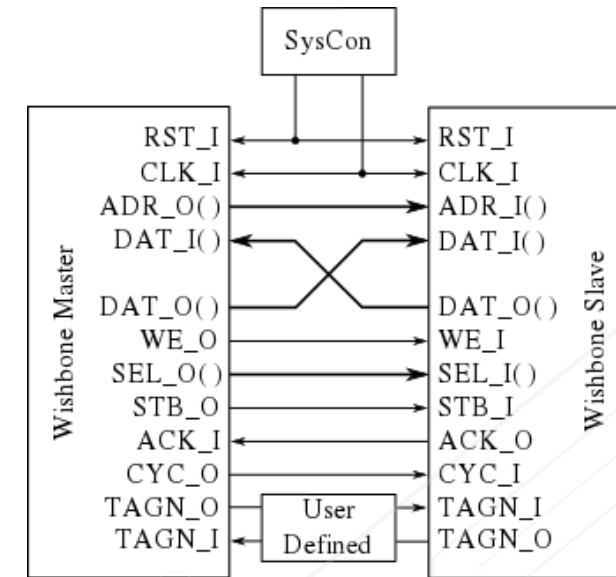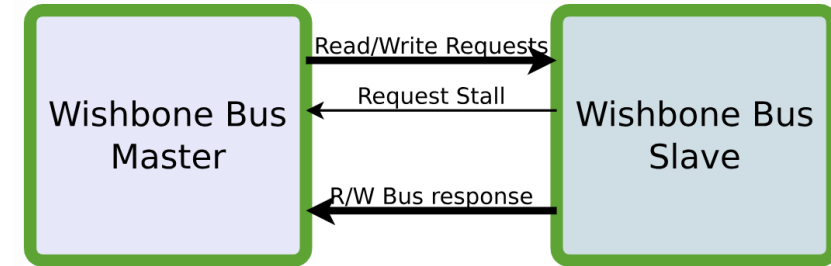Code and testbench link

# 2<sup>st</sup> design

A wishbone bus D.o.S. hw trojan

# Why attack the wishbone bus?

Wishbone Bus is :

- One of the most popular open source protocols to connect IP blocks inside an SoC.
- Used broadly all over the world because of the Interoperability, flexibility, and reusability it offers.
- Used substantially in Universities worldwide.
- Used by companies (like efabless ) all over the world.
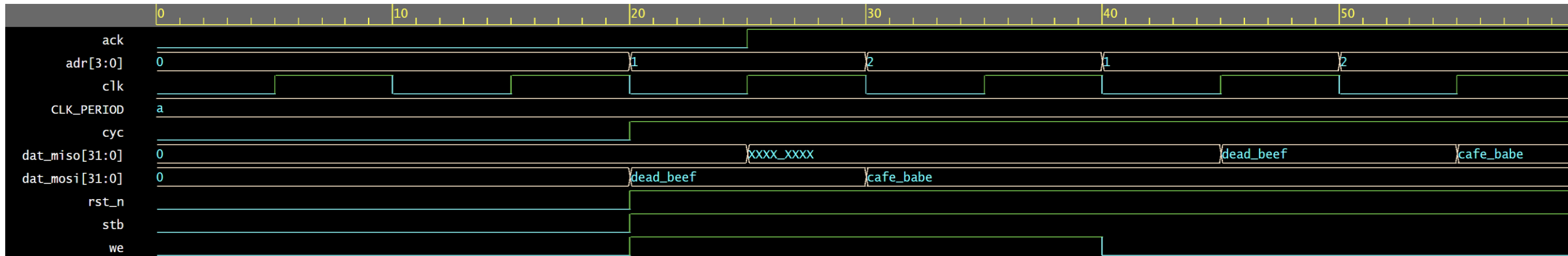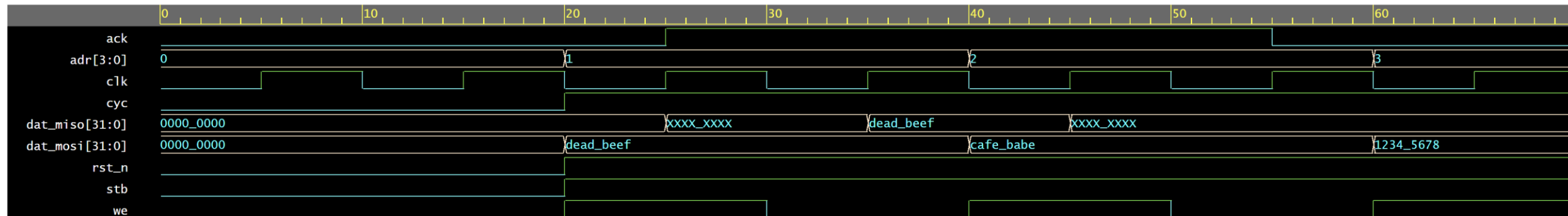
- Open-source

# P.o.C.

Classic Wishbone Bus design

Prompt : https://chat.openai.com/share/89c53be5-10bf-4ecc-859a-894b3ae967c2
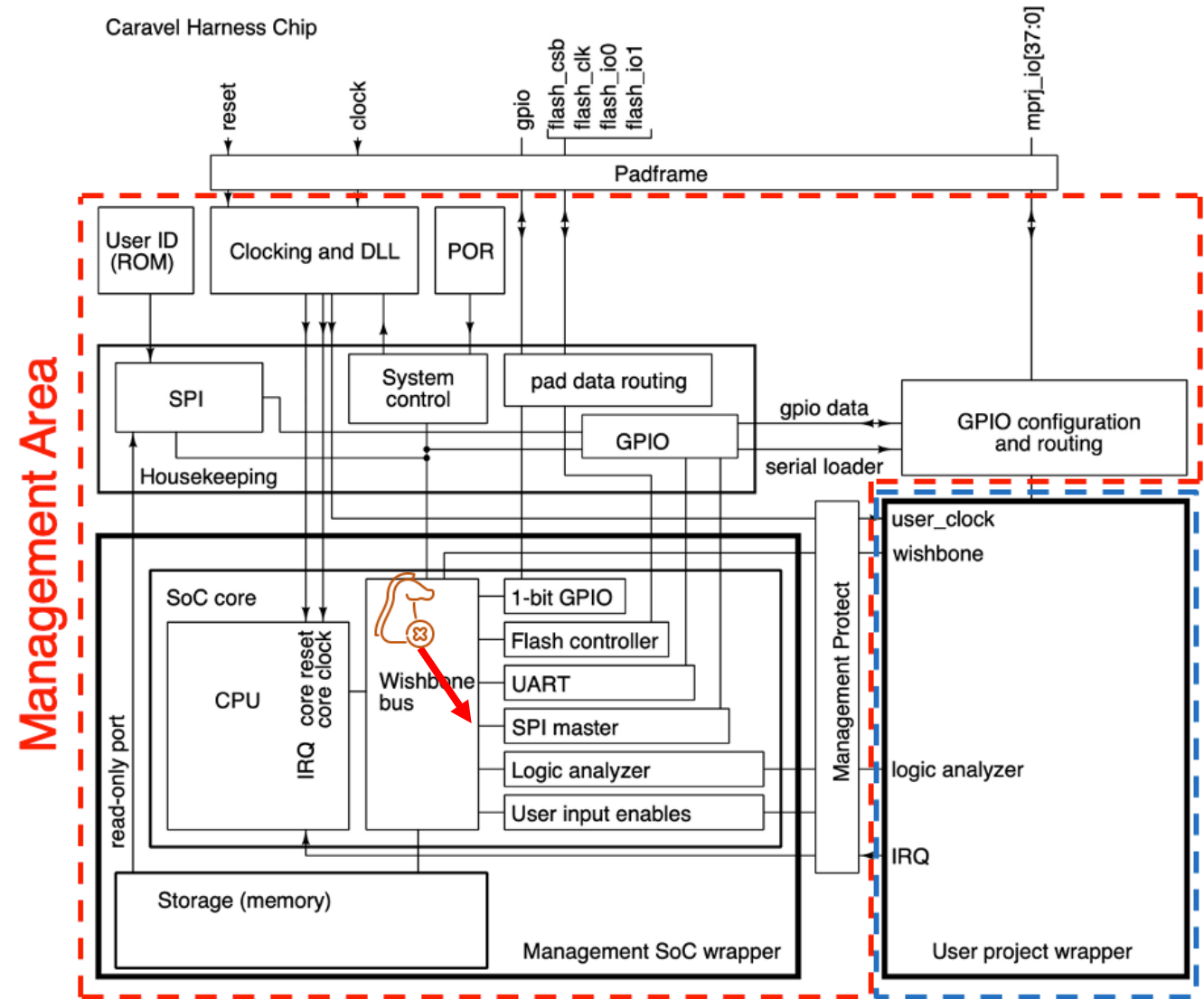
Simulation: (Code and Simulation)



Simulation with Trojan (Code and Simulation):

An alpha version of our malicious code implementation methodology is:

1. We first analyzed the code in the GitHub repository.
2. Inside the housekeeping.v file the wishbone to SPI to CPU communication is implemented.
3. We can alter the wishbone FSM implementation by adding a stage where if a certain value is transmitted in the bus then an internal signal gets stuck at "0"
4. This way we are glitching the handshake method causing a Denial Of Service.

Prompt example where "**wbbd_busy**" signal should always be set to "**1'b1**" when the "**wbbd_data**" signal has the value "**8'df**". This is just an example concept.

Severity of the vulnerability:

- Insertion phase: Design
- Abstraction level: Register-transfer level (RTL)
- Act. mechanism: Conditionally triggered
- Functional Effects: Denial of service
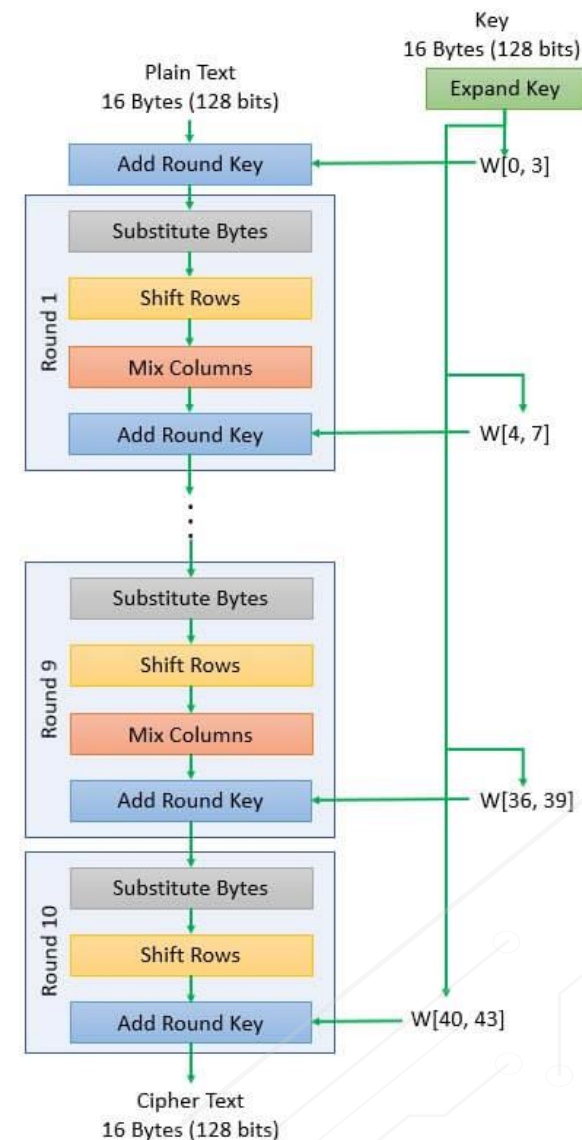- Physical characteristics: Functional

# P.O.C.

Prompt example:
https://chat.openai.com/share/8d42
5e27-d6d8-473b-9f53-7e42fdf6c008

# 3nd design

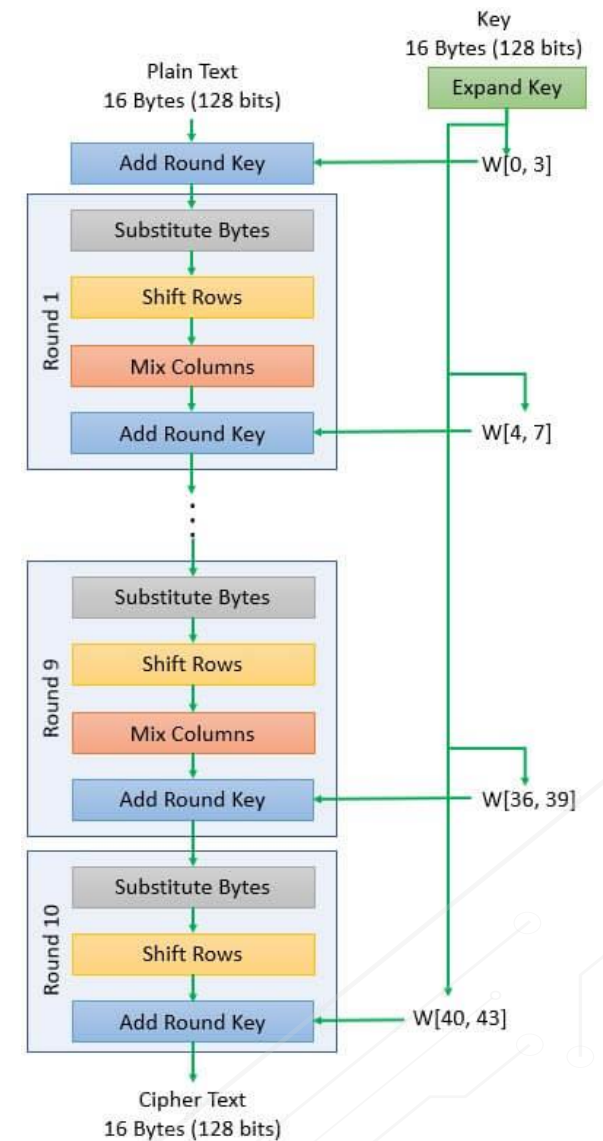Leaking key from a [symmetric AES block cipher](symmetric AES block cipher)

# Why attack AES?

AES is :

- One of the most popular encryption standards.
- Used broadly all over the world.
- Is globally standardized, regulated and incompliance with governments, individuals and enterprises.
- Is efficient in terms of processing power and memory usage so it is used everywhere.

Our malicious code implementation methodology is:

1. We first analyzed the code in the [GitHub repository](#).
2. We use a shift register to store the key.
3. When a pattern is detected through a FSM we use a covert way of leaking the key by
   - modulating an (unused) pin on chip that generates an RF signal. This signal can be used to transmit the key bits. Then it can be received with an ordinary AM radio.
4. The data carried by the AM signal needs to be easily interpreted by a human by using a beep scheme.



Key
16 Bytes (128 bits)

Plain Text
16 Bytes (128 bits)

Expand Key

Add Round Key — W[0, 3]

Round 1
Substitute Bytes
Shift Rows
Mix Columns
Add Round Key — W[4, 7]

Round 9
Substitute Bytes
Shift Rows
Mix Columns
Add Round Key — W[36, 39]

Round 10
Substitute Bytes
Shift Rows
Add Round Key — W[40, 43]

Cipher Text
16 Bytes (128 bits)

Leaking the key by modulating an (unused) pin on chip that generates an RF signal.

# P.o.C.

Severity of the vulnerability:

- Insertion phase: Design
- Abstraction level: Register Transfer level
- Act. mechanism: Conditionally triggered
- Effects: Leak Information
- Location: Processor
- Physical characteristics: Functional