# Deep Learning Steganography

## to Hide Malware in Web Content

**REDOCS**

**THALES**
Together • Safer • Everywhere
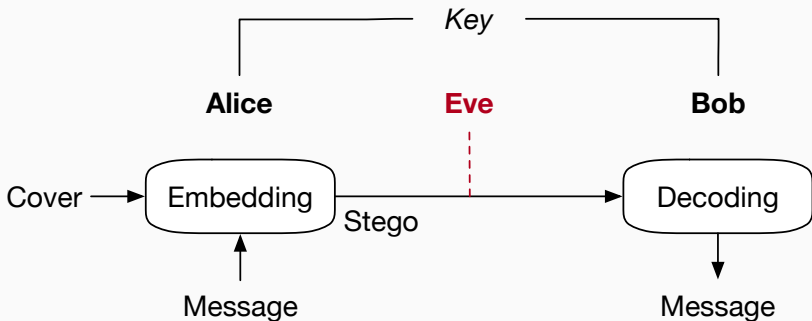
**N. Fournaise** (Univ. Limoges)
**H. Nguyen** (ENS Lyon)
**F. Recoules** (CEA LIST)
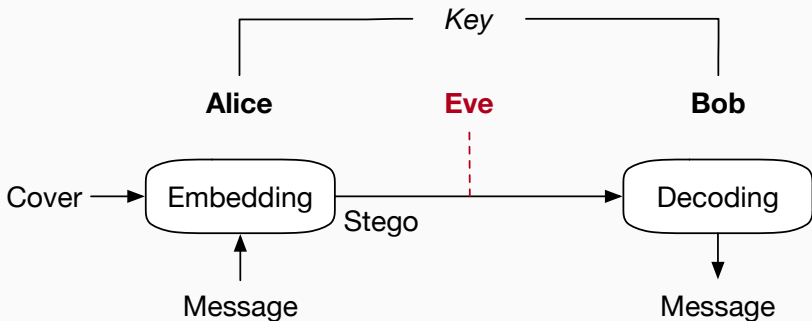**T. Taburet** (Centrale Lille, Univ. Lille, CNRS)
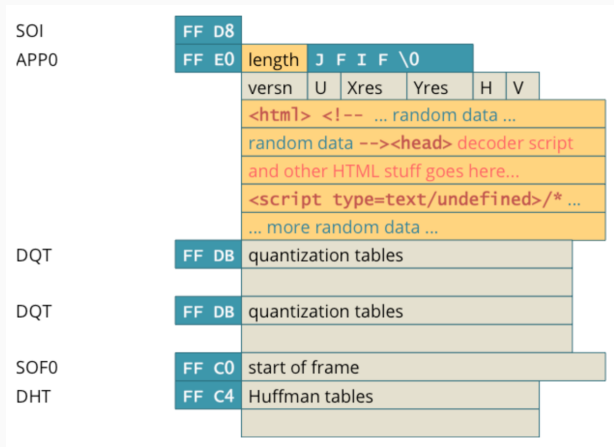
October 25, 2019

# Steganography in a nutshell
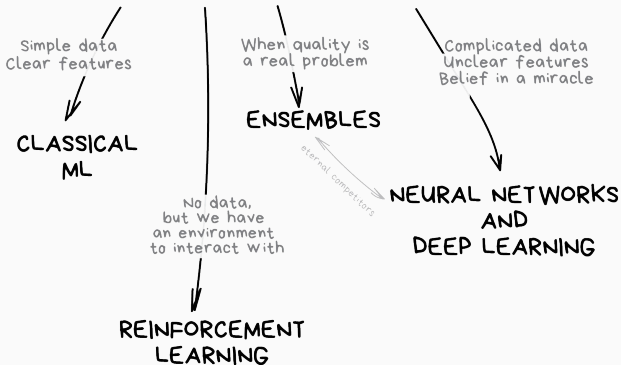
**Message can be a malicious code**

**Polyglot** *(Noun)* :
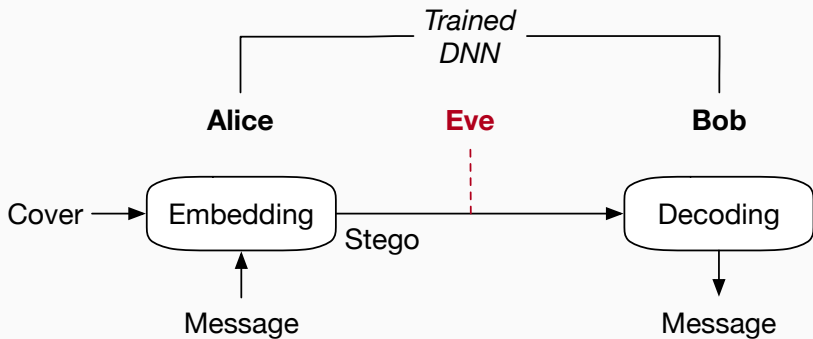
a person who knows and is able to use several languages.

source: *https://vas3k.com/blog/machine_learning/*



THE MAIN TYPES OF MACHINE LEARNING

Simple data
Clear features

When quality is
a real problem

Complicated data
Unclear features
Belief in a miracle

CLASSICAL
ML

ENSEMBLES

No data,
but we have
an environment
to interact with

eternal competitors

NEURAL NETWORKS
AND
DEEP LEARNING

REINFORCEMENT
LEARNING

**Huge claims about capacity & security**

Can we implement a Machine Learning-based steganographic decoder using web technologies?

**Disclaimer**

- No GPUs in our laptops
- No Machine-Learning Background

**Embedding decoder in a browser**

## SteganoGAN

- A Steganography algorithm from a Generative Adversarial Network
- Unpublished but public article
- Implementation Available
- Huge claims about capacity and security!

**(a)** Original image     **(b)** Basic encoder     **(c)** Dense encoder

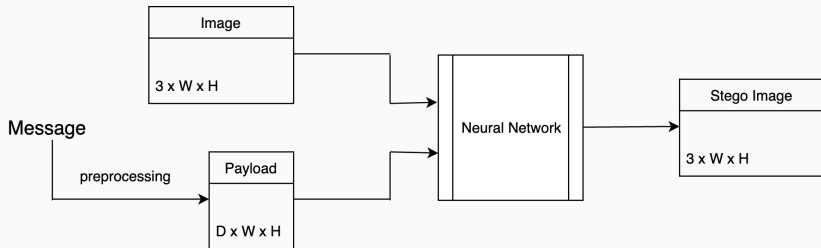## Goal

Adapt the SteganoGAN Decoder part to browser-compatible technologies

One candidate:

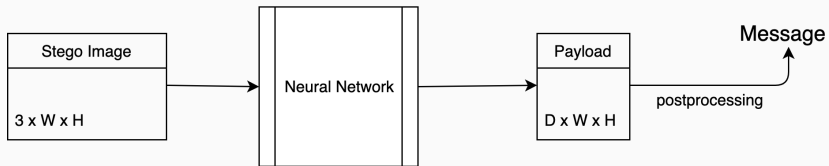**Three** already trained versions of the Neural Network

## Essential components

Embed the decoding part *only*

Components:

- Tensor from Image: Python $\Rightarrow$ JavaScript
- Tensor manipulation: PyTorch $\Rightarrow$ TensorFlow
- Neural Network inference: PyTorch $\Rightarrow$ TensorFlow $\Rightarrow$ TensorFlow.js
- Message extraction: Python $\Rightarrow$ JavaScript

Before the encoding part:

RS_ECC(zlib_compression(utf-8_encoding(**Message**)))  |  0x00000000

Repeated until no more space in a vector of size D x W x H

After the (decode) neural network inference:

PIS_ECC(Nb_compression)utf-8_encoding **Message** ||  0x00000000

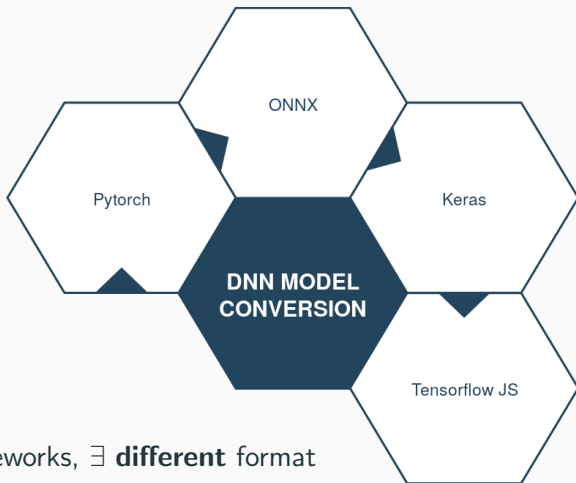Repeated until no more space in a vector of size W x H

## How we improved it

- Define the separators as 4 bytes with a *hamming_distance* $< n$
- Use the separators to deduce message length
- From message length, compute the most common bytes for every message character
  - zlib optional
  - Reed-Solomon ECC no longer needed
  - Much faster decoding for large images

## How we improved it

- Define the separators as 4 bytes with a *hamming_distance* $< n$
- Use the separators to deduce message length
- From message length, compute the most common bytes for every message character
  - zlib optional
  - Reed-Solomon ECC no longer needed
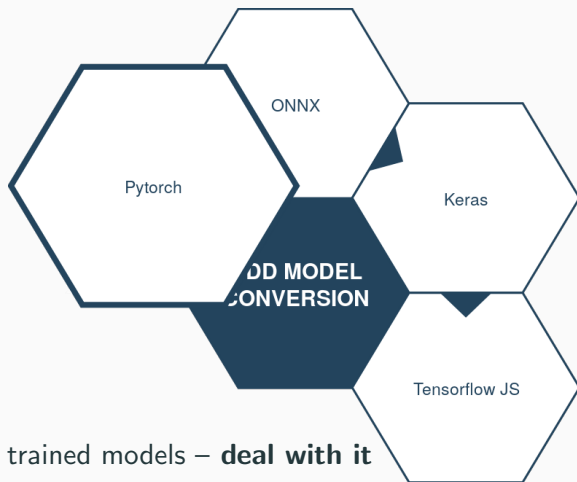  - Much faster decoding for large images

  **Easier to translate to JS**

ONNX

Pytorch

Keras

**DNN MODEL CONVERSION**

Tensorflow JS

$\forall$ frameworks, $\exists$ **different** format

Already trained models – **deal with it**

**ONNX** graph was **easy** to get

**Keras** model *import* failed – until fixed

Fixed *loading* but **wrong** *inference*

*Happiness*

**The end**

Stand-alone
TensorFlow
Decoder

...but all
components
work
independently

*JavaScript
time!!!*

Merci
Julien !

Stand-alone
PyTorch
Decoder

The NN model
cannot be
correctly
imported...

Installing
SteganoGAN

Exporting NN
model to PyTorch
and ONNX

Unusable without fixing
tiny piece of code

*Time*

NN model from
ONNX to Keras

20

# Benchmarking SteganoGAN

**Decoding time as a function of image sizes**



Stealthy decoding CPU implementation implies small images.

**CPU usage as a function of image sizes**



**VRM usage as a function of image sizes**



The footprint of CPU/VRM is not sneaky.

## Size of javascript exploits

|  |  | mean size (B) | min size (B) | max size (B) |
|---|---|---|---|---|
| without compression | | 427.4 | 42 | 3871 |
| (total: 179 exploits) | | | | |
| uglifyjs | compressed | 263.8 | 40 | 2025 |
| (42 exploits) | not compressed | 357.5 | 42 | 2839 |

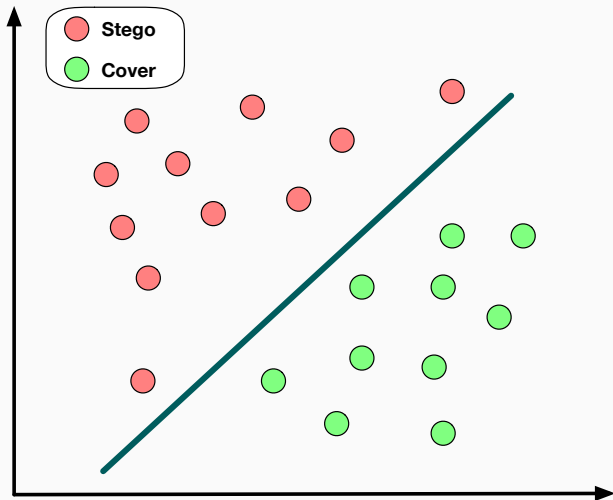- average compression gain: 26.2%
- average compressed files size: 315.4 B

## Stegananalysis (1/2): Hand crafted features sets

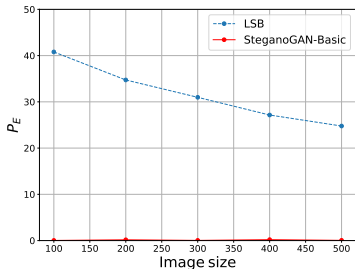|  | Dimensions | Domain |
|---|---|---|
| SRM | 34671 | Spatial |
| SRMQ1 | 12753 | Spatial |
| maxSRM | 34671 | Spatial |
| DCTR | 8000 | JPEG |
| GFR | 17000 | JPEG |
| ... + 20 others | ... | ... |

*Available on dde.binghamton.edu*
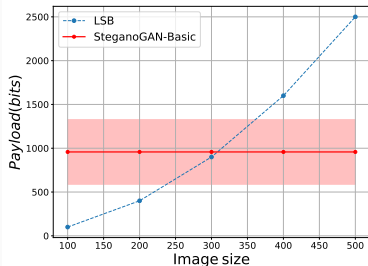
# Steganalysis using trained adversary (SRMQ1 features set)
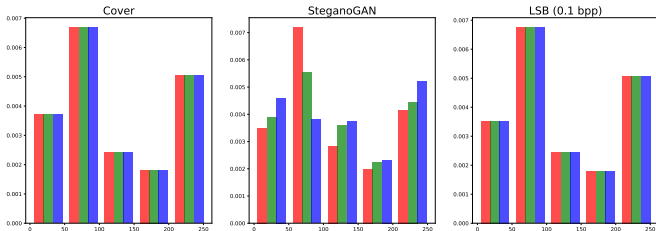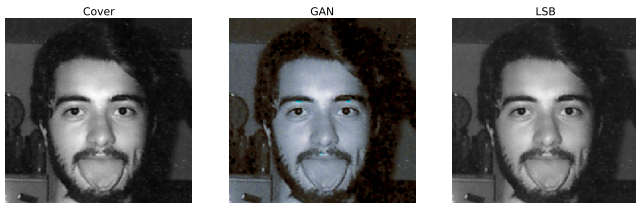
## PE : Probability of error



## Message size



$$P_E = \frac{1}{2}(F_P + M_D)$$

**Bonuses**

# SteganoGAN's weakness (one of them)

## Stegosploit countermeasure (WebMaster POV)

If users are allow JPEG/PNG/... to upload file on the website :

- Place these assets on a separate domain.
- Rewrite the JPEG header to ensure no code is sneaked in there and remove all JPEG comments.
- Refuse requests whose type is "script" and source has a MIME type that starts match an image format

**Literature about Steganography/Steganalysis**

1. Traditional algorithms
   - Spatial-domain
   - Transform-domain
2. Deep learning-based algorithms

## Discussion

- How to achieve a balance between security and capacity?
- How to improve the quality of steganographic image from the ML based large capacity steganography algorithm?
- How to consider complexity?

## Conclusion

**Already done**

- one error away from full **POC**
- some clues about (SteganoGAN)
    - bad performance
    - bad security

**Future works**

- train a proper model for TensorFlow JS
- is steganography relevant for exploits?