

Security of Neural Networks : Attacks, Defenses and Evaluation methods

Rémi BERNHARD (CEA Tech)
Pierre-Alain MOELLIC (CEA Tech)
Jean-Max DUTERTRE (MSE)

*Laboratoire de Sécurité des Architectures et des Systèmes,
Centre CMP, Equipe Commune CEA-Tech Mines Saint-Etienne,
F-13541 Gardanne France*

June 24, 2020



- **Neural networks:** state-of-the art performances in various complex tasks (e.g., image recognition, speech translation)
 - Growing use of neural networks
 - Growing will to deploy models on embedded systems



- **Adversarial machine learning:**
 - Critical decision systems (health, defense and security, ...)
 - Autonomous car
- **Privacy** issues

Serious threats require efficient countermeasures

Security of Machine Learning systems

Security of Machine Learning systems

Threat Model

EXTRACT INFORMATION

Training data (medical, financial, biometric, classified...)

Model (IP, limited authorization)



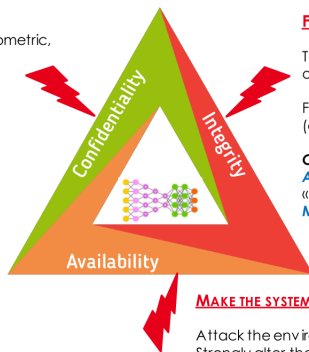
ADVERSARY'S KNOWLEDGE / CAPACITY

Attack at learning / inference time ?

What knowledge about the model ?

→ White box / Black box paradigm

Probing / Querying the model



FOOL A MODEL

The output prediction is not the expected one (i.e. correctly learned)

Fool a model *under the radar*, i.e. in a (almost) imperceptible way

Critical cases:

Autonomous vehicle « Stop » recognizes as « 130 km/h » sign.

Malware detection

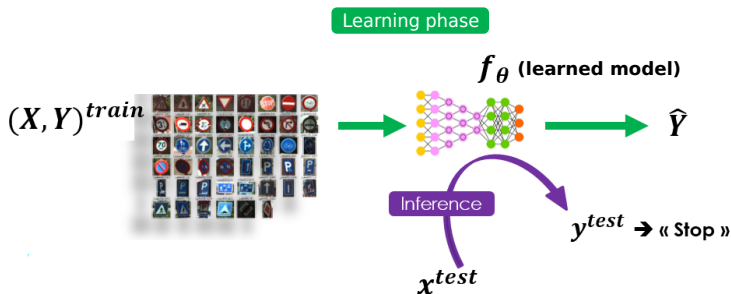
MAKE THE SYSTEM USELESS

Attack the environment (e.g. classical DoS)
Strongly alter the performance of the model

Figure: CIA threat model for a Machine Learning system

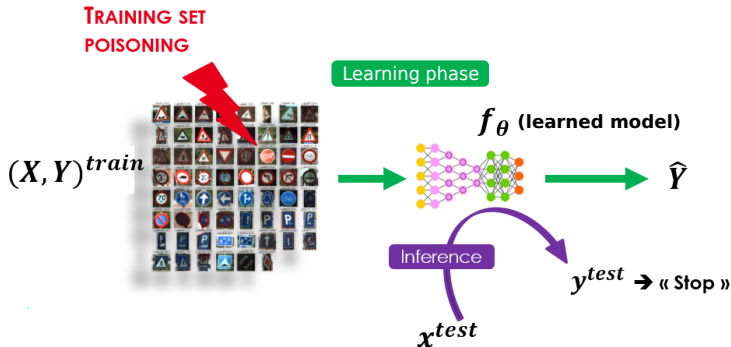
Security of Machine Learning Systems

Striking the ML pipeline



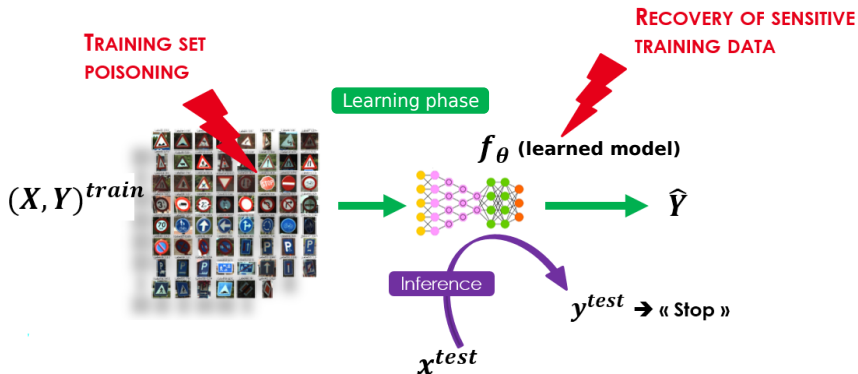
Security of Machine Learning Systems

Striking the ML pipeline



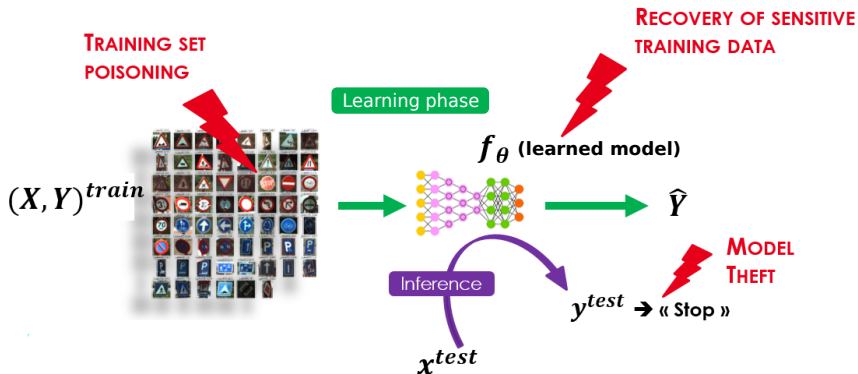
Security of Machine Learning Systems

Striking the ML pipeline



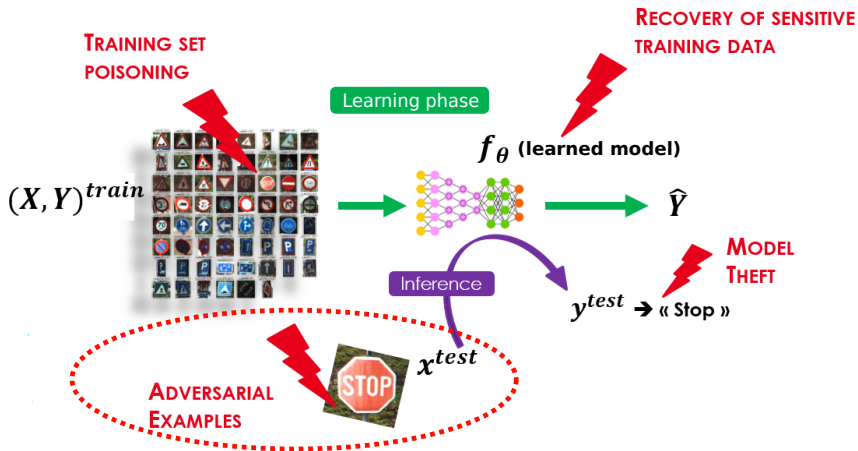
Security of Machine Learning Systems

Striking the ML pipeline



Security of Machine Learning Systems

Striking the ML pipeline



Security of Machine Learning Systems

Adversarial examples

Principle: Craft maliciously modified examples to fool a model.

Adversarial example = Clean example + Adversarial perturbation

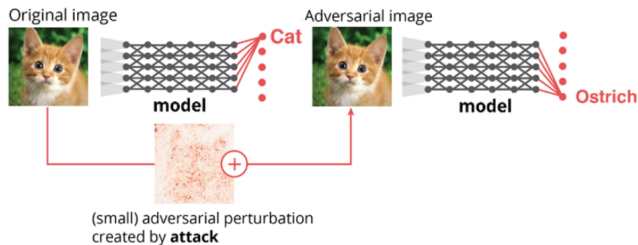


Figure: NIPS 2018 Adversarial Vision Challenge

Security of Machine Learning Systems

Settings and transferability

Threat model:

- **White-box** setting: the adversary has a total access to the target model; He can compute gradients: *gradient-based attacks*
- **Black-box** setting: the adversary has an obstructed access to the target model (score only, label only, etc.). He can't compute gradients: *score/decision based attacks, or transferability*

Principle of transferability:

Adversarial examples crafted on a substitute model transfer to the target model.

→ Powerful tool for an adversary in the black-box setting.

Security of embedded neural networks

Security of embedded neural networks

Quantization methods for embedded systems

Memory footprint:

Parameters storage

Energy cost:

Efficient inference methods

} → *Quantization methods*

Quantization-aware training:

Learn a model with quantized weights and/or activation values during the training process.

Issues: Non-differentiability of quantization functions, difficulty of training, ...

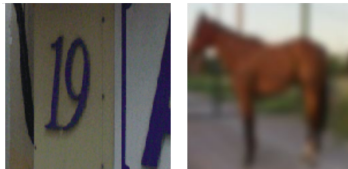
Security of embedded neural networks

Complete study of quantized models vulnerabilities

How does quantization influence robustness against adversarial example ?

Two Data sets:

- SVHN (73,257/26,032)
- CIFAR10 (50,000/10,000)



Experiences:

- Quantization: Activation and Weight / Weight *quantization*: 1,2,3,4 bits
- Techniques: Courbariaux et al. (2015, 2016), Zhou et al. (2016)
- Various threat models considered

Security of embedded neural networks

Results

Results:

- Detection of some gradient masking issue (false impression of security)
→ Quantization is not a robust "*natural*" defense when facing advanced attacks
- But, interestingly, gradient misalignment issues and *quantization shift phenomenon* cause poor transferability
- This enables to build a defense based on an ensemble of quantized models

Best Paper Award at *IEEE Conference on Cyberworlds, 2019*

Luring of adversarial perturbations

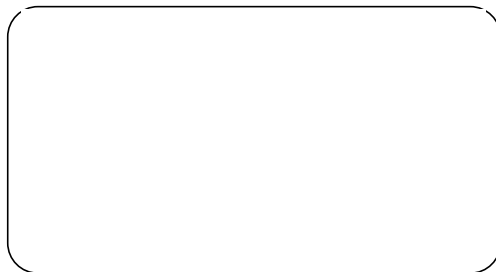
Luring of adversarial perturbations

Context

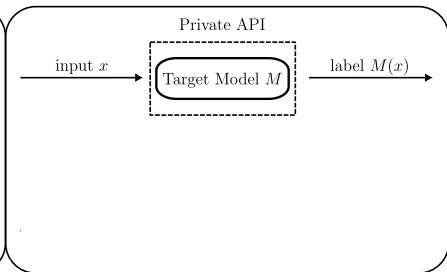
Motivation:

An API provider has trained a **proprietary model** M . He wants to release a **public version** T of this model. *But*, adversarial examples crafted on T should not transfer to M .

Gray-box setting



Black-box setting



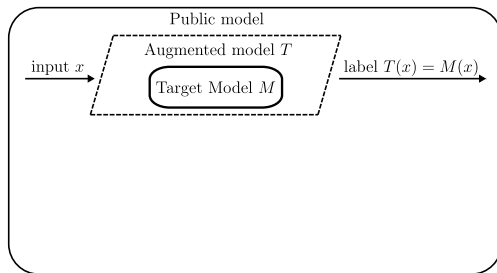
Luring of adversarial perturbations

Context

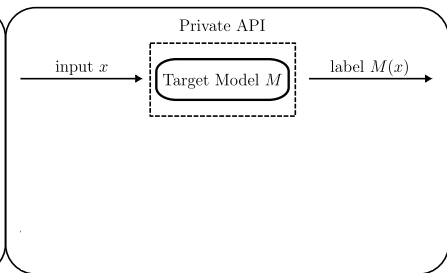
Motivation:

An API provider has trained a **proprietary model** M . He wants to release a **public version** T of this model. *But*, adversarial examples crafted on T should not transfer to M .

Gray-box setting



Black-box setting

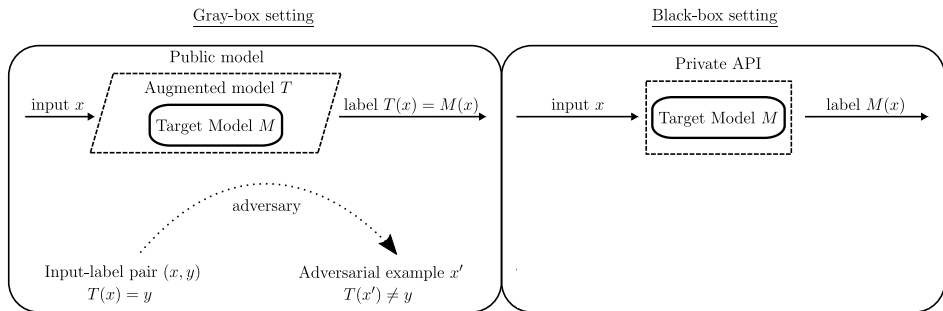


Luring of adversarial perturbations

Context

Motivation:

An API provider has trained a **proprietary model** M . He wants to release a **public version** T of this model. *But*, adversarial examples crafted on T should not transfer to M .



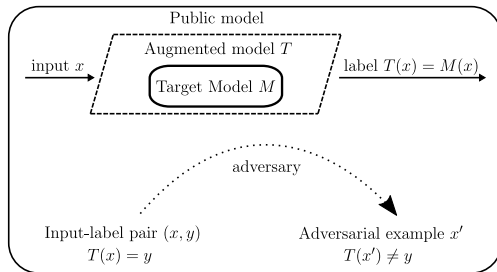
Luring of adversarial perturbations

Context

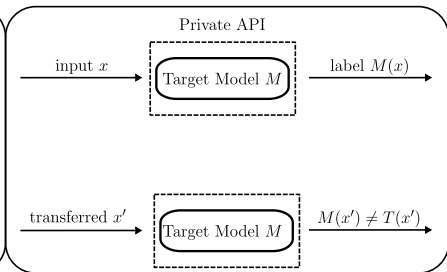
Motivation:

An API provider has trained a **proprietary model** M . He wants to release a **public version** T of this model. *But*, adversarial examples crafted on T should not transfer to M .

Gray-box setting



Black-box setting



Luring of adversarial perturbations

Objective and design

Objective: Lure the adversary

Augment M with a neural network component P to form $T = M \circ P$, so that:

- M and $M \circ P$ agree on clean examples: $M(x) = M \circ P(x)$
- M and $M \circ P$ disagree on adversarial examples: $M(x') \neq M \circ P(x')$

Design:

P is trained so that $M \circ P$ presents **different sensitive features** than M

→ P is designed to fool the adversary (*luring effect*)

→ P is not a preprocessing component aiming at cleaning the adversarial example: it is based on the way M performs prediction.

Luring of adversarial perturbations

Results

Study:

The effectiveness of the method at thwarting an adversary is verified with:

- Three data sets (MNIST, SVHN and CIFAR10)
- State-of-the-art transferability attacks
- Large perturbations allowed for the adversary

Conclusion:

A novel and effective approach to defend against transferred adversarial examples.

Submitted to Usenix Security Symposium 2021

Bio-inspired approach: exploiting frequencies to defend against adversarial examples

Exploiting frequencies against adversarial examples

First results

Objective:

Develop a bio-inspired method to defend against adversarial examples.

Preliminary results:

Take advantage of data sets frequency properties

- Low transferability between models trained on low-pass and high-pass filtered data sets
- Adding frequency specific constraints to the loss function induces non-trivial white-box robustness.

Partnership between the CEA and the university of Grenoble (LPNC)

Timeline and contacts

Planned progress of the Ph.D. :

- *Now - September 2020*: Bio-inspired approach for robustness
- *September 2020 - May 2021*: Link between robustness and vulnerability to M.I.A (Membership Inference Attacks)
- *May 2021 -* :Redaction of the thesis manuscript

Contact:

Secure Architectures and Softwares, SAS

Centre de Microélectronique Provence, Gardanne (13)

Rémi Bernhard: *remi.bernhard@emse.fr*

Pierre-Alain Moellic: *pierre-alain.moellic@cea.fr*

Jean-Max Dutertre : *dutertre@emse.fr*