

## Panorama des Attaques sur Réseaux de Neurones

Quentin Le Roux  
Yannick Teglia  
@ Thales DIS, Cybersecurity Hardware Lab  
& INRIA (partenariat)

Conférence AMUSEC – 25 et 26 Avril 2023

*Inria*



## Thales Digital Identity and Security (DIS)

- Produits et services sécurisés
- Solutions de paiement, Identité & biométrie, Internet des objets (IoT), téléphonie mobile...

## Cybersecurity Hardware Lab

- Garantir la sécurité des produits Thales DIS
- Domaines de compétences:
  - Attaques par canaux auxiliaires (« side-channel »), injection de fautes
  - Attaques micro-architecturales
  - Architectures matérielles sécurisées
  - Cryptographie appliquée
  - Intelligence Artificielle pour l'analyse de sécurité et **sécurité de l'Intelligence Artificielle**

## Thèse en cours avec INRIA, Rennes *Inria*

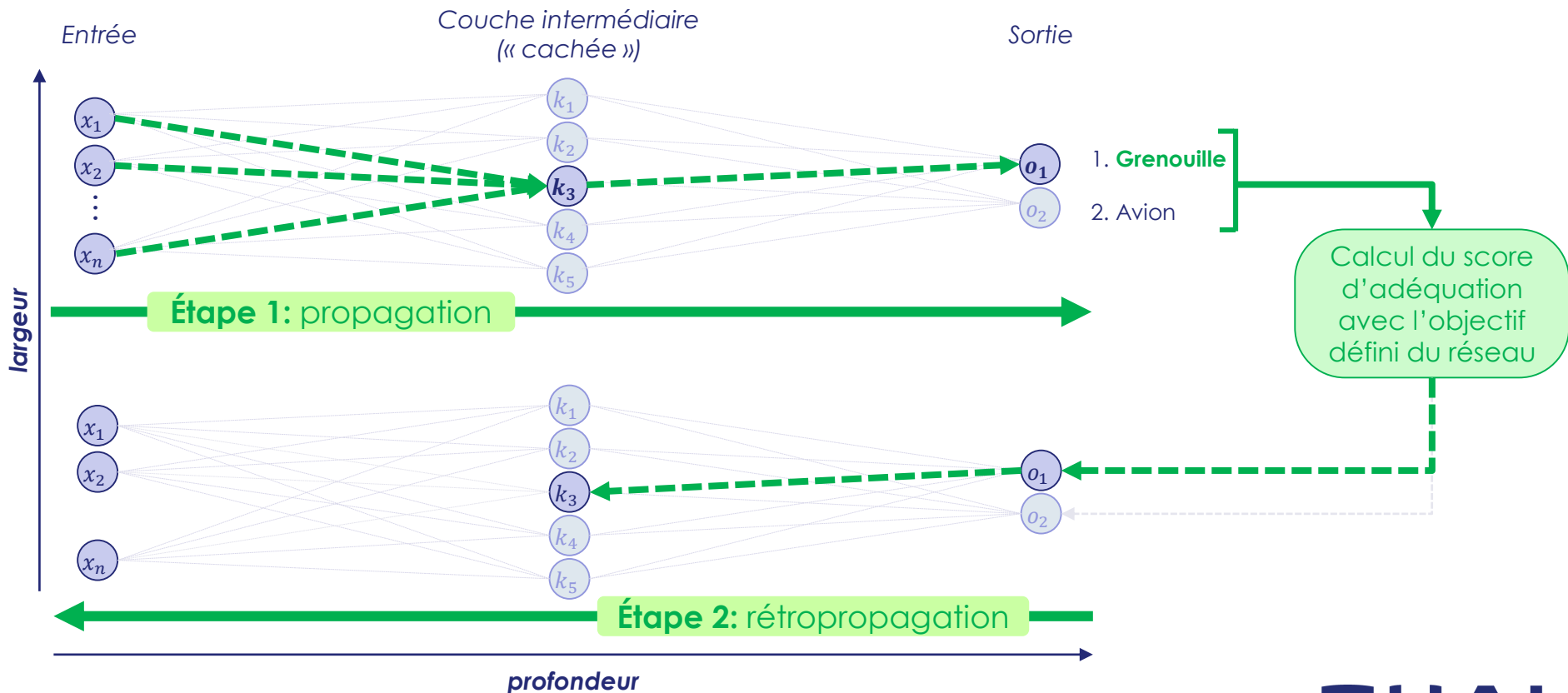
- Intégrité des réseaux de neurones face aux portes dérobées et attaques adverses

# Les Réseaux de Neurons

**Quoi ?** Programme informatique prédictif (classification, détection, etc.)

**Comment ?** Ensemble structuré d'additions, multiplications et de fonctions non-linéaires dont les coefficients sont appris par rétropropagation

**Exemple:**

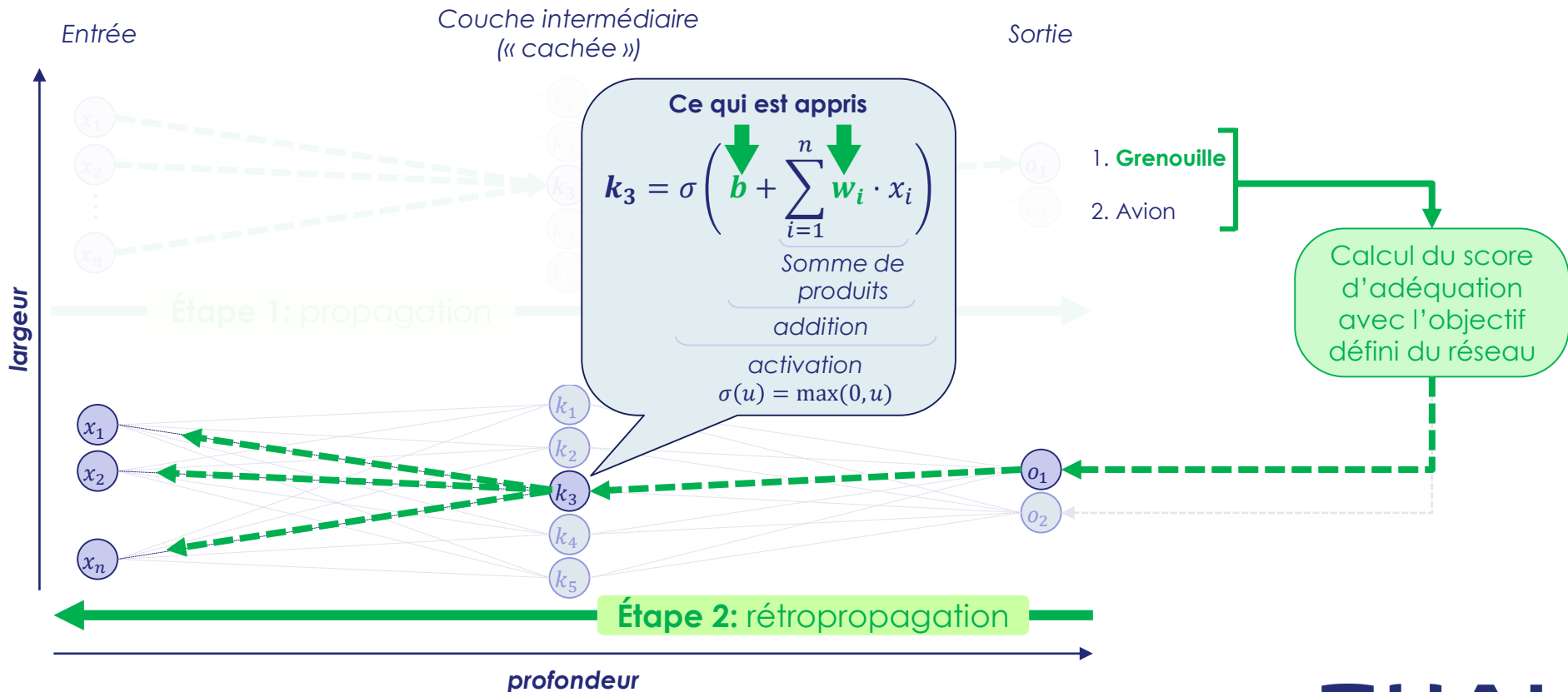


# Les Réseaux de Neurons

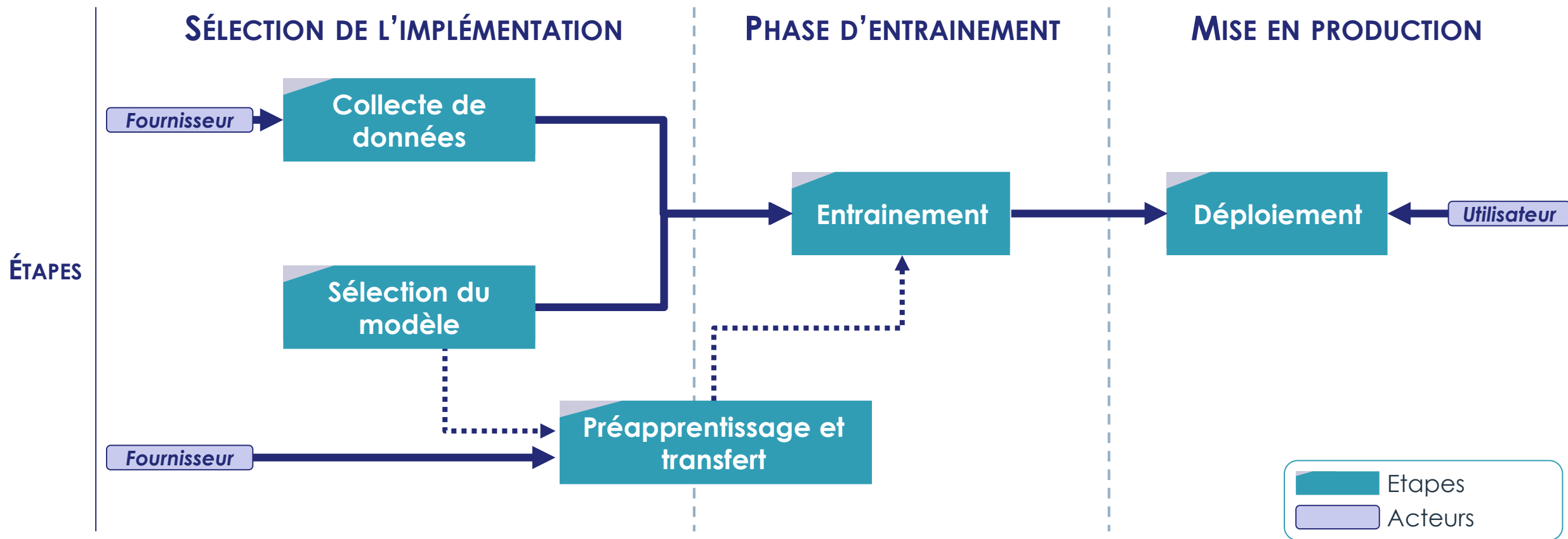
**Quoi ?** Programme informatique prédictif (classification, détection, etc.)

**Comment ?** Ensemble structuré d'additions, multiplications et de fonctions non-linéaires dont les coefficients sont appris par rétropropagation

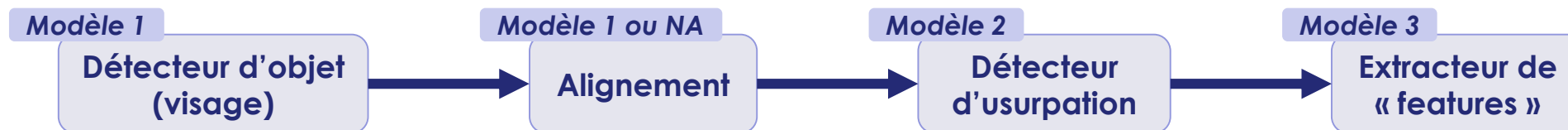
**Exemple:**



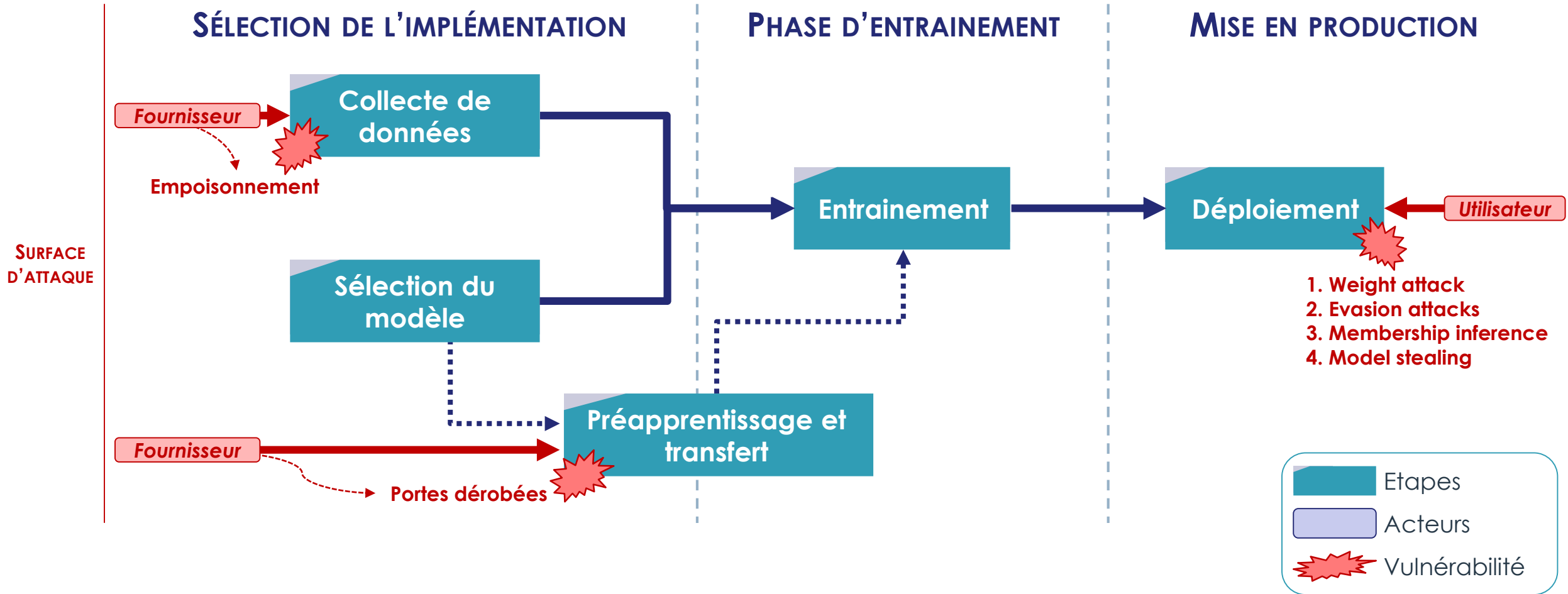
# Cycle de vie d'un réseau de neurones



## PIPELINE EN RECONNAISSANCE FACIALE



# Risques associés au long du cycle de vie (1/3)



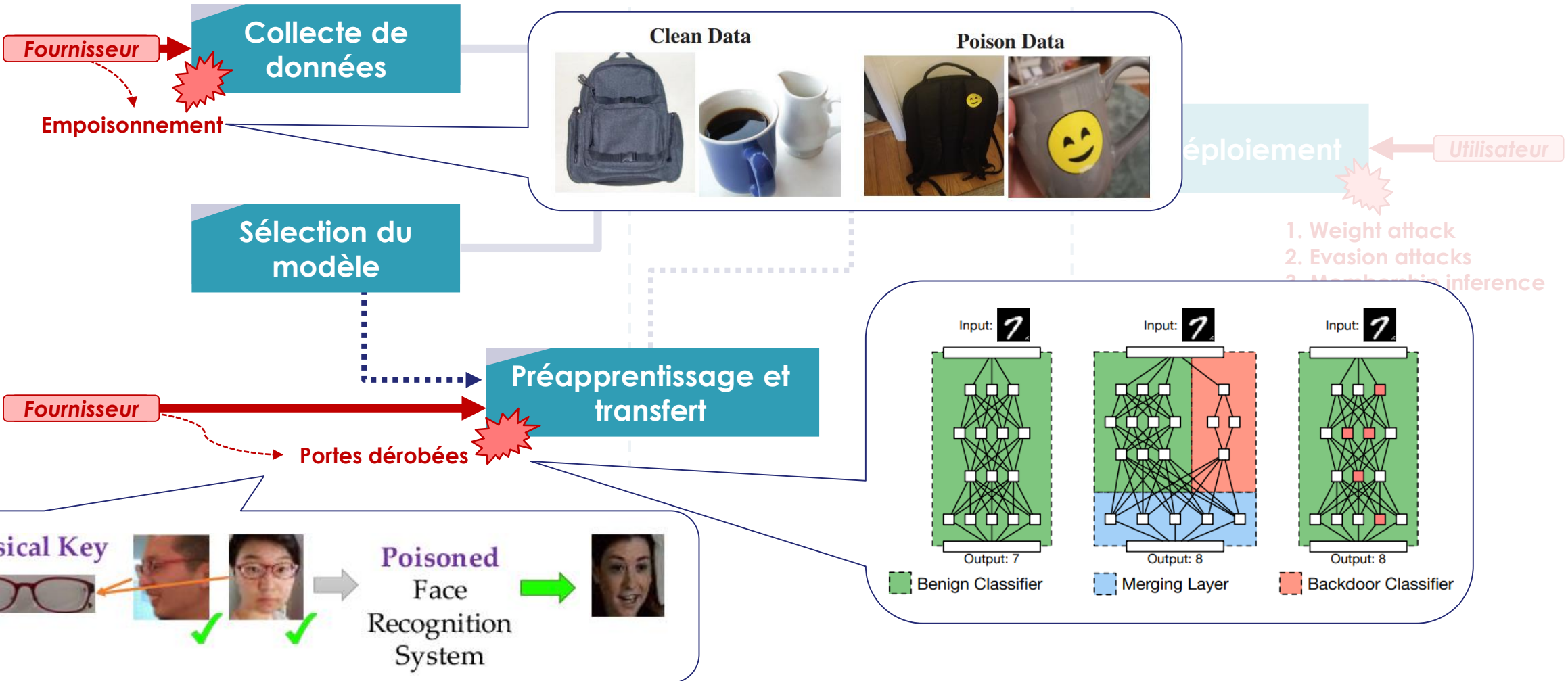
# Risques associés au long du cycle de vie (2/3)

## SÉLECTION DE L'IMPLÉMENTATION

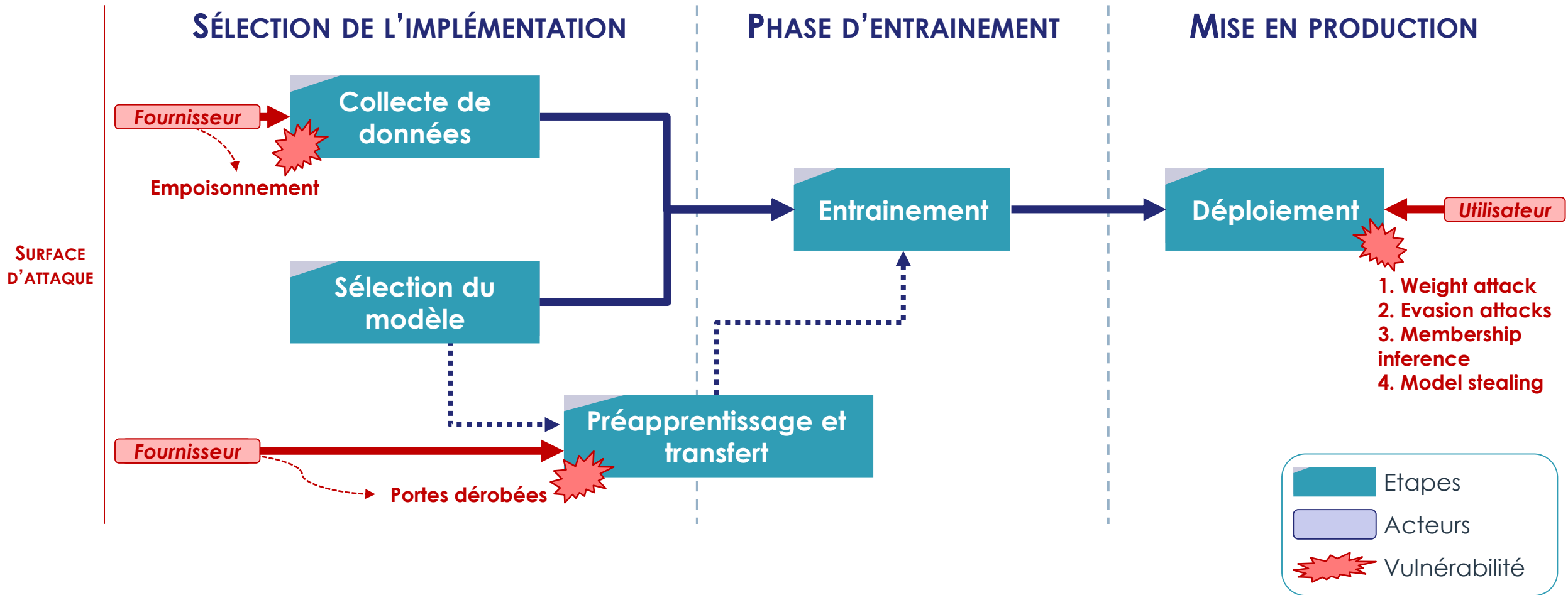
## PHASE D'ENTRAÎNEMENT

## MISE EN PRODUCTION

SURFACE D'ATTAQUE



# Risques associés au long du cycle de vie (1/3)

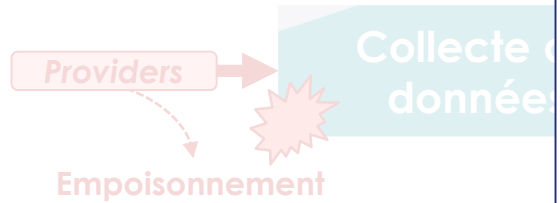




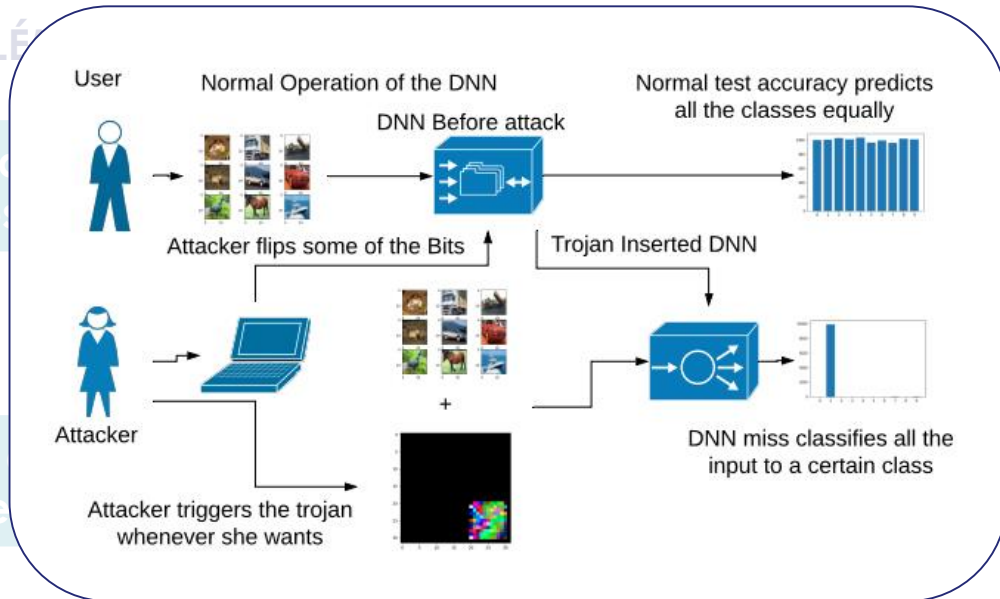
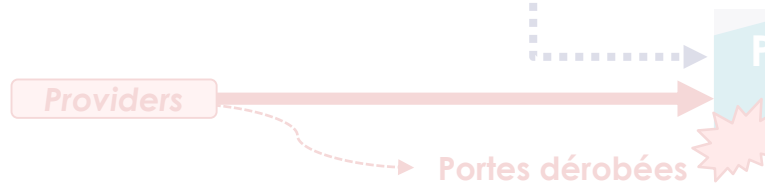
# Risques associés au long du cycle de vie (3/3)

SURFACE D'ATTAQUE

SÉLECTION DE L'IMPLÉ



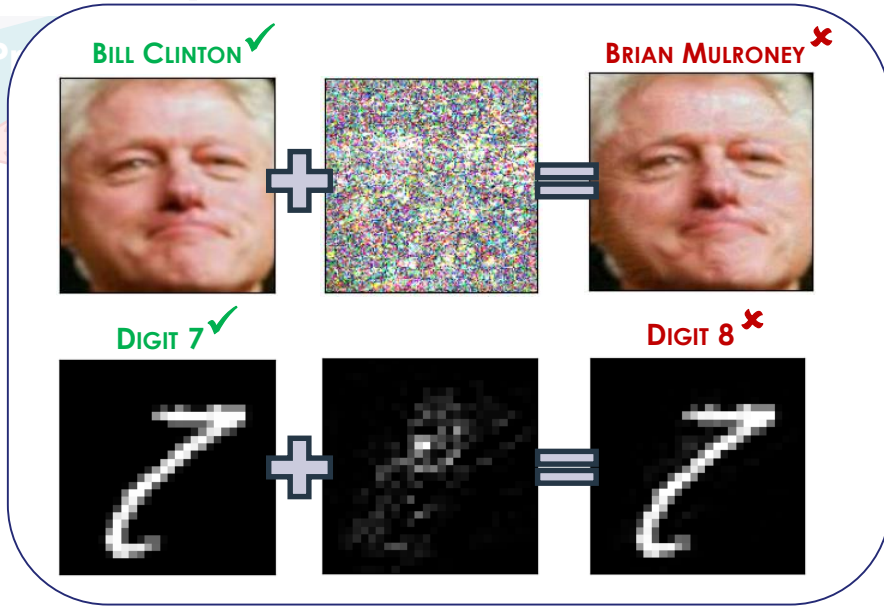
Sélection modèle



MISE EN PRODUCTION



- 1. Weight attack
- 2. Evasion attacks
- 3. Membership inference
- 4. Model stealing



# Exemple: construire une attaque adverse



Basé sur le cadre de la formulation de l' « Empirical Risk Minimization »:

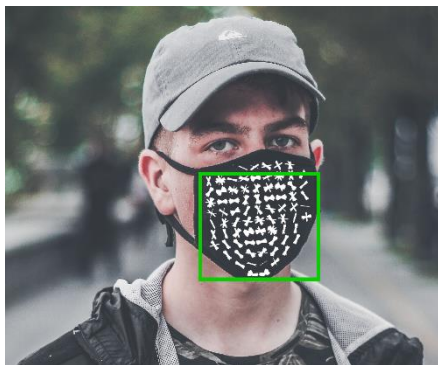
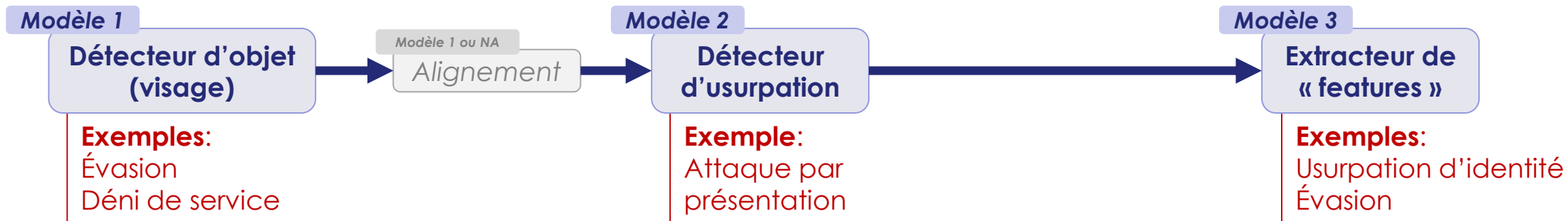
$$\theta^* = \arg \min_{\theta \in \Theta} \mathbb{E}_{\{x,y\} \sim D} [L(f_{\theta}(x), y; \theta)] \quad \text{(ERM)}$$

$$\delta = \arg \min_{z \in \Delta} f(x + z) = y' \neq f(x) \quad \text{(méthode avec classe ciblée)}$$

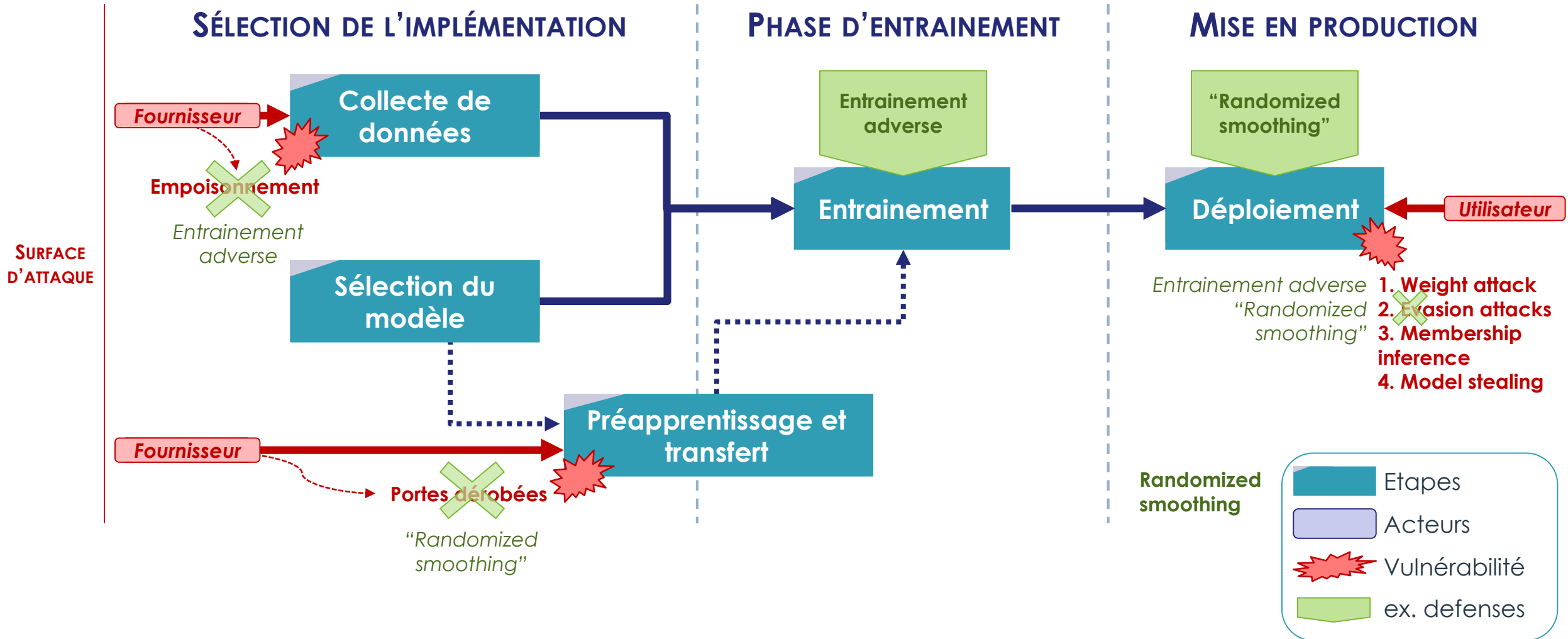
Un agent malicieux peut construire un exemple adverse avec une multitude de méthodes (boîtes noires, boîtes blanches, méthodes apprises, algorithmes génétiques, etc.)

# Cycle de vie d'un réseau de neurones

## PIPELINE EN RECONNAISSANCE FACIALE



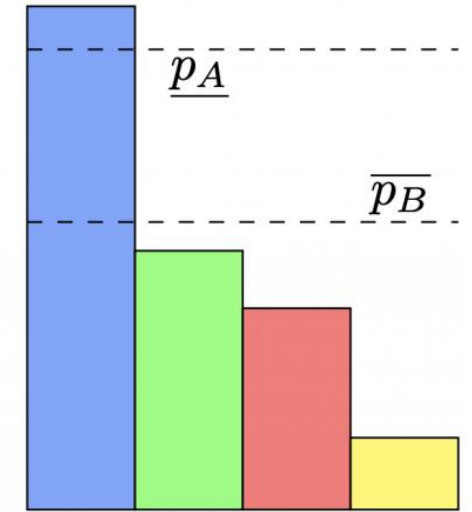
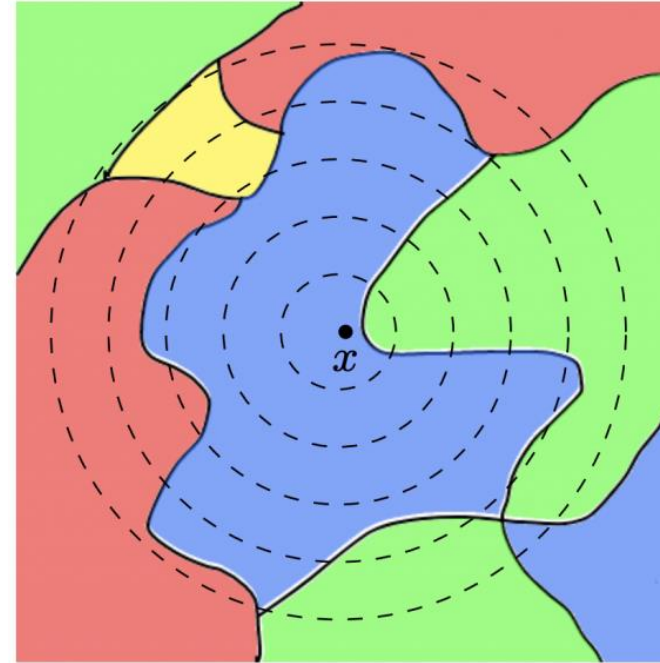
# Risques associés au long du cycle de vie (1/3)



# “Randomized Smoothing,” une défense en phase de test

Un modèle  $g$  est construit à partir d'un réseau  $f$

- $g$  duplique une entrée  $x$  donnée et applique à chaque duplicata un bruit donné (ex. Gaussien)
- La classe avec la valeur de prédiction la plus élevée en moyenne (ex. probits) est celle prédite par le modèle
- La méthode **certifie** un intervalle de confiance autour de l'entrée  $x$



L'espace autour de l'entrée  $x$  est intégrée selon une notion de distance et, sachant un encadrement estimé vis-à-vis de la probabilité de voir la seconde plus haute classe, le modèle  $g$  retourne la classe la plus prédite

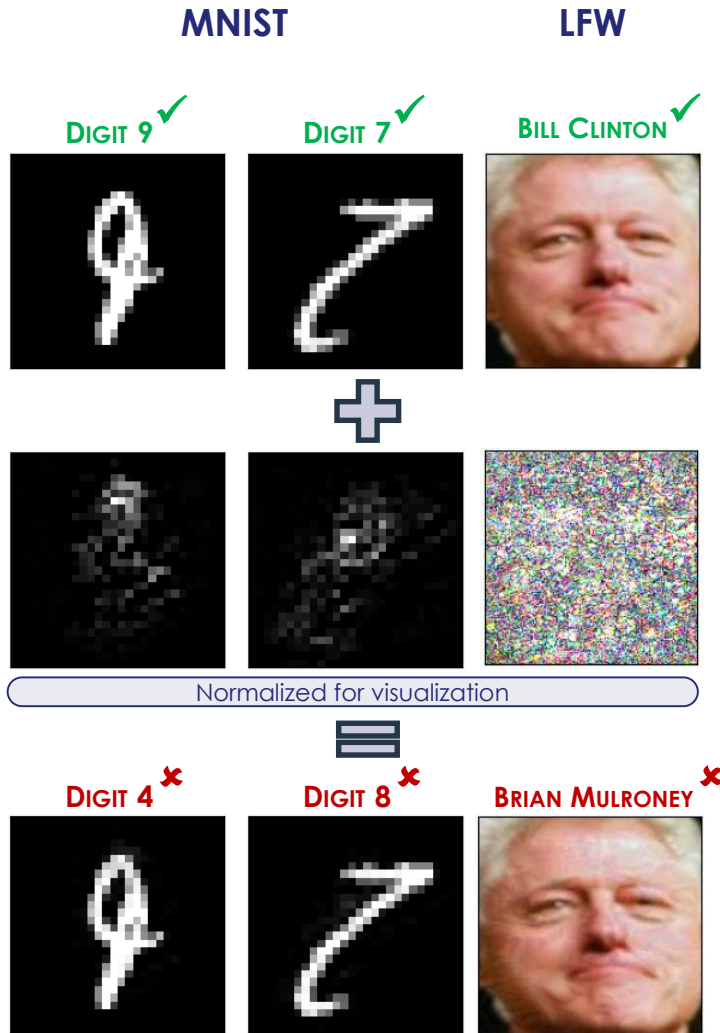
# Conclusions et perspectives

- **Un réseau de neurones est un modèle versatile mais vulnérable**
- **Le cycle de vie entier du réseau est susceptible à des manipulations par des agents malicieux**
- **Comprendre les différents vecteurs d'attaque est primordial dans le cadre d'applications critiques**
- **Présentation centrée sur l'intégrité des réseaux de neurones**
- **D'autres menaces existent:**
  - **Confidentialité, disponibilité, respect de la vie privée**
- **Nombreux travaux en cours chez Thales**

# Annexe

# Attack Examples

## Adversarial examples

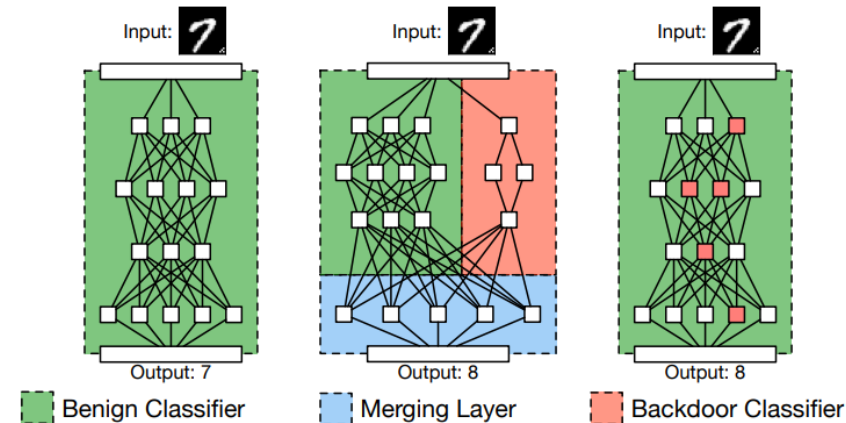


## Model backdoor

"BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain," Gu et al. (2019)



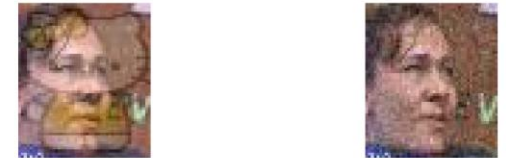
Figure 8. Real-life example of a backdoored stop sign near the authors' office. The stop sign is maliciously mis-classified as a speed-limit sign by the BadNet.



## Poisoned data

### Digital trigger

"Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning," Chen et al. (2017)



(a) An image blended with the Hello Kitty pattern. (b) An image blended with the random pattern.

Fig. 6: Poisoning instances blended with different patterns. In both images, the blended ratio  $\alpha = 0.2$ .

### Physical trigger

"Backdoor Attacks Against Deep Learning Systems in the Physical World," Wenger et al. (2021)

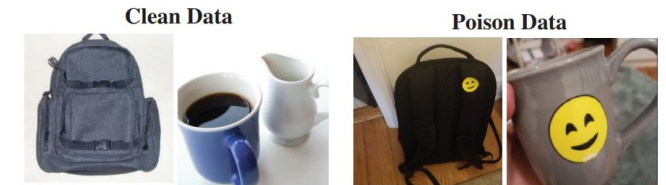
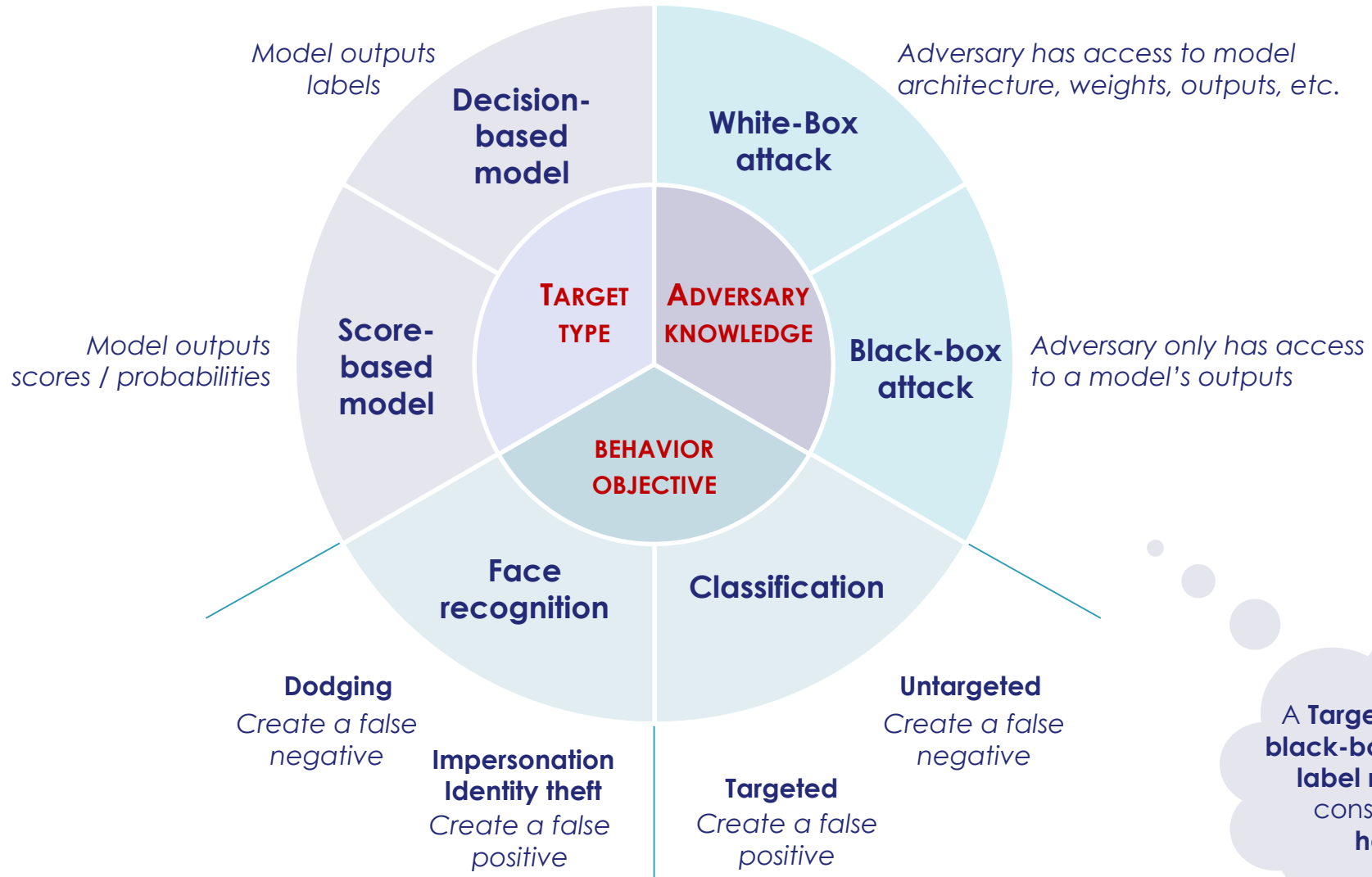


Figure 11: Examples of clean and poison data used in the object recognition experiments

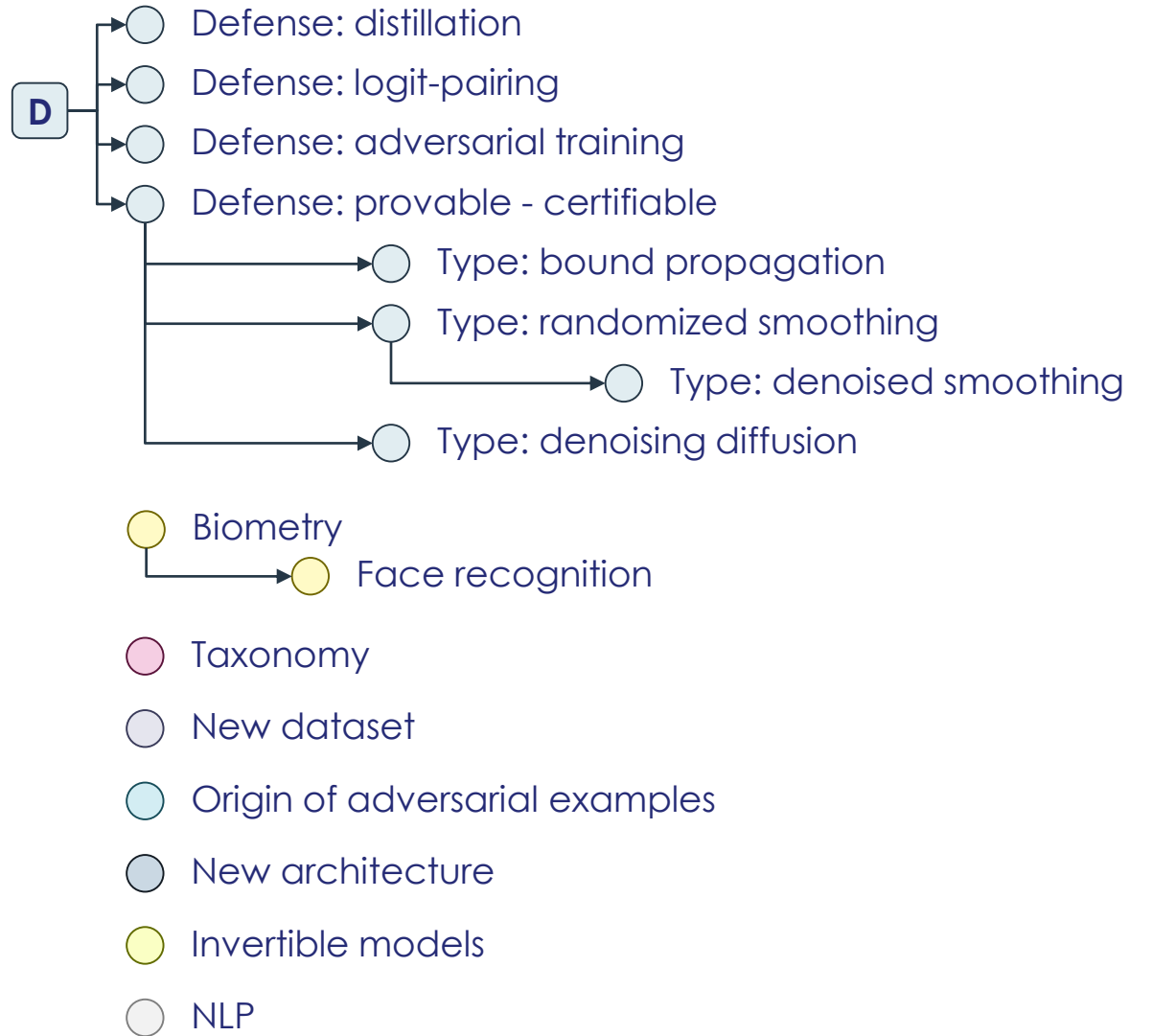
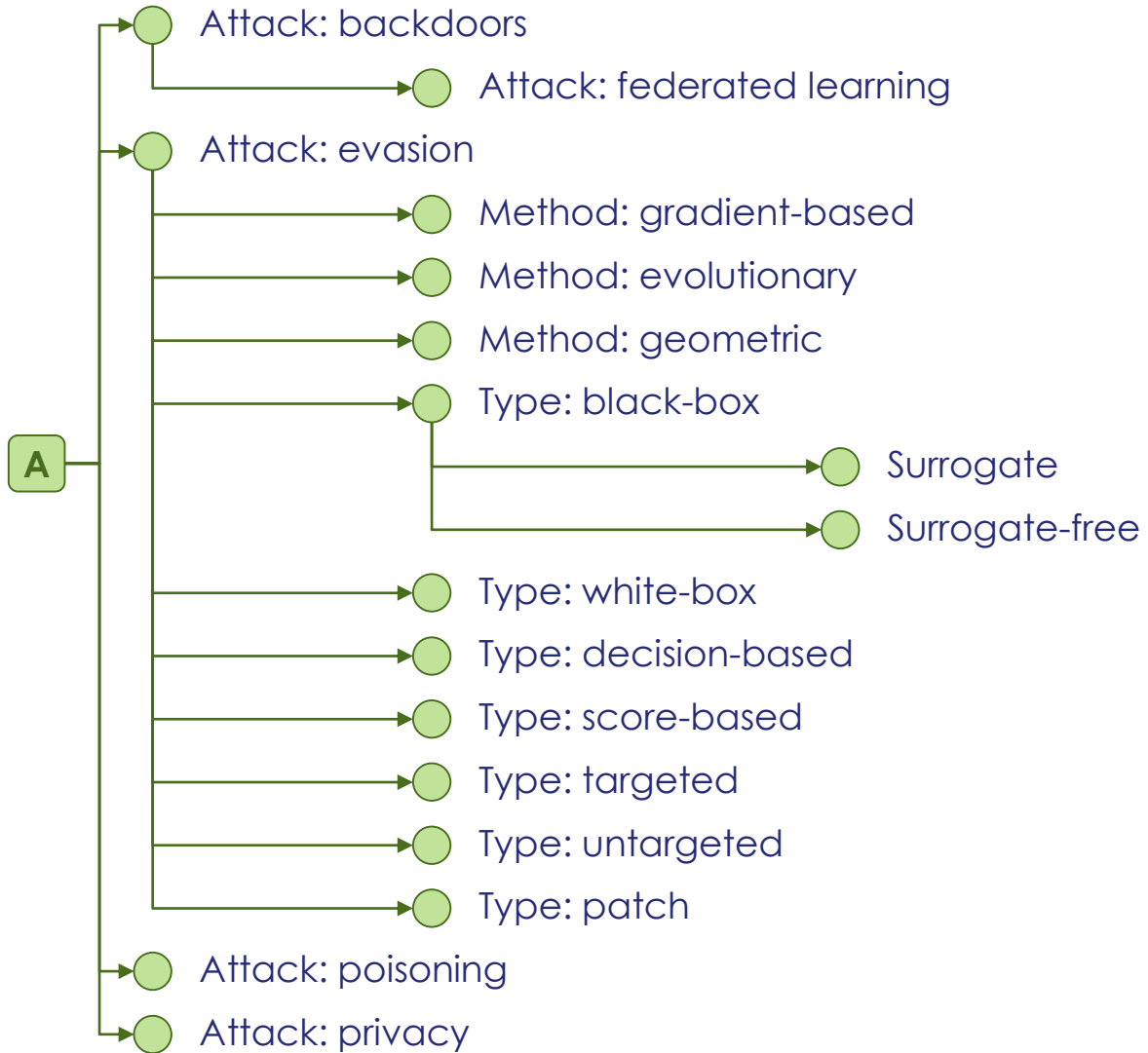


# Threat model (2/2)

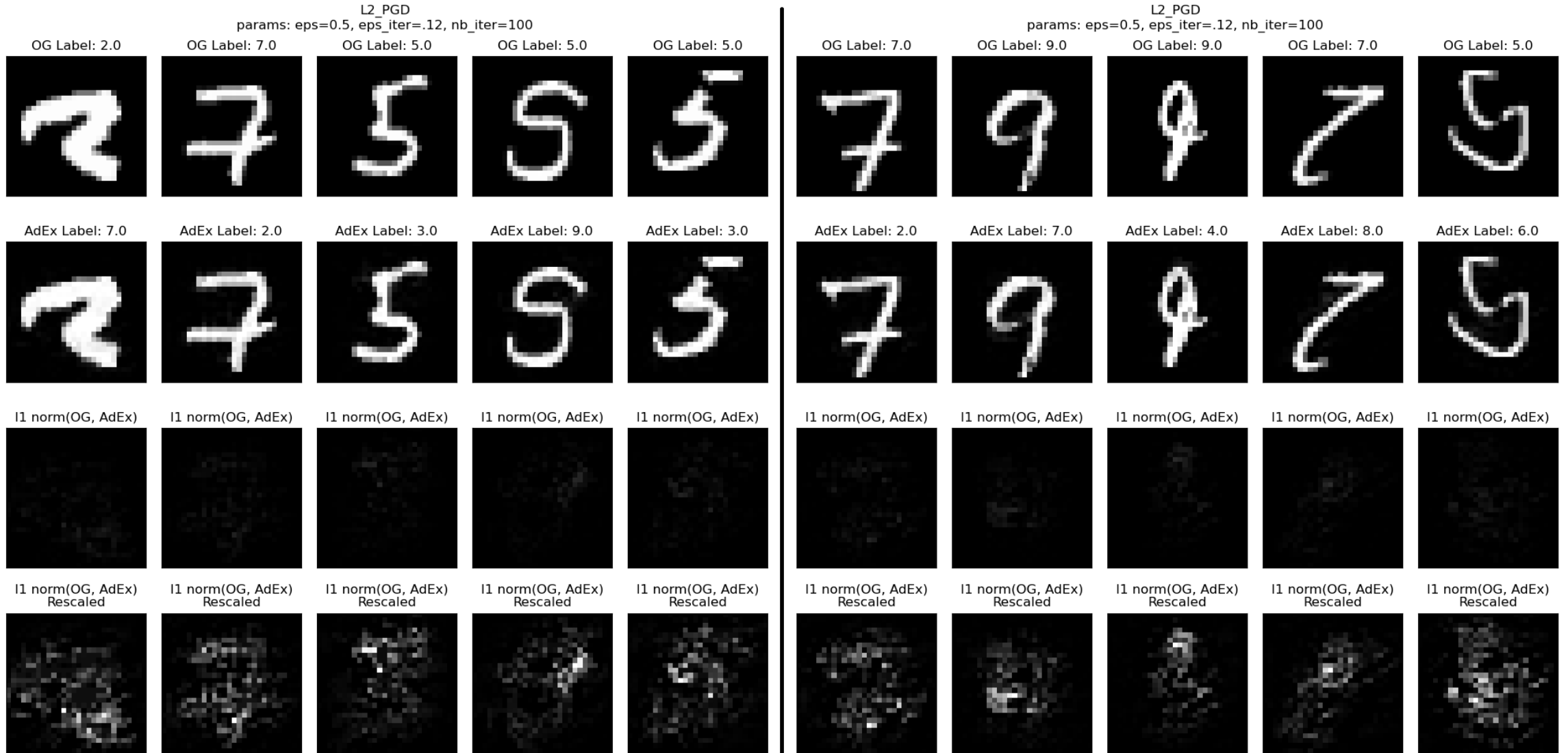


A Targeted/impersonation black-box attack on a top-1 label model is generally considered to be the hardest setting

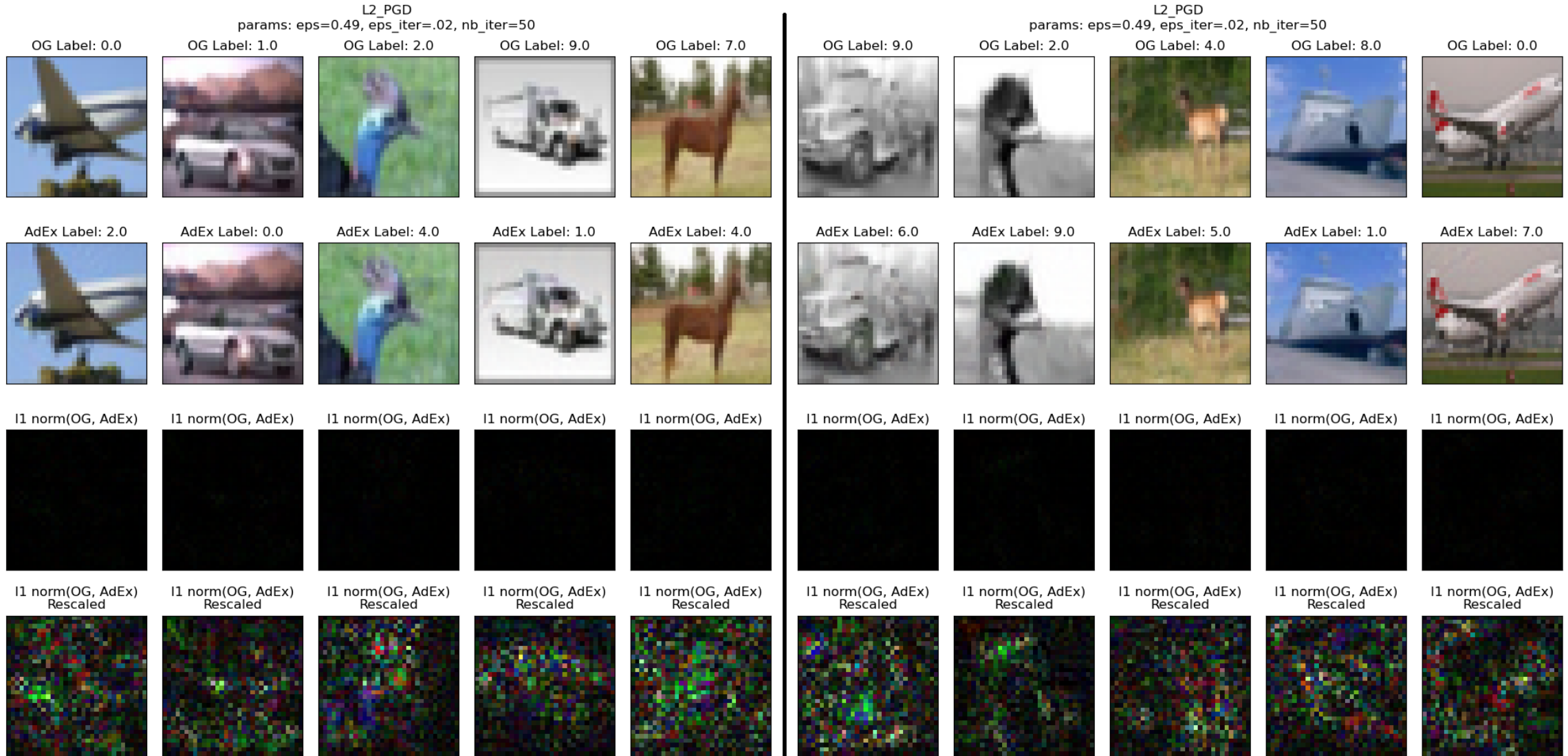
# Nomenclature



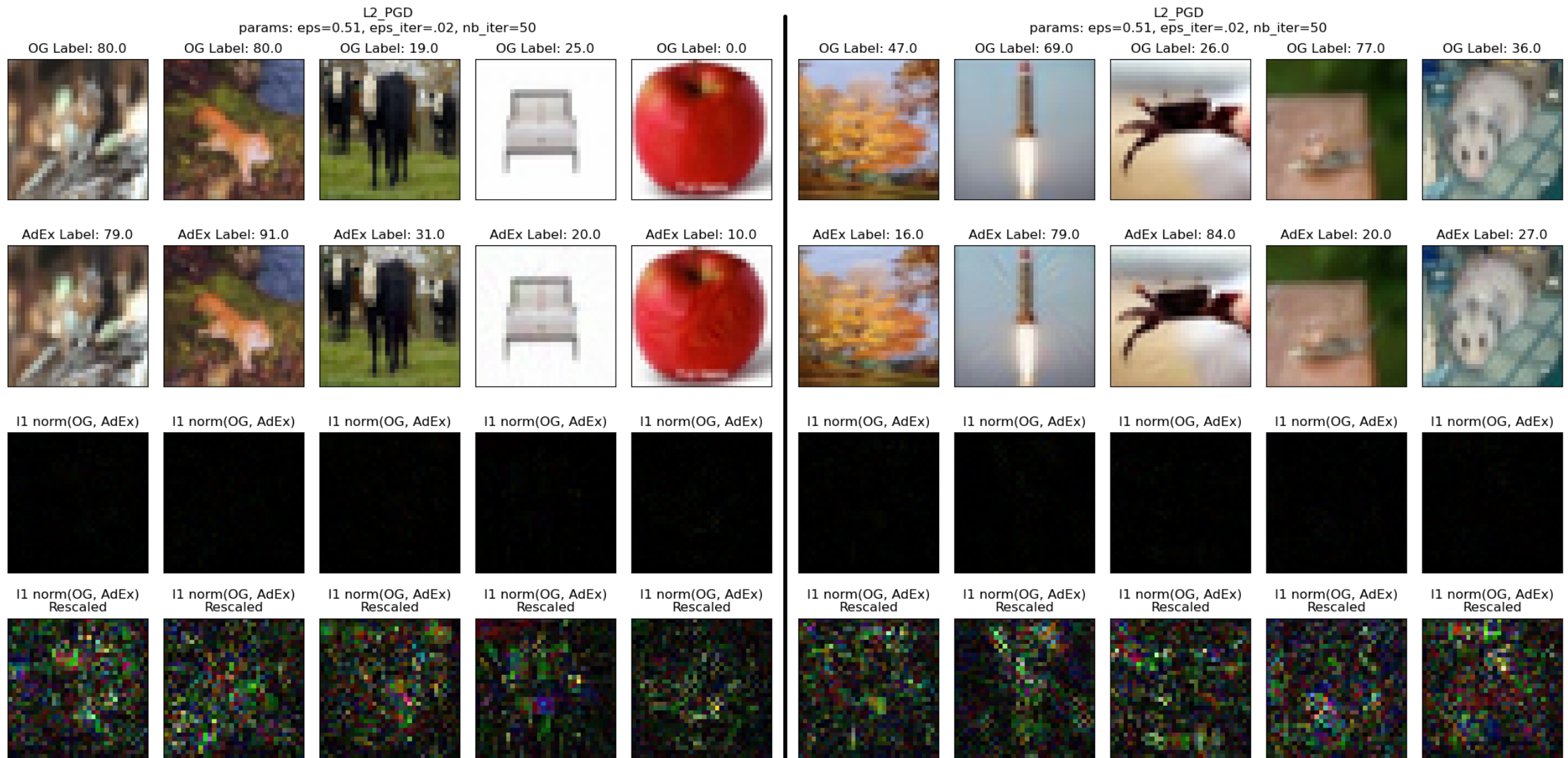
# What is an evasion/adversarial attack?



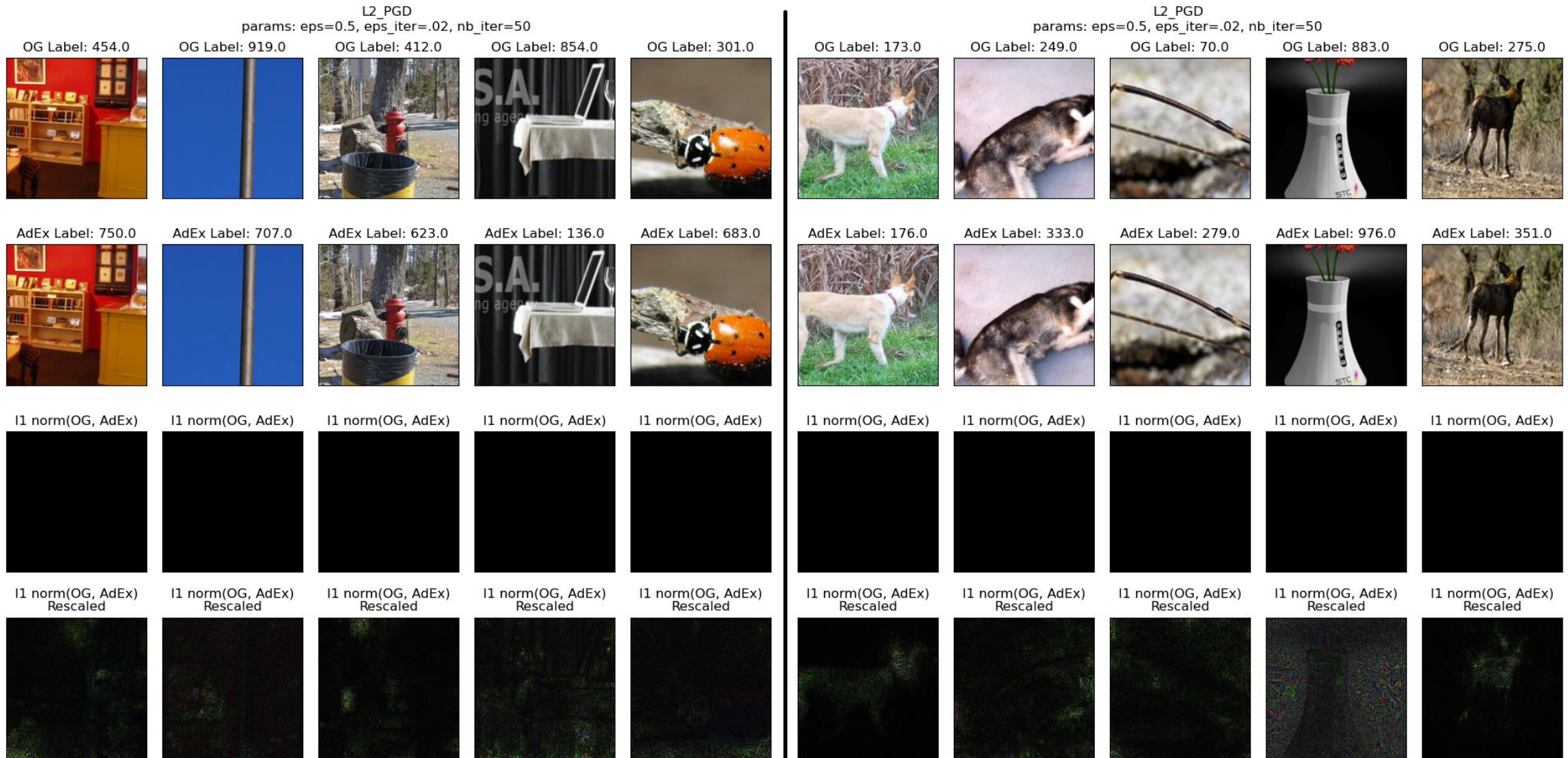
# What is an evasion/adversarial attack?



# What is an evasion/adversarial attack?



# What is an evasion/adversarial attack?



# What is an evasion/adversarial attack?

L2\_PGD  
params: eps=2.5, eps\_iter=.02, nb\_iter=50

OG Label: 115.0

OG Label: 50.0

OG Label: 311.0

OG Label: 419.0

OG Label: 318.0



AdEx Label: 198.0

AdEx Label: 13.0

AdEx Label: 318.0

AdEx Label: 455.0

AdEx Label: 374.0



l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)



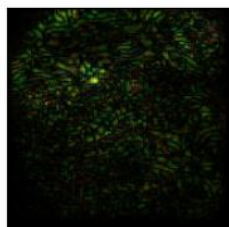
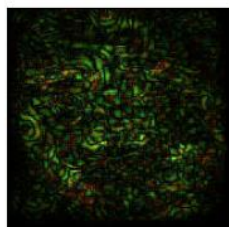
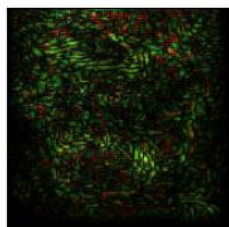
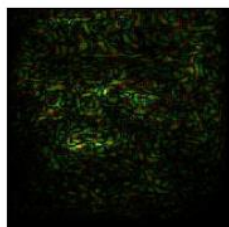
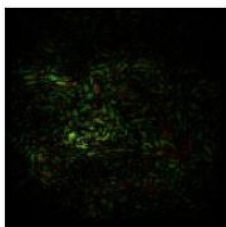
l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled



L2\_PGD  
params: eps=2.125, eps\_iter=.02, nb\_iter=50

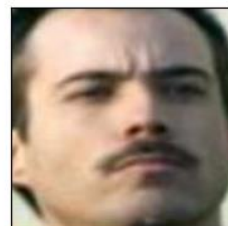
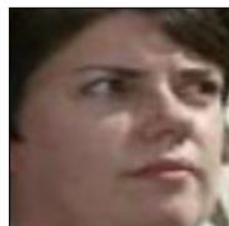
OG Label: 253.0

OG Label: 279.0

OG Label: 133.0

OG Label: 248.0

OG Label: 113.0



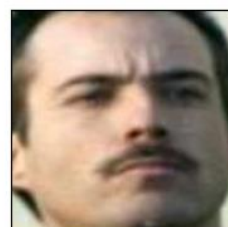
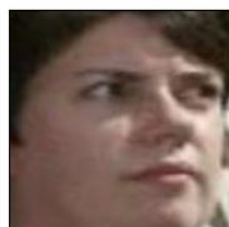
AdEx Label: 255.0

AdEx Label: 226.0

AdEx Label: 198.0

AdEx Label: 250.0

AdEx Label: 186.0



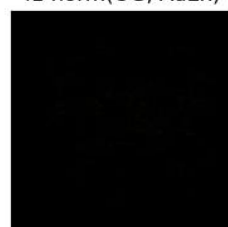
l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)

l1 norm(OG, AdEx)



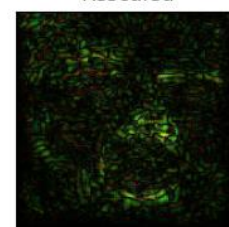
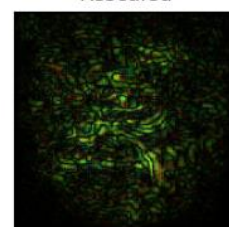
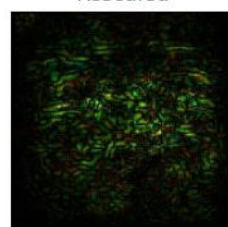
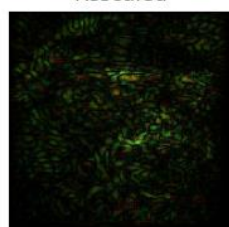
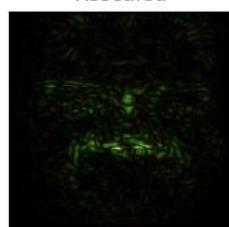
l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled

l1 norm(OG, AdEx)  
Rescaled



# Adversarial Training (training-time defense)

## Modification of a neural network's SGD update step

- The loss is computed not only on the input data  $\mathbf{X}$  but also its noisy/adversarial version  $\mathbf{X}_{\text{adv}}$
- $\mathbf{X}_{\text{adv}}$  is sampled with a noise distribution or an iterative generation at each epoch

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\text{loss}(f_{\theta}(x), y)] \quad (\text{general case})$$

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{\delta \in \Delta} \text{loss}(f_{\theta}(x + \delta), y)] \quad (\text{worst-case performance})$$

Source: Frost, "Certifiable Robustness to Adversarial Attacks" (2020)