



RÉPUBLIQUE
FRANÇAISE

*Liberté
Égalité
Fraternité*

ONERA

THE FRENCH AEROSPACE LAB

Empoisonnement de données

23/09/2022

Adrien CHAN-HON-TONG

ONERA



RÉPUBLIQUE
FRANÇAISE

*Liberté
Égalité
Fraternité*

ONERA

THE FRENCH AEROSPACE LAB

Ce document est la propriété de l'ONERA. Il ne peut être communiqué à des tiers et/ou reproduit sans l'autorisation préalable écrite de l'ONERA, et son contenu ne peut être divulgué.
This document and the information contained herein is proprietary information of ONERA and shall not be disclosed or reproduced without the prior authorization of ONERA.

IA de confiance à l'ONERA



Objectif

→ l'objectif de cette présentation est de présenter une cartographie des risques d'empoisonnement mais aussi des contre-mesures

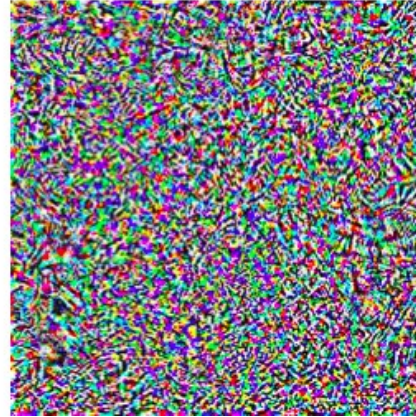
Exemple 1

La dangerosité des exemples adversaires invisibles ?

“pig”



+ 0.005 x



=

“airliner”

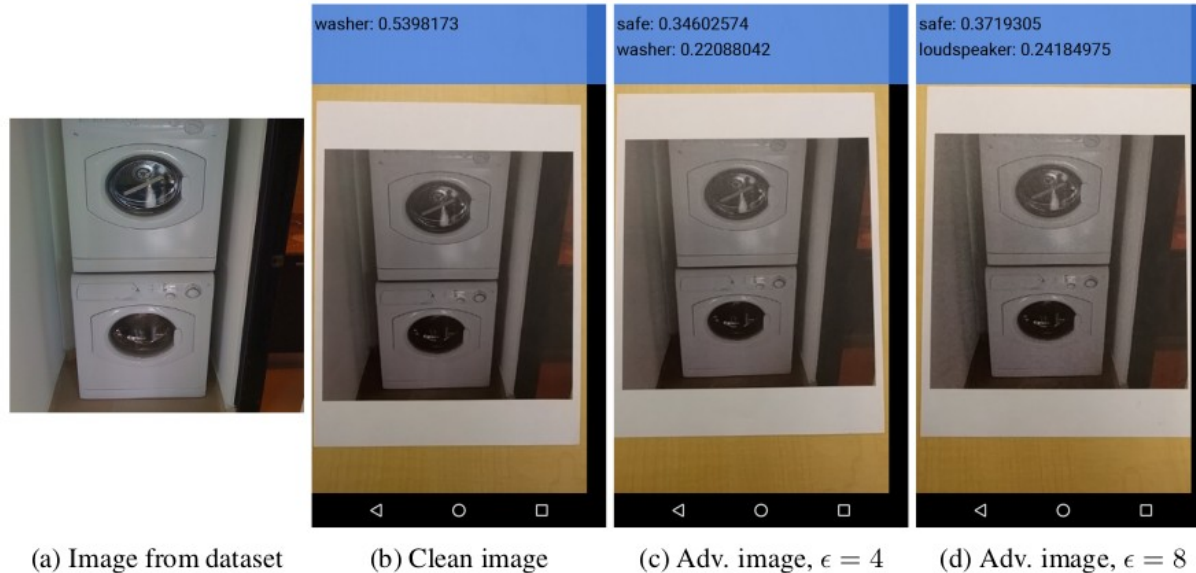


Exemple 1

La dangerosité des exemples adversaires invisibles ?

1 → Ce type d'attaque semble difficilement réalisable dans le monde physique

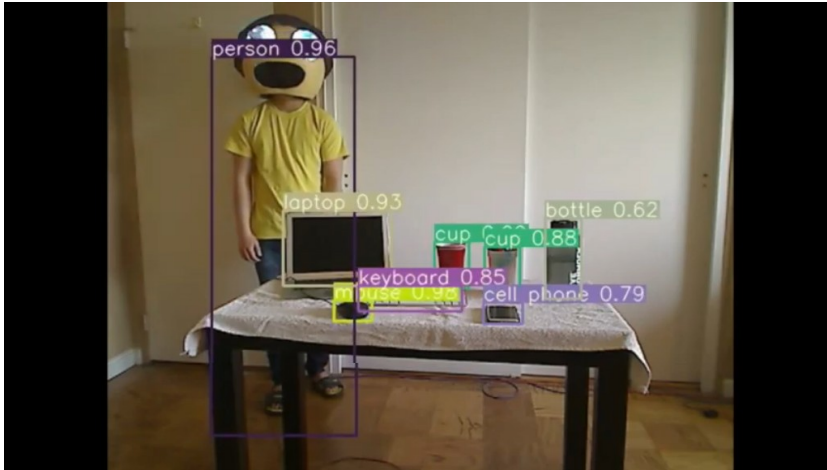
malgré



2 → Un réentraînement avec PGD donne une robustesse non négligeable

Exemple 2

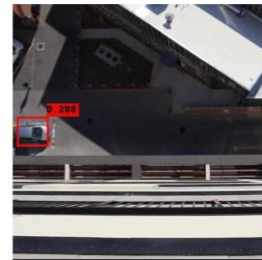
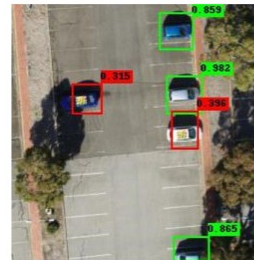
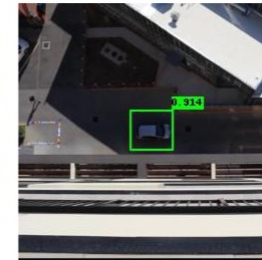
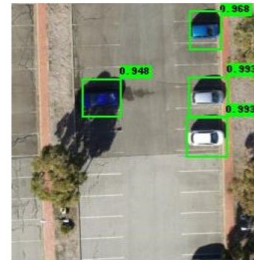
La dangerosité des exemples adversaires patches ?



(a) Adversarial patch on the roof of a car.



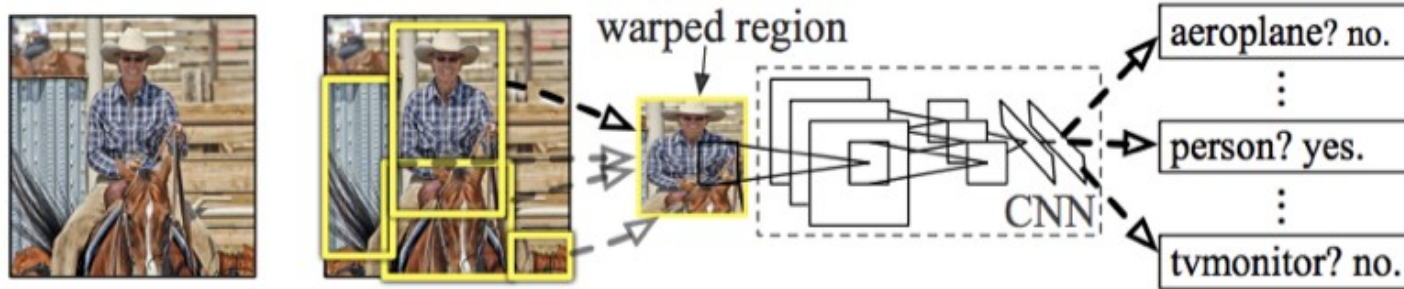
(b) Adversarial patch off-and-around a car.



Exemple 2

La dangerosité des exemples adversaires patches ?

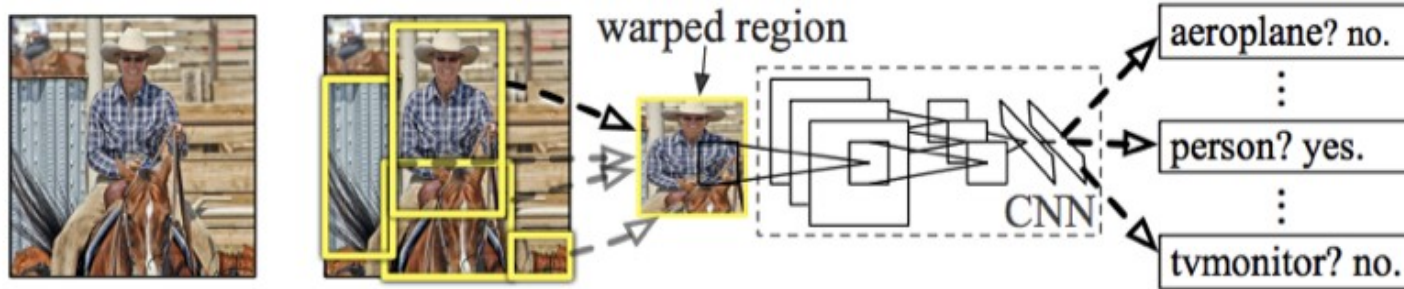
R-CNN: *Regions with CNN features*



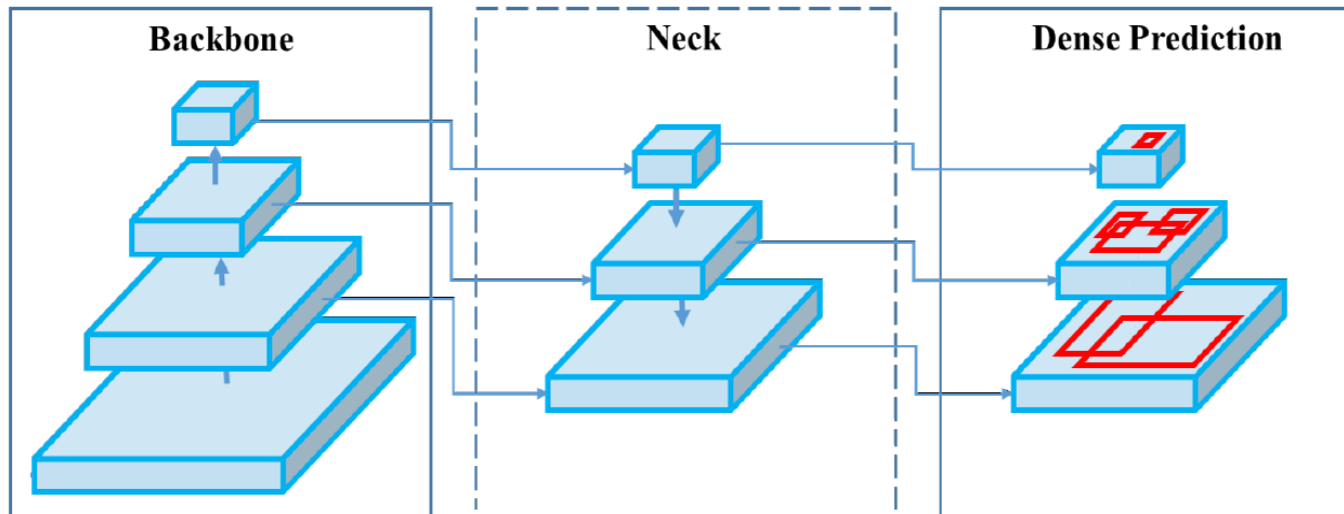
Exemple 2

La dangerosité des exemples adversaires patches ?

R-CNN: *Regions with CNN features*

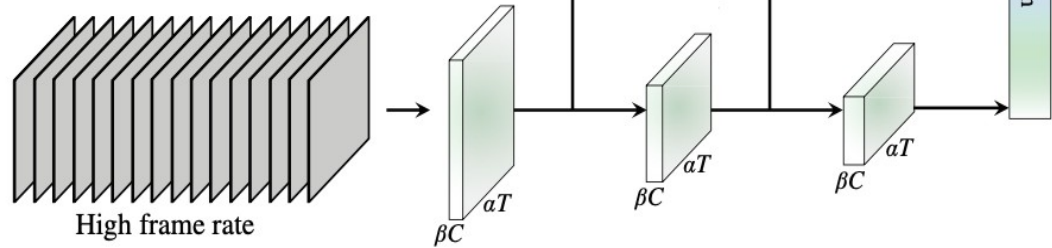
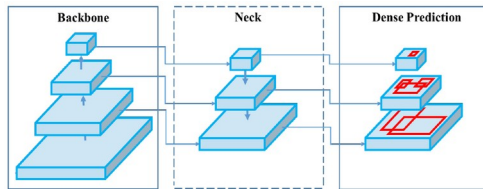
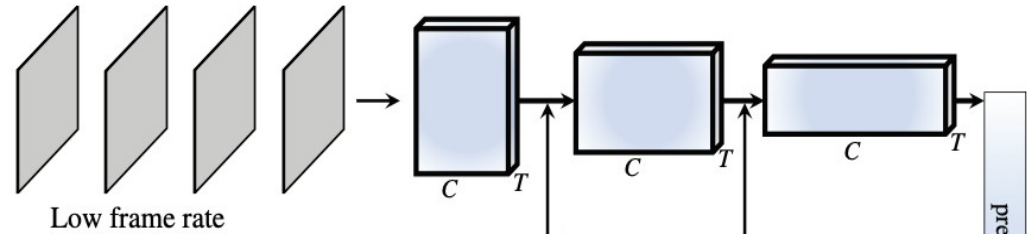
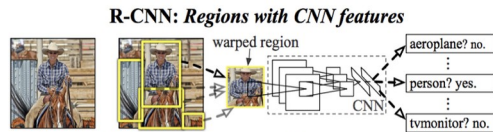


medium.com



Exemple 2

La dangerosité des exemples adversaires patches ?

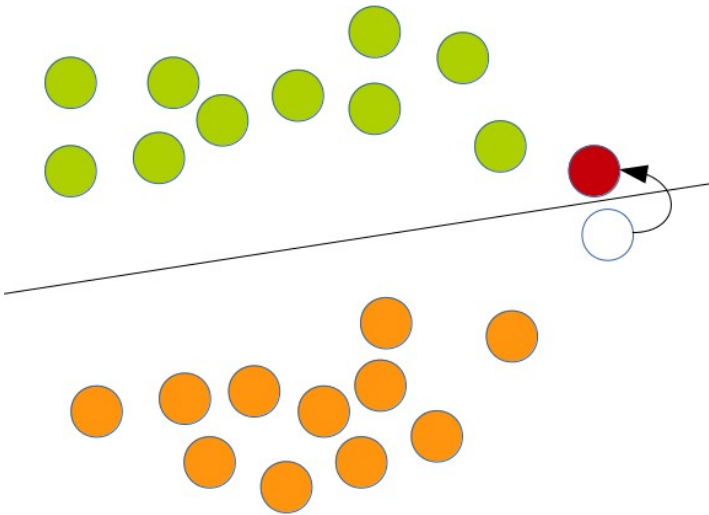


Plan

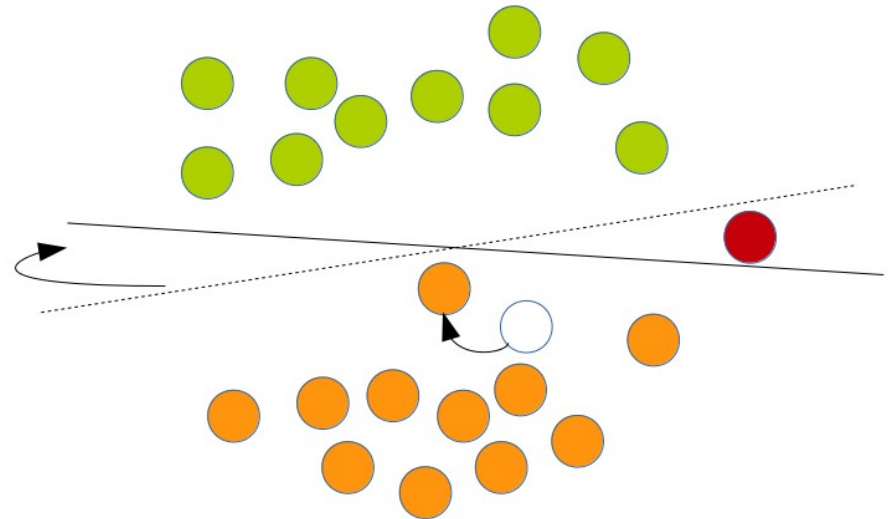
- scénario d'empoisonnement
- contre-mesure à une empoisonnement invisible
- perspectives

empoisonnement vs attaque évasion

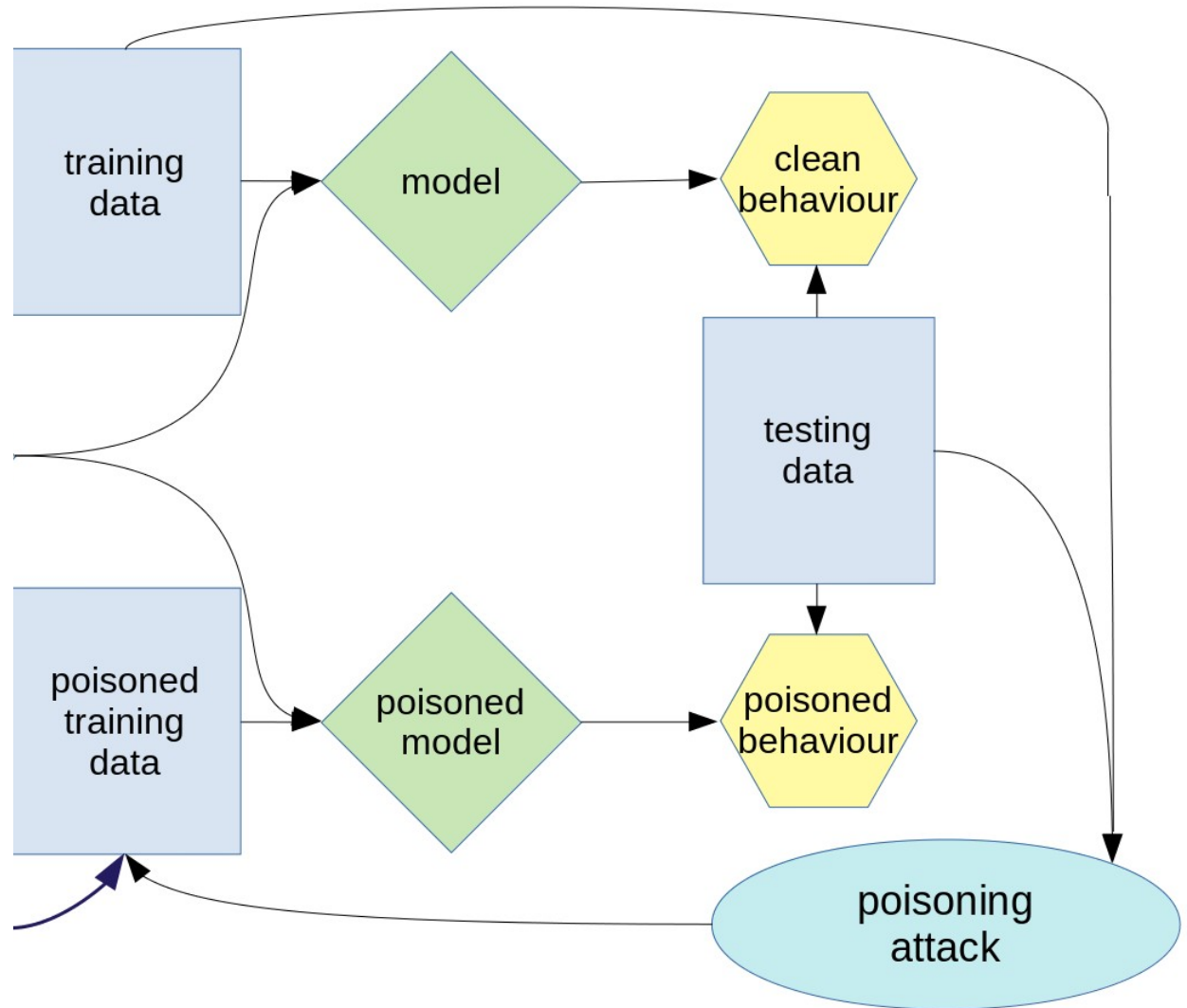
Attaque adversaire



Empoisonnement



empoisonnement

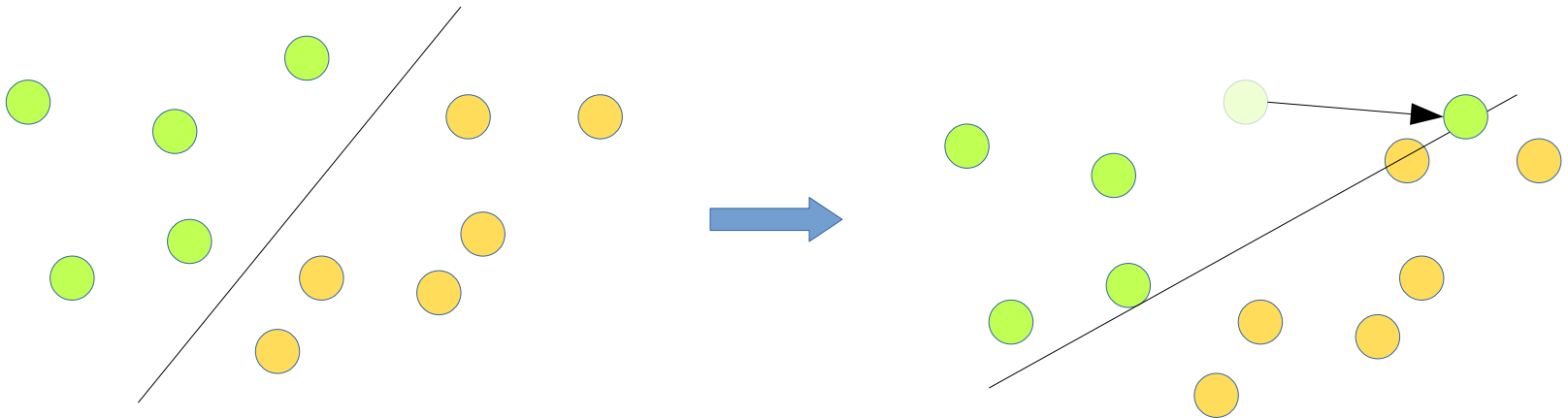
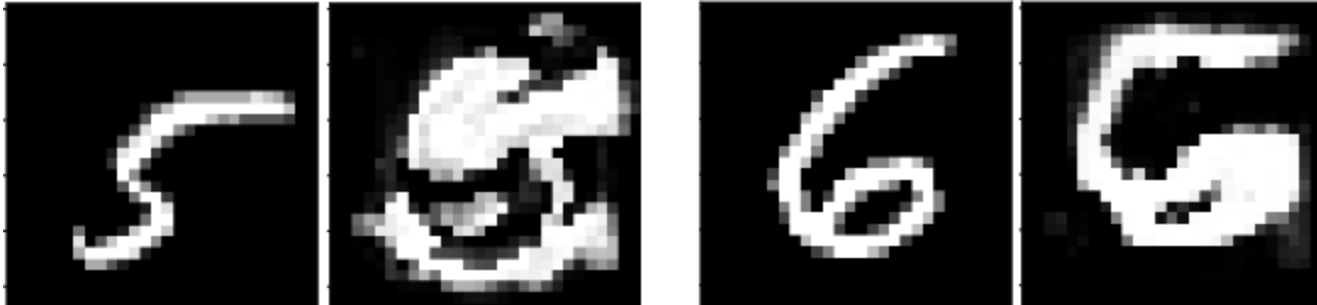


empoisonnement vs attaque évasion

	Attaque adverse	Empoisonnement
Avantages	Facilité d'exécution	Modifie le comportement de la cible partout
Inconvénients	Modification du comportement locale	Nécessite un accès aux données d'apprentissage.

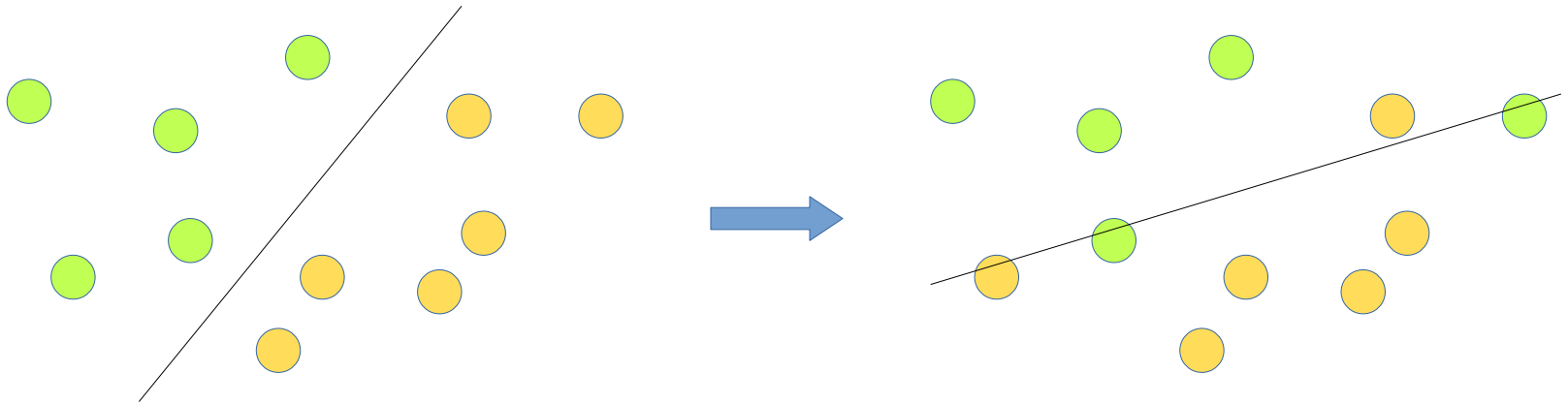
empoisonnement classique vs empoisonnement invisible

Image corruption



empoisonnement classique vs empoisonnement invisible

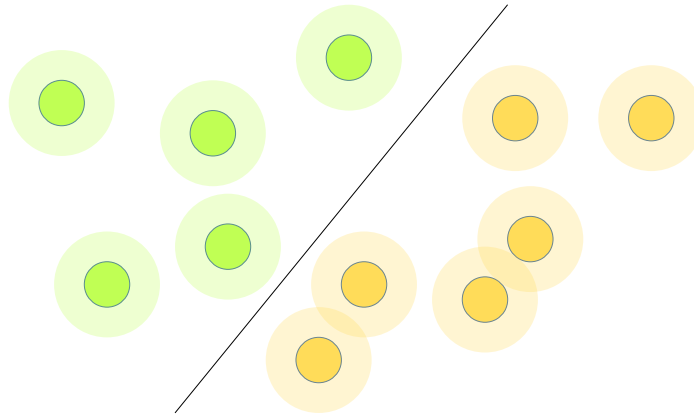
Label flip



empoisonnement classique vs empoisonnement invisible

Clean label

Modification des images contraintes à être indétectable à l'oeil



Scénarios d'empoisonnement

Soit l'apprentissage est fait sans contrôle

Soit la modification doit passer inaperçue

Scénarios d'empoisonnement

Soit l'apprentissage est fait sans contrôle

Ce qui paraît pas très réaliste

Soit la modification doit passer inaperçue

Scénarios d'empoisonnement

Soit l'apprentissage est fait sans contrôle

Ce qui paraît pas très réaliste

(Sauf pseudo labeling, online learning, watermarking, falsification)

Soit la modification doit passer inaperçue

Scénarios d'empoisonnement

Soit l'apprentissage est fait sans contrôle

Ce qui paraît pas très réaliste

(Sauf pseudo labeling, online learning, watermarking, falsification)

Soit la modification doit passer inaperçue

→ faible nombre de modification

→ OU modification invisible

Scénarios d'empoisonnement



Modifications fortes mais peu nombreuses



Modifications invisibles

Modifications fortes mais minoritaires vs partition

Training data

Partitionnement

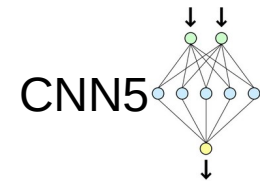
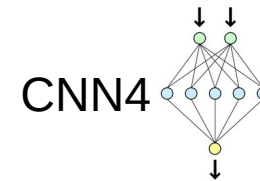
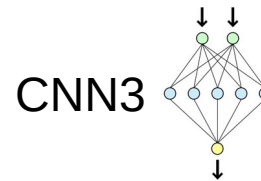
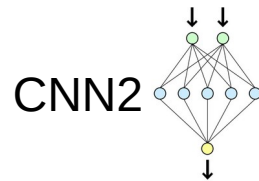
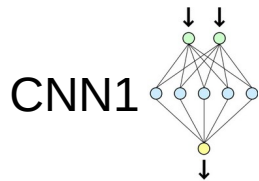
P1

P2

P3

P4

P5



Modifications fortes mais minoritaires vs partition

Training data

Partitionnement

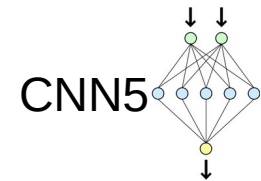
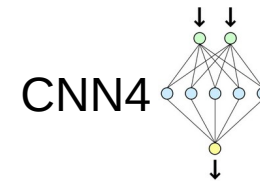
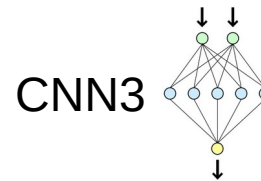
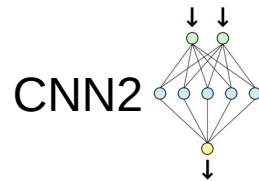
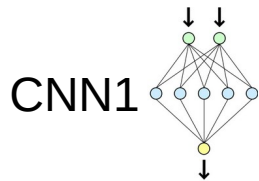
P1

P2

P3

P4

P5



Si on modifie 1 données, seul 1 des CNN est modifié

→ Si au moins 4 CNN sont d'accord sur la classe, alors cet accord n'a pas pu être influencé par la présence d'1 donnée empoisonnée

Modifications fortes mais minoritaires vs partition

Training data

Partitionnement

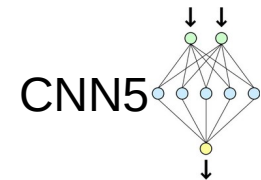
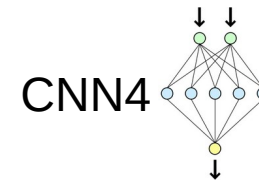
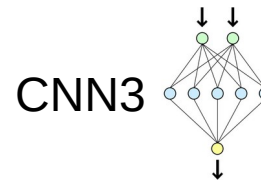
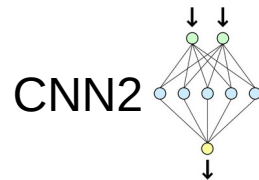
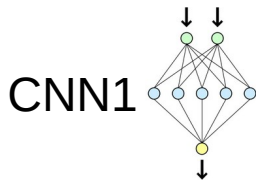
P1

P2

P3

P4

P5

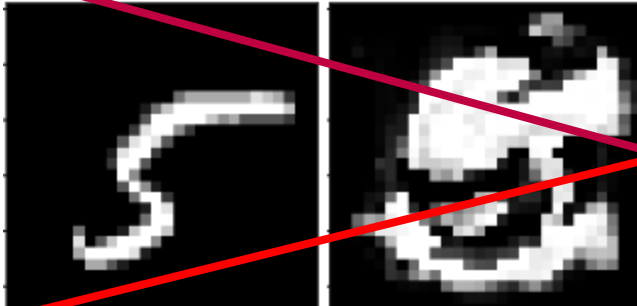


Si on modifie 1 données, seul 1 des CNN est modifié

→ Si au moins 4 CNN sont d'accord sur la classe, alors cet accord n'a pas pu être influencé par la présence d'1 donnée empoisonnée

→ bien entendu cette défense est difficile mais peut s'implémenter de façon souple
(RANSAC)

Scénarios d'empoisonnement



Modifications fortes mais peu nombreuses



Modifications invisibles

Plan

- scénario d'empoisonnement
- contre-mesure à une empoisonnement invisible
- perspectives

empoisonnement invisible

Poison Frog

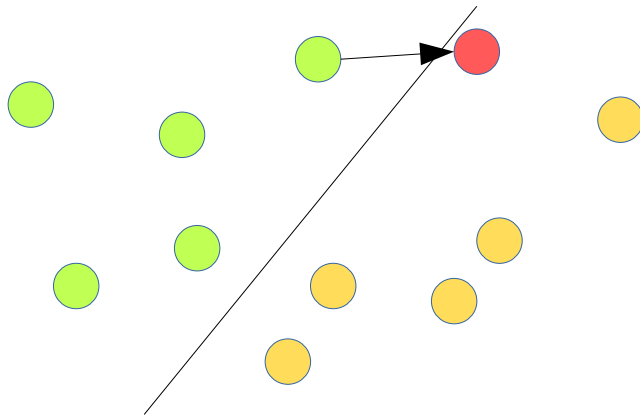
2.1 Crafting poison data via feature collisions

Let $f(\mathbf{x})$ denote the function that propagates an input \mathbf{x} through the network to the penultimate layer (before the softmax layer). We call the activations of this layer the *feature space* representation of the input since it encodes high-level semantic features. Due to the high complexity and nonlinearity of f , it is possible to find an example \mathbf{x} that “collides” with the target in feature space, while simultaneously being close to the base instance \mathbf{b} in input space by computing

$$\mathbf{p} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2 + \beta \|\mathbf{x} - \mathbf{b}\|_2^2 \quad (1)$$

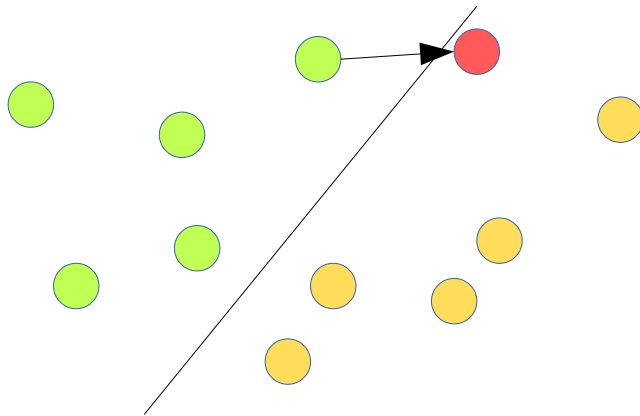
empoisonnement invisible

Poison Frog



empoisonnement invisible

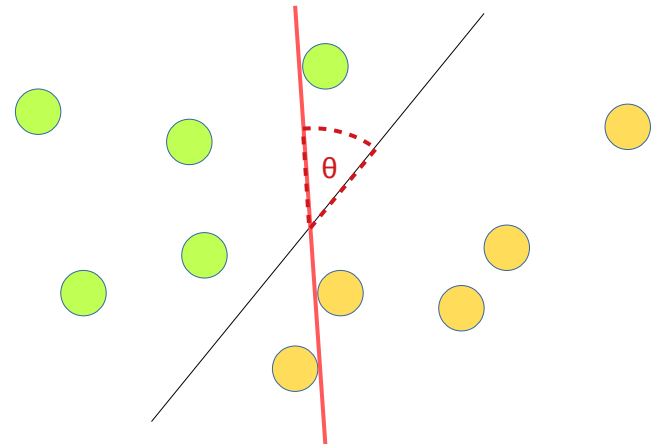
Poison Frog



Success rate 100 %
(Inception v3 + svm, Imagenet dog vs fish)

empoisonnement invisible

Proxy based



empoisonnement invisible

Proxy based

$$u \rightarrow u + \tau(\theta_{fair}^T u)\theta_{poison} - \tau(\theta_{poison}^T u)\theta_{fair}$$

$$\theta_{hack} \approx \theta_{fair} + \tau(\theta_{fair}^T \theta_{fair})\theta_{poison} - \tau(\theta_{poison}^T \theta_{fair})\theta_{fair}$$

empoisonnement invisible

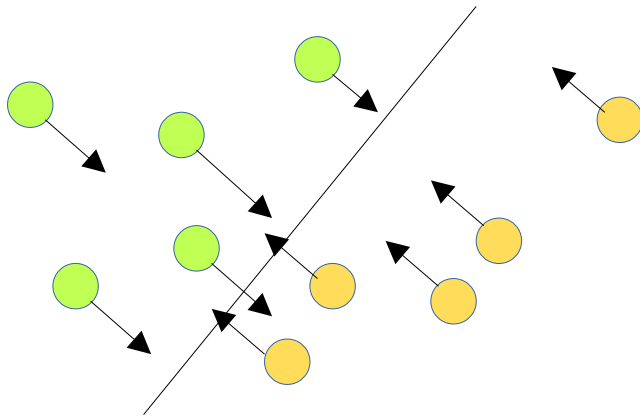
Proxy based

$$u \rightarrow u + \tau(\theta_{fair}^T u)\theta_{poison} - \tau(\theta_{poison}^T u)\theta_{fair}$$

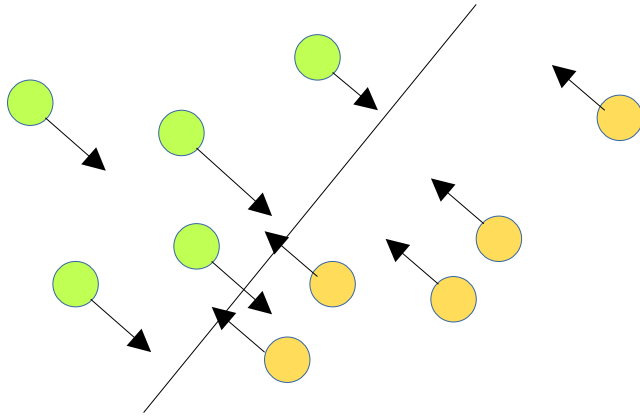
$$\theta_{hack} \approx \theta_{fair} + \tau(\theta_{fair}^T \theta_{fair})\theta_{poison} - \tau(\theta_{poison}^T \theta_{fair})\theta_{fair}$$

CIFAR 75 % \rightarrow 63 % accuracy (vgg + svm)

empoisonnement invisible

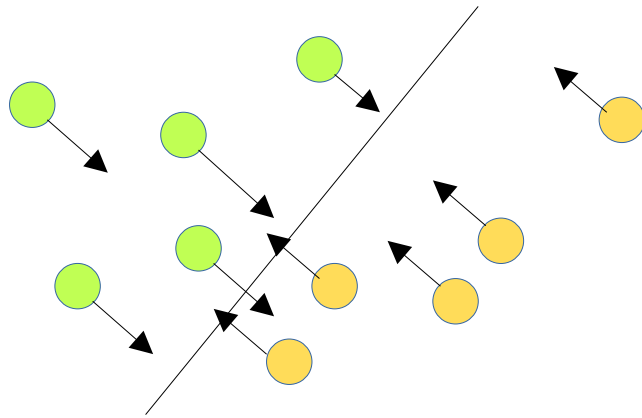


empoisonnement invisible



Un empoisonnement non coordonné ne conduit à une modification de comportement

empoisonnement invisible

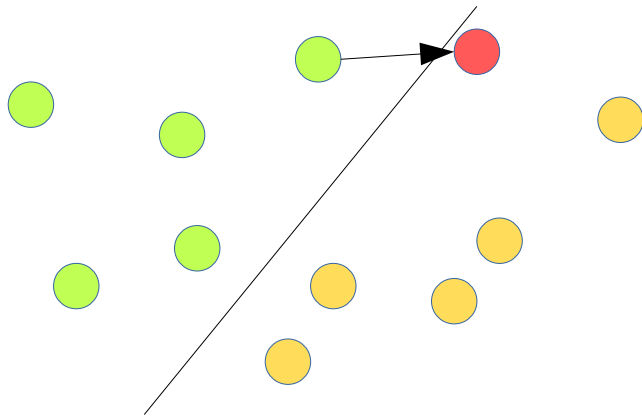


$SGD_{\theta}(f, Test)$	$\ll 87\%$ (-31% in [18])
$SGD_{\theta}(f, Train)$	$\approx 87\%$ (0% in [18])
$-w_{imagenet}$	$\approx 87\%$
$w_{imagenet}$	$\approx 87\%$
$-SGD_{\theta}(f, Train)$	$\approx 87\%$ (-1% in [18])
$-SGD_{\theta}(f, Test)$	$\gg 87\%$ (+7% in [18])

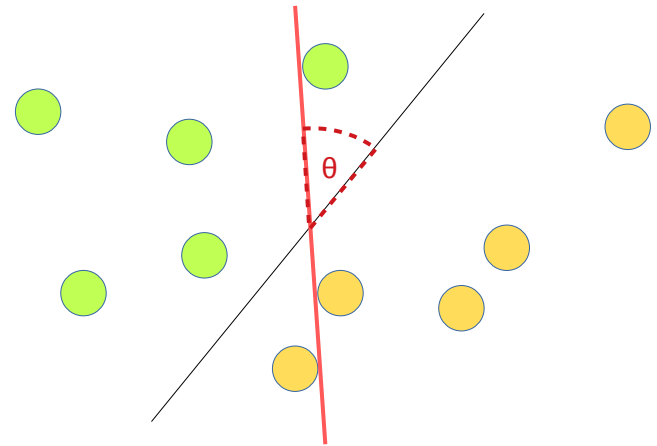
Un empoisonnement non coordonné ne conduit à une modification de comportement

empoisonnement invisible

Poison Frog



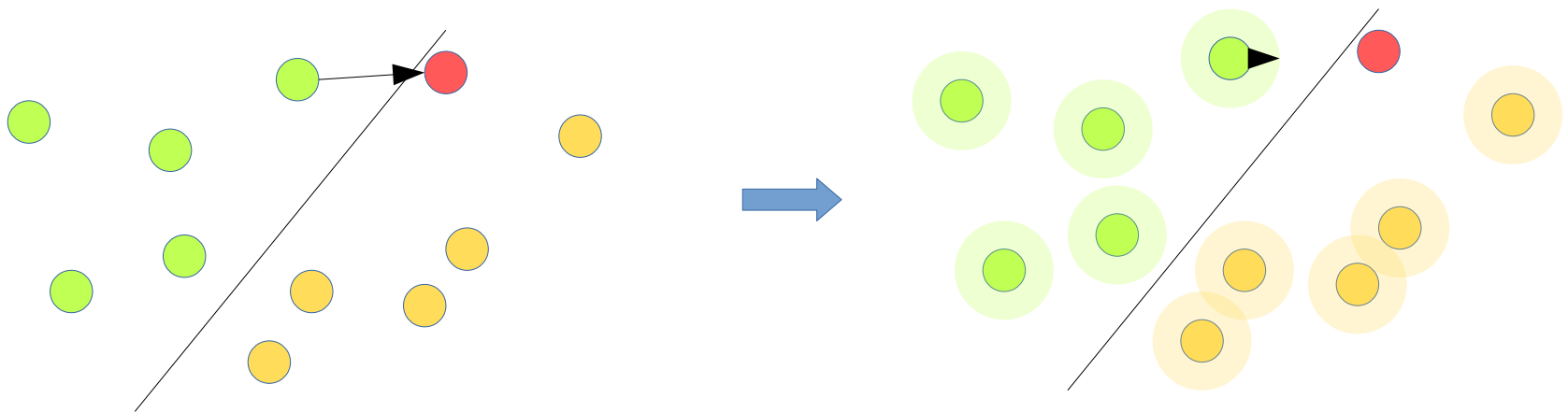
Proxy based



empoisonnement invisible

Clean label

Modification des images contraintes à être indétectable à l'oeil

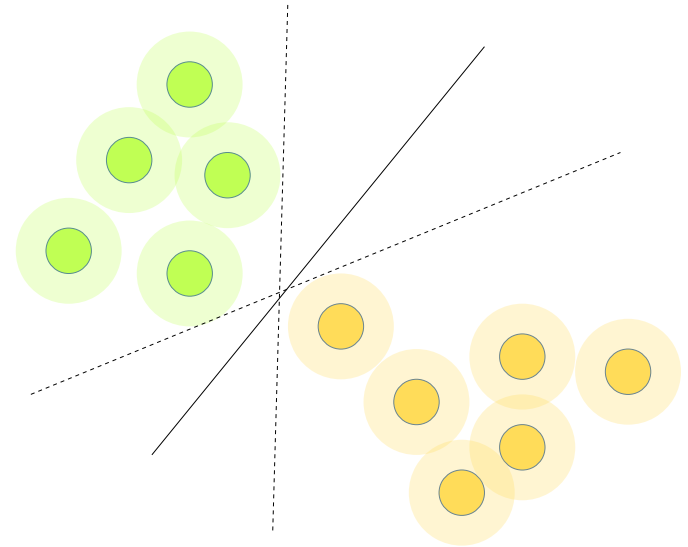
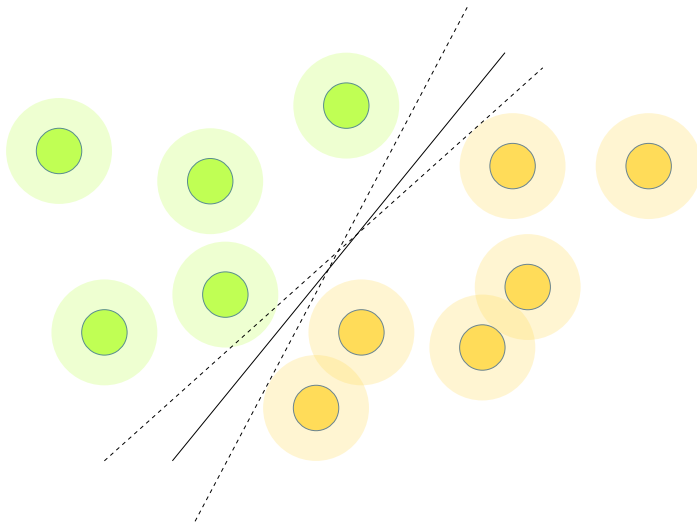


→ Si l'apprentissage est robuste, y a-t-il toujours sensibilité à cet empoisonnement ?

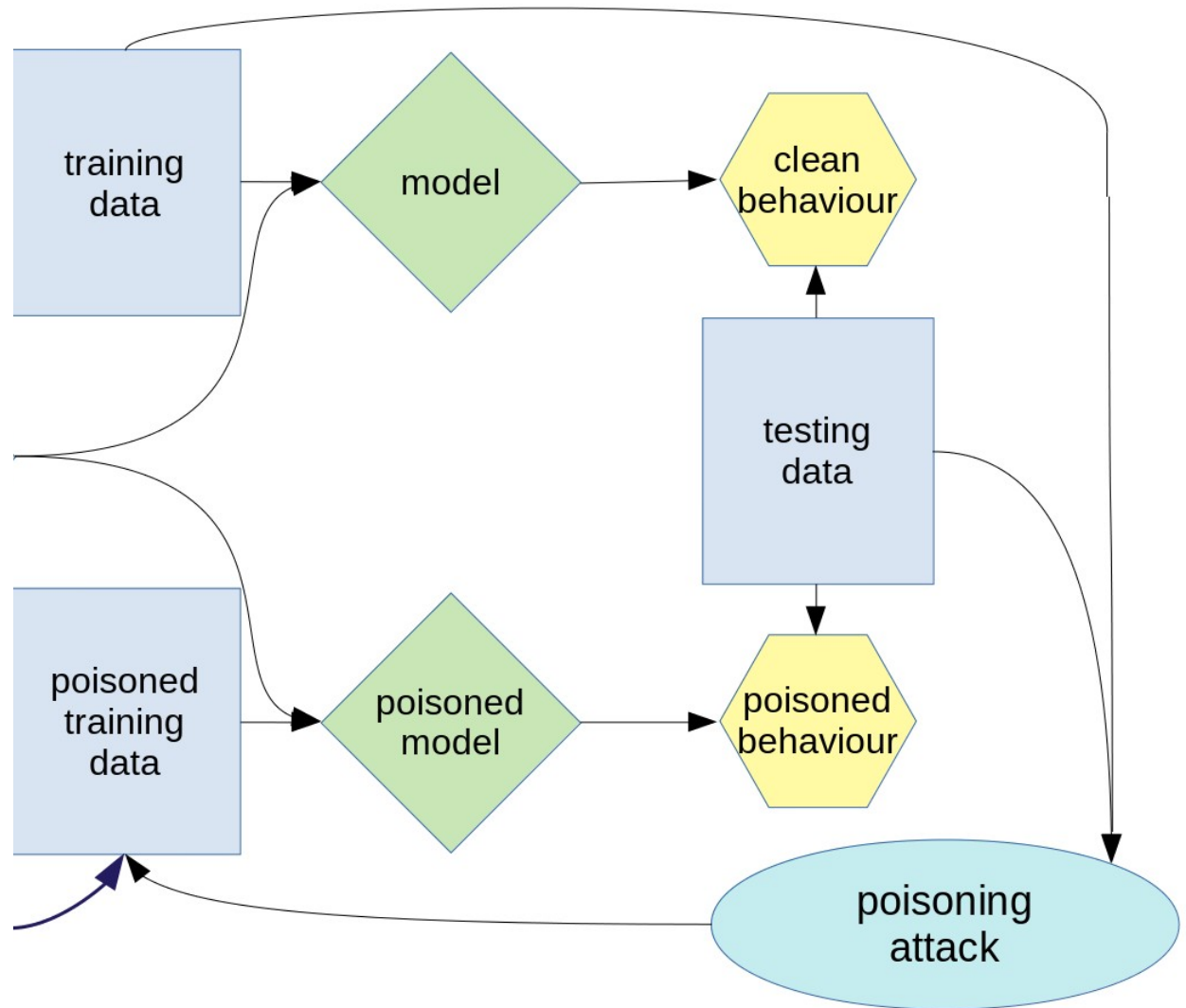
empoisonnement invisible et apprentissage robuste

→ Si l'apprentissage est robuste, y a-t-il toujours sensibilité à cet empoisonnement ?

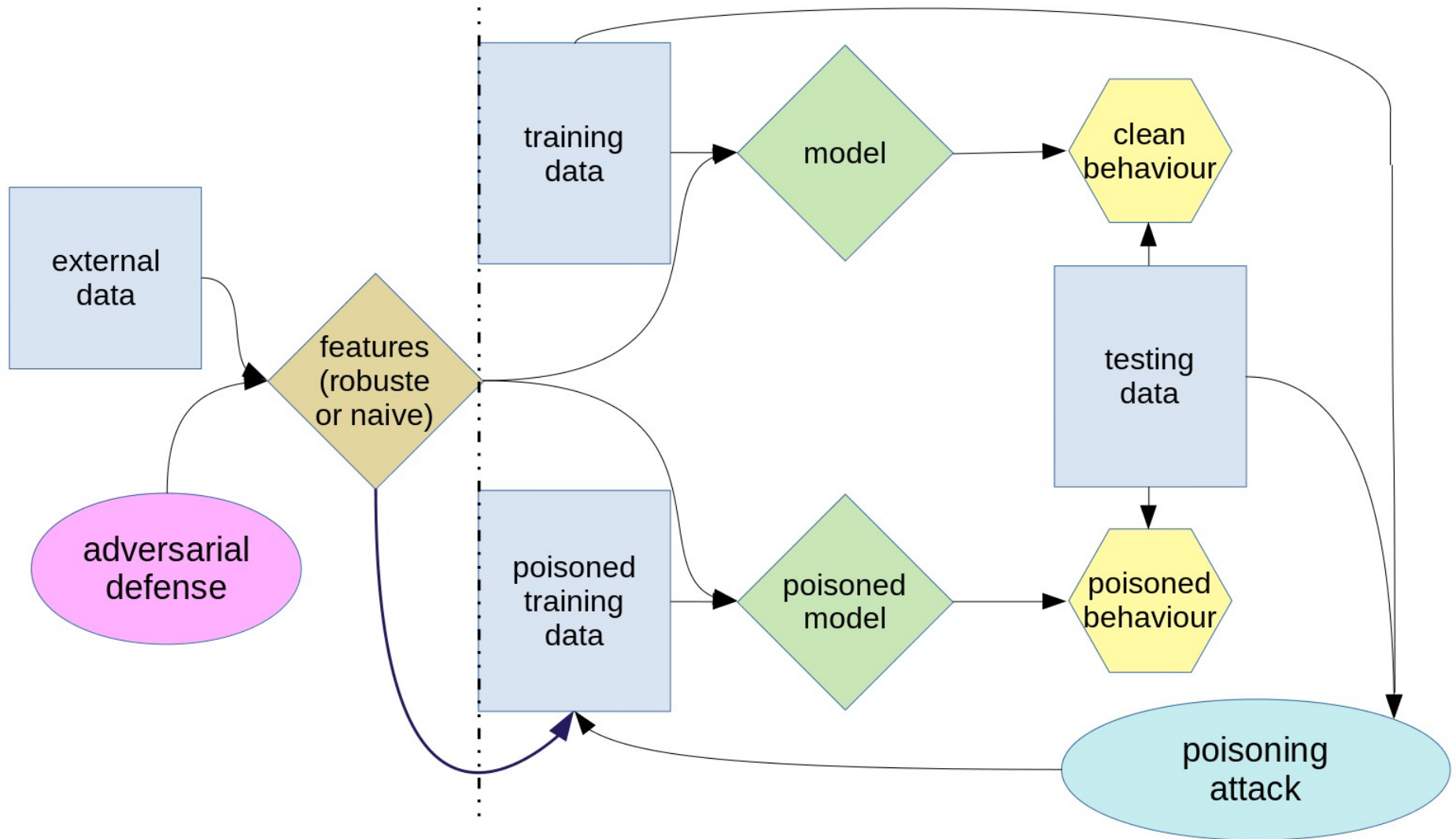
La question n'est pas tout à fait triviale :
on peut avoir des points margés mais plus ou moins dispersés



empoisonnement



empoisonnement invisible et apprentissage robuste



empoisonnement invisible et apprentissage robuste

Dataset	CIFAR	MNIST
AD on naive feature	24%	68%
AD on FSGM feature	30%	93%
AD on PGD feature	34%	95%

feature	CIFAR
PF on naive feature	85%
PF on FSGM feature	53%
PF on PGD feature	16%

Poison Frog (reimplémentation) passe de 85 % de succès à 16 % (CIFAR)

L'accuracy sous attaque AD qui avait chuté à 68 % remonte à 95 % (MNIST)

Relecture de la base + « RANSAC » + caractéristique robuste

Relecture → protection contre trop de données fortement corrompues

RANSAC → protection contre peu de données fortement corrompues

Robuste → protection contre des données faiblement corrompues

Conclusion

Quelque soit le scénario, l'empoisonnement ne paraît pas un risque majeur.

Conclusion

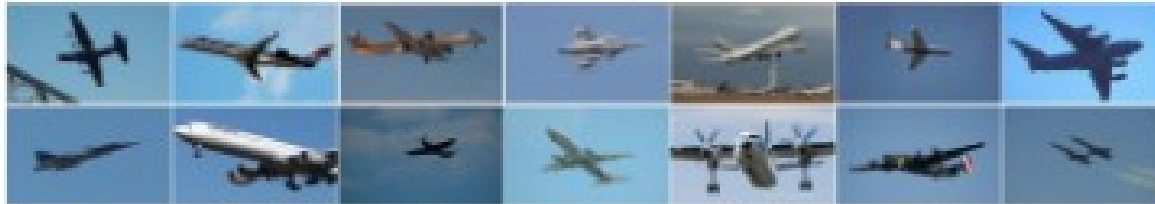
Quelque soit le scénario, l'empoisonnement ne paraît pas un risque majeur.

Sauf si l'apprentissage a lieu sans relecture de la base !

Perspectives

Les dangers du pseudo-labelling

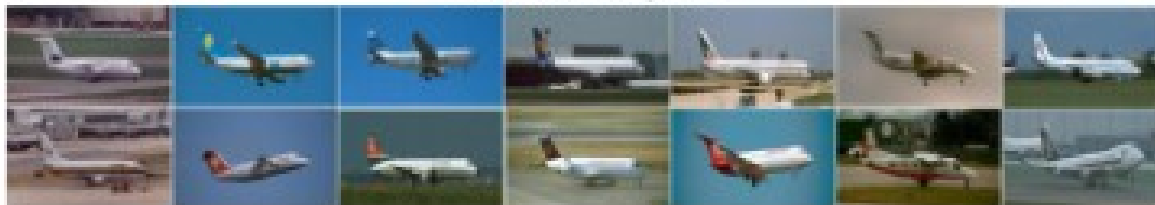
PASCAL airplane



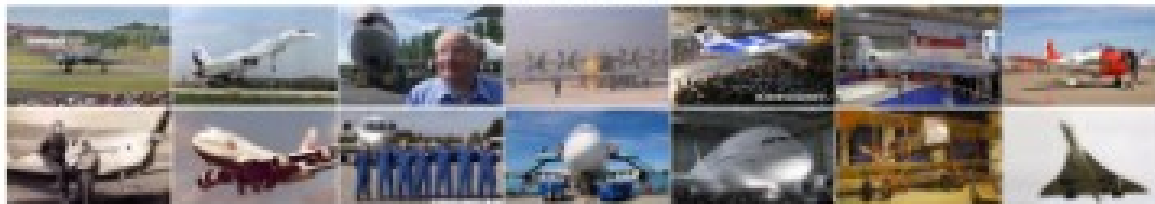
SUN airplane



Caltech101 airplane



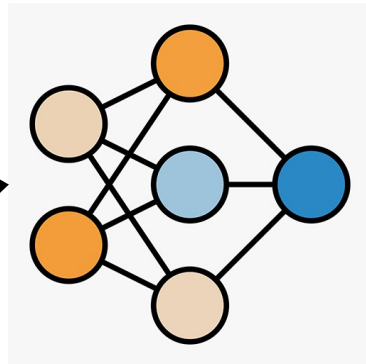
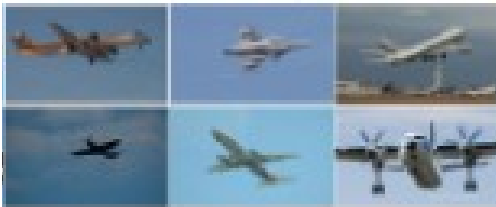
ImageNet airplane



Perspectives

Les dangers du pseudo-labelling

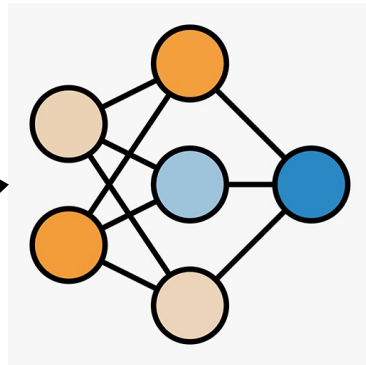
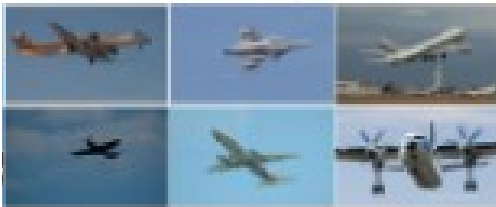
Domaine 1



Perspectives

Les dangers du pseudo-labelling

Domaine 1



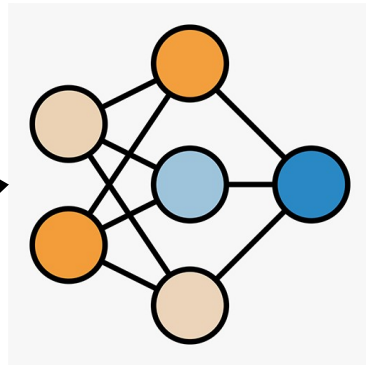
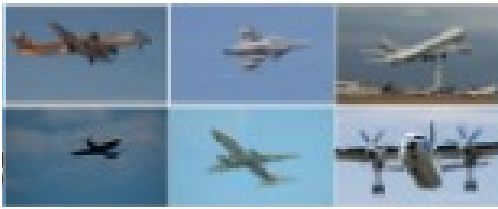
Domaine 2



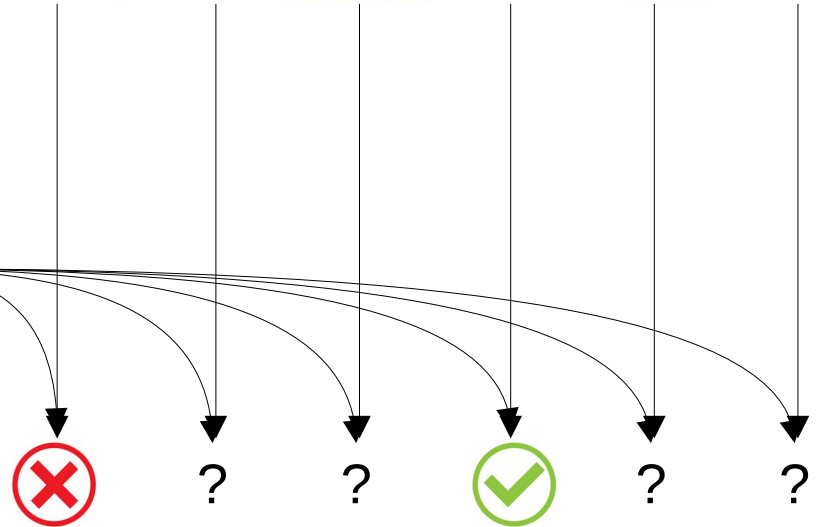
Perspectives

Les dangers du pseudo-labelling

Domaine 1



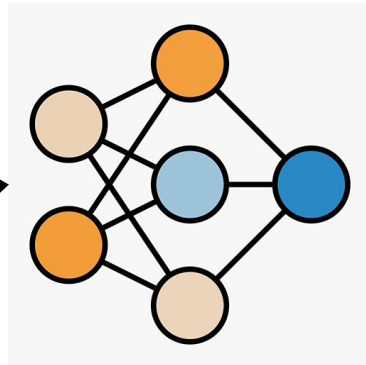
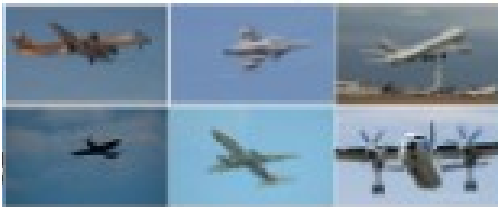
Domaine 2



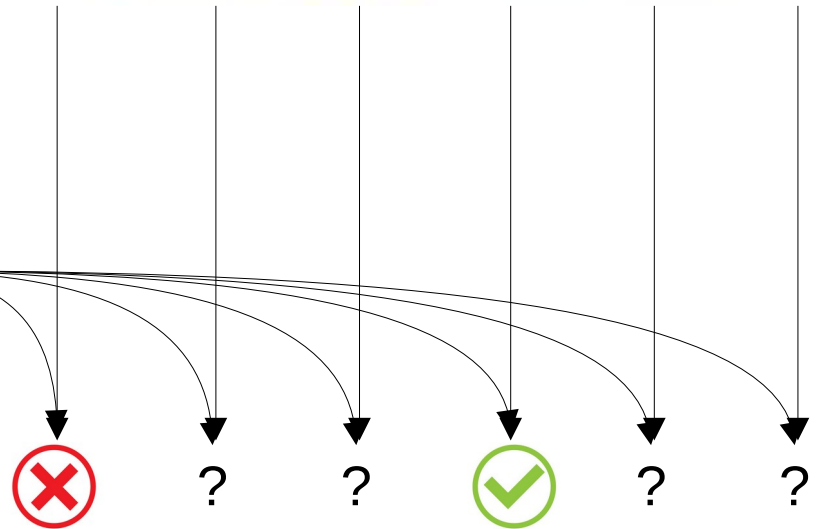
Perspectives

Les dangers du pseudo-labelling

Domaine 1



Domaine 2

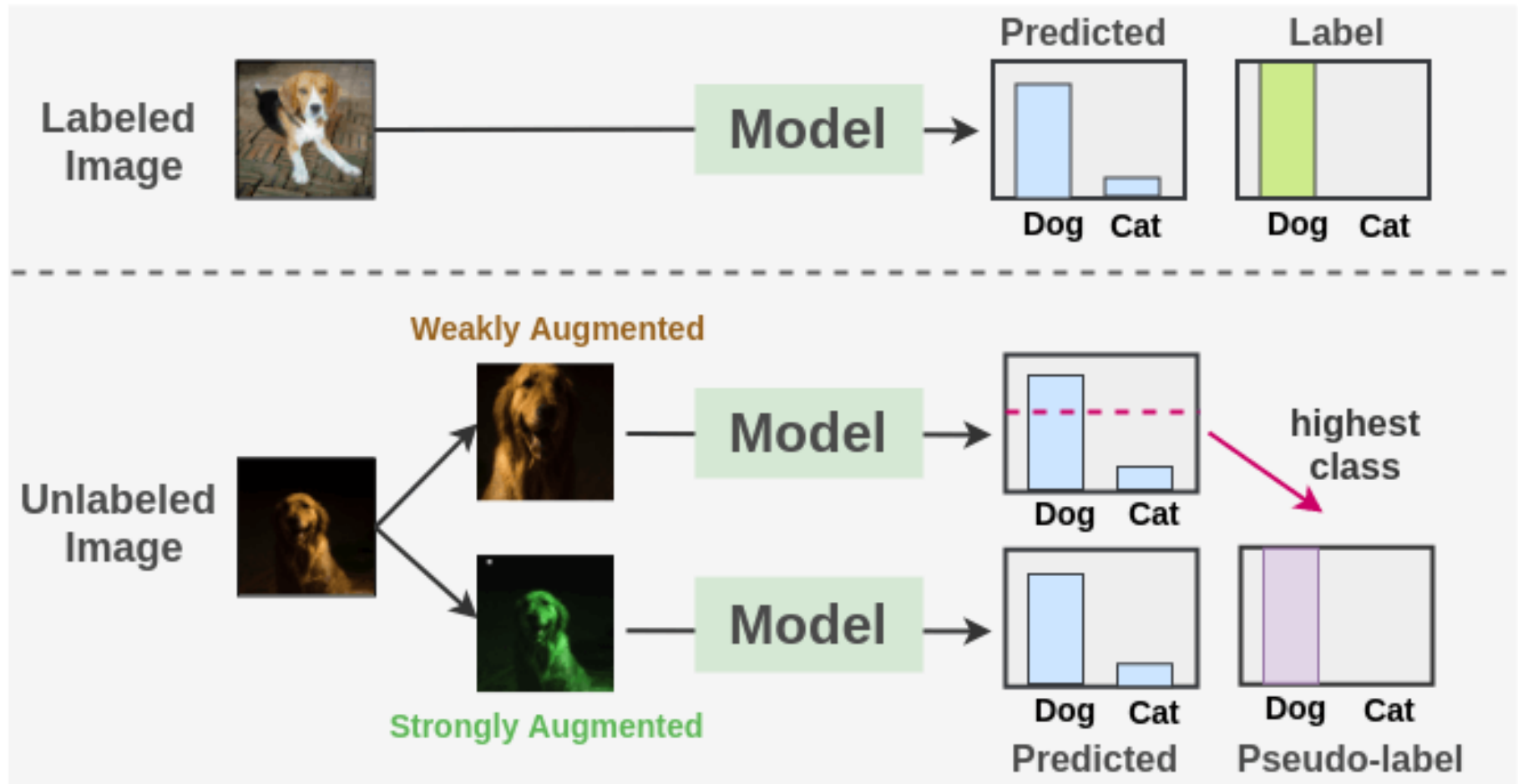


Confiance + consistance !

Perspectives

Les dangers du pseudo-labelling

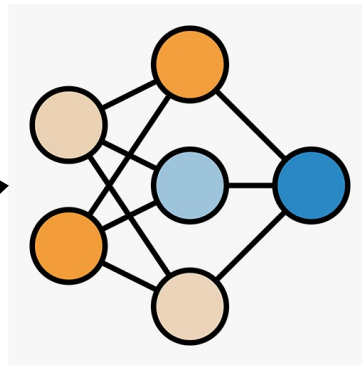
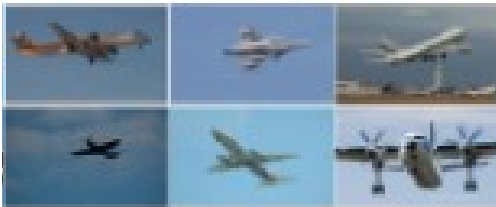
FixMatch Pipeline



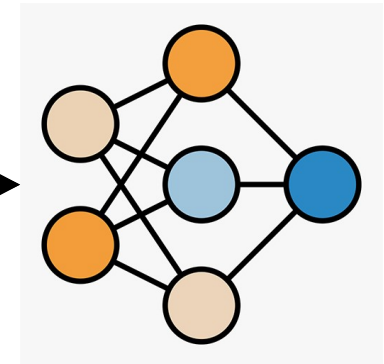
Perspectives

Les dangers du pseudo-labelling

Domaine 1



Domaine 2



Perspectives

Les dangers du pseudo-labelling

Est-il possible de créer une perturbation qui

Augmente la confiance

Résiste à un test de consistance

Perspectives

Les dangers du pseudo-labelling

Est-il possible de créer une perturbation qui

Augmente la confiance

Résiste à un test de consistance

Et qui une fois dans la base d'apprentissage provoque un empoisonnement ???

Perspectives

Les dangers du pseudo-labelling

Est-il possible de créer une perturbation qui

Augmente la confiance

Résiste à un test de consistance

Et qui une fois dans la base d'apprentissage provoque un empoisonnement ???

À suivre !

Perspectives

Les limites de ne considérer que « deep feature + SVM »

Là où les attaques adversaires sont apparues avec le DL,
l'empoisonnement est plus facile avec une approche type SVM !

Perspectives

Les limites de ne considérer que « deep feature + SVM »

Là où les attaques adversaires sont apparues avec le DL,
l'empoisonnement est plus facile avec une approche type SVM !

→ typiquement une approche SVM est beaucoup plus sensible au bruit des labels !
(en fonction de « C »)

Perspectives

Les limites de ne considérer que « deep feature + SVM »

Là où les attaques adversaires sont apparues avec le DL,
l'empoisonnement est plus facile avec une approche type SVM !

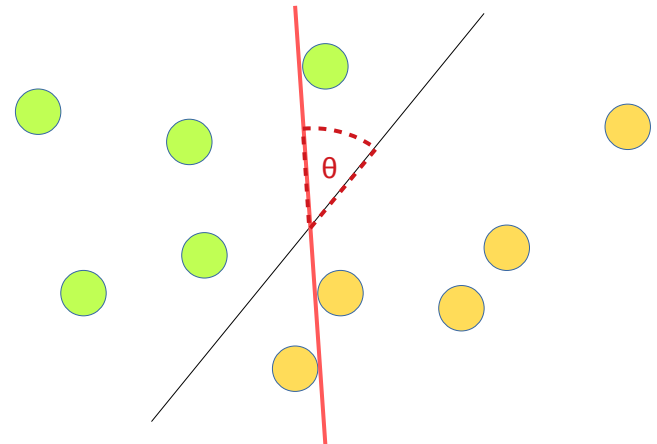
→ typiquement une approche SVM est beaucoup plus sensible au bruit des labels !
(en fonction de « C »)

→ empoisonner un CNN est « plus difficile »
mais ouvre potentiellement la porte à d'autres conclusions ...

Perspectives

Les limites de ne considérer que « deep feature + SVM »

Proxy based



Perspectives

Les limites de ne considérer que « deep feature + SVM »

proxy used Eq.3	testing accuracy	desired
$SGD_{\theta}(f, Test)$	27%	$\ll 87\%$ (-31% in [18])
$SGD_{\theta}(f, Train)$	34%	$\approx 87\%$ (0% in [18])
$-w_{imagenet}$	64%	$\approx 87\%$
$w_{imagenet}$	58%	$\approx 87\%$
$-SGD_{\theta}(f, Train)$	73%	$\approx 87\%$ (-1% in [18])
$-SGD_{\theta}(f, Test)$	77%	$\gg 87\%$ (+7% in [18])
Original accuracy	87%	-

Perspectives

Les limites de ne considérer que « deep feature + SVM »

proxy used Eq.3	testing accuracy	desired
$SGD_{\theta}(f, Test)$	27%	$\ll 87\%$ (-31% in [18])
$SGD_{\theta}(f, Train)$	34%	$\approx 87\%$ (0% in [18])
$-w_{imagenet}$	64%	$\approx 87\%$
$w_{imagenet}$	58%	$\approx 87\%$
$-SGD_{\theta}(f, Train)$	73%	$\approx 87\%$ (-1% in [18])
$-SGD_{\theta}(f, Test)$	77%	$\gg 87\%$ (+7% in [18])
Original accuracy	87%	-

Mais depuis Imagenet seulement (sinon ça marche comme prévu) !

Perspectives

Les limites de ne considérer que « deep feature + SVM »

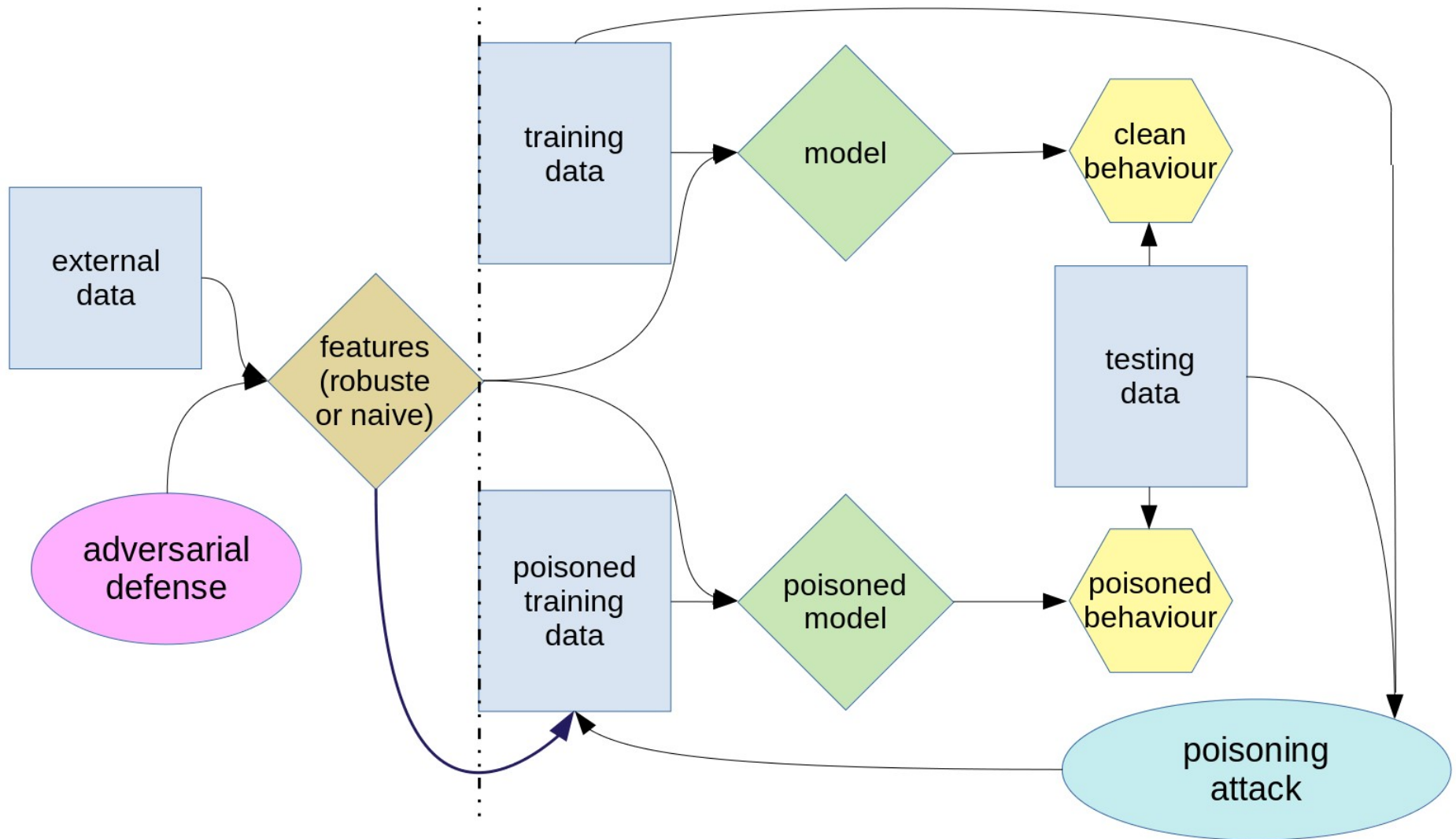
proxy used Eq.3	testing accuracy	desired
$SGD_{\theta}(f, Test)$	27%	$\ll 87\%$ (-31% in [18])
$SGD_{\theta}(f, Train)$	34%	$\approx 87\%$ (0% in [18])
$-w_{imagenet}$	64%	$\approx 87\%$
$w_{imagenet}$	58%	$\approx 87\%$
$-SGD_{\theta}(f, Train)$	73%	$\approx 87\%$ (-1% in [18])
$-SGD_{\theta}(f, Test)$	77%	$\gg 87\%$ (+7% in [18])
Original accuracy	87%	-

Mais depuis Imagenet seulement (sinon ça marche comme prévu) !

(l'impact sur le trajet est plus important que l'impact sur l'arrivé)

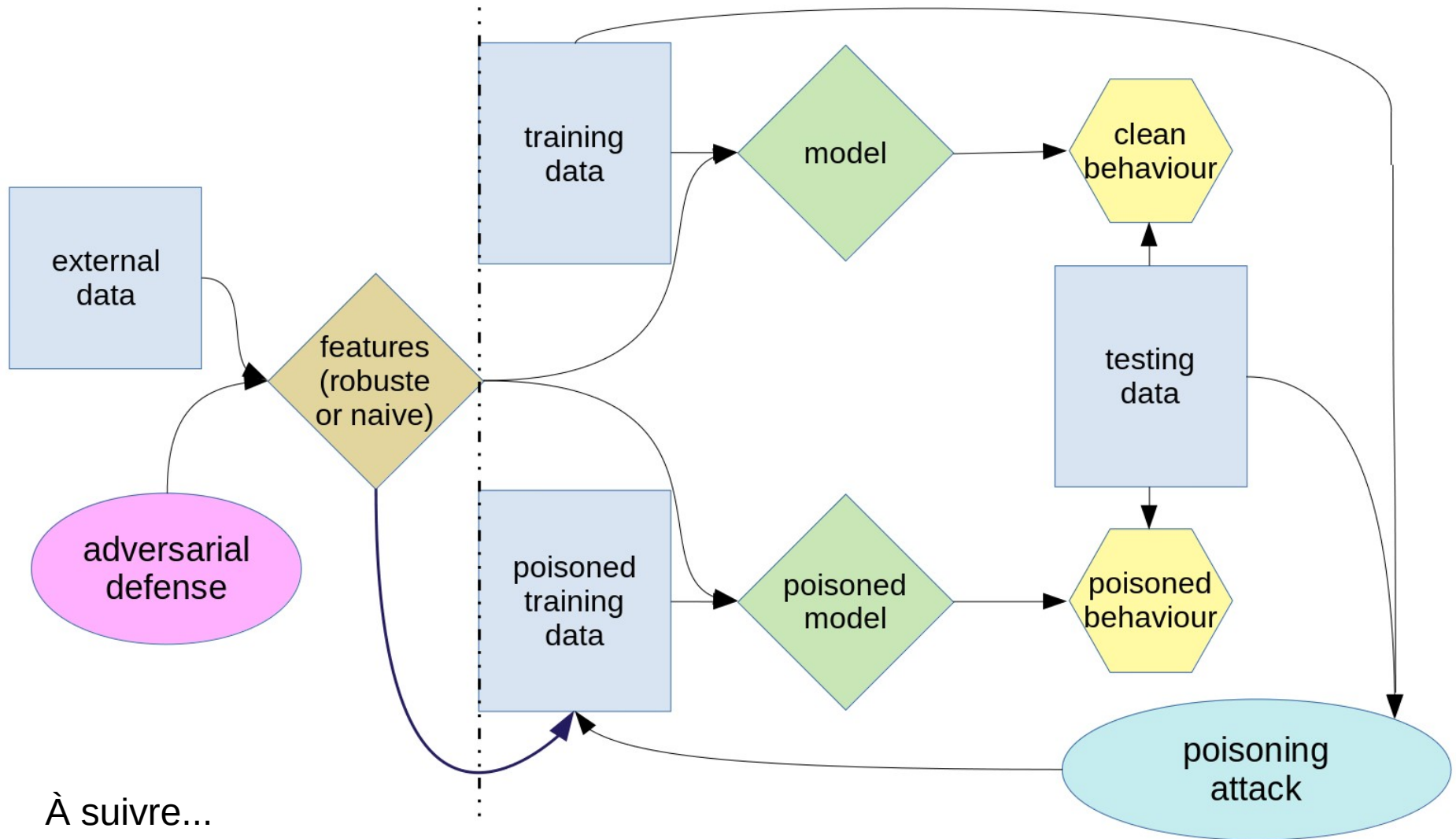
Perspectives

Les limites de ne considérer que « deep feature + SVM »



Perspectives

Les limites de ne considérer que « deep feature + SVM »



Merci pour votre attention.

Les codes sont disponibles à github.com/achanhon/AdversarialModel