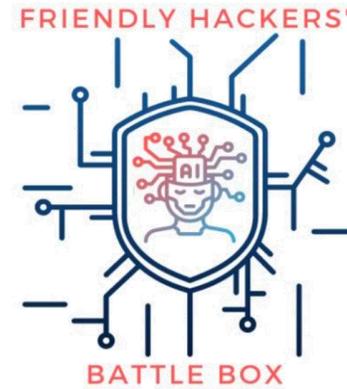
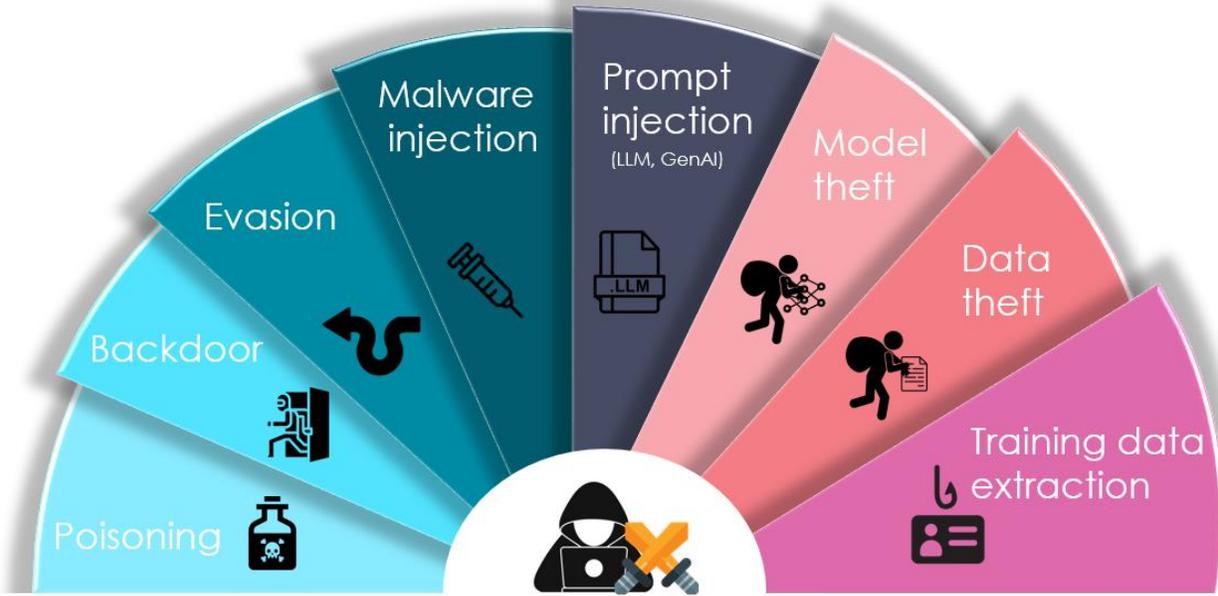


AI Friendly Hackers

Thales initiative on
secure AI

Explore intrinsic AI vulnerabilities to develop better countermeasures



Model backdooring



Model inversion

OPEN



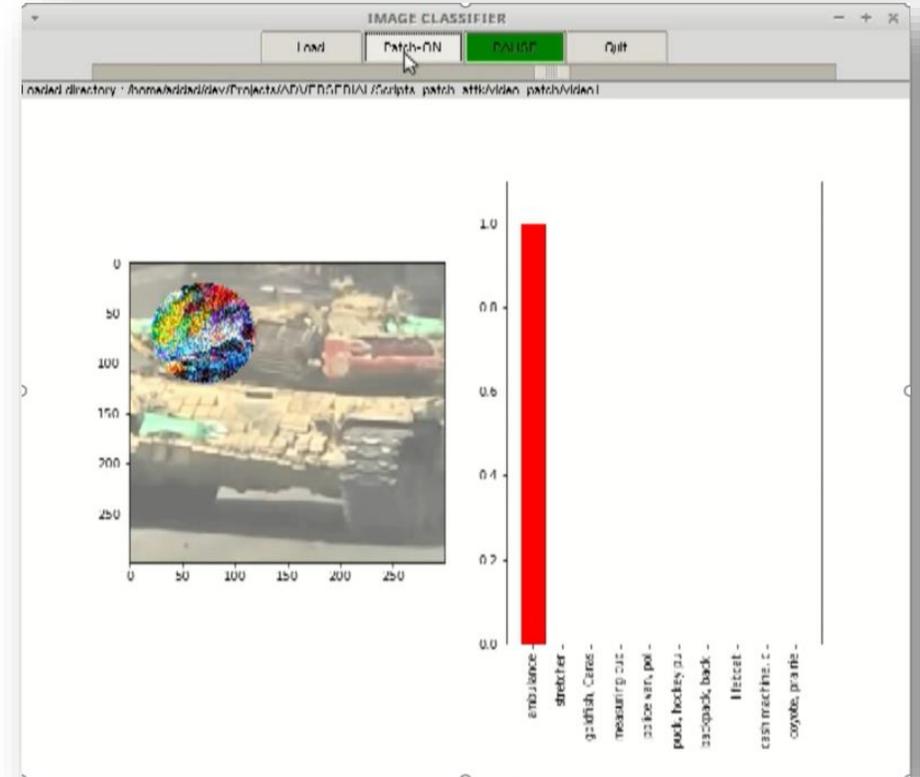
Demos

OPEN

Evasion (adversarial) attack



Patch attack



Crafting inputs in a specific way to get the wrong result from the model

Unlearning Challenge @CAID 2023



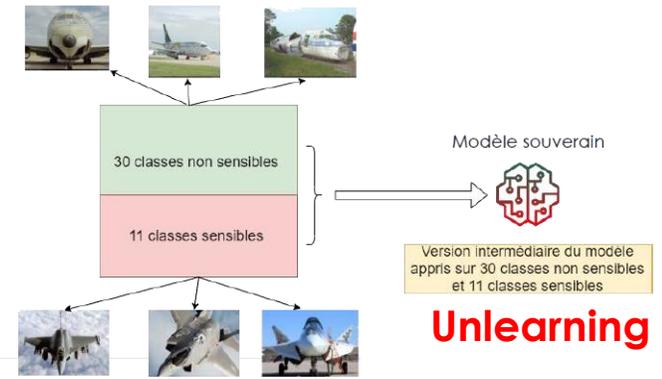
> Thales Friendly Hackers team won the first prize of CAID 2023 challenge by crafting two novel AI privacy attacks

> Membership attack:

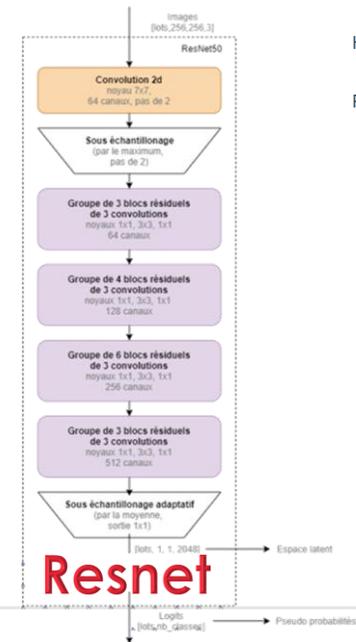
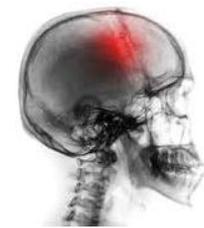
▶ identifying data samples that were in the training dataset

> Unlearning attack:

▶ recovering sensitive information from a model that passed through an unlearning procedure



DISCLAIMER
Multiple models had suffered during the challenge...



Leaderboard

Tâche A : Membership Attack

Friendly hackers	Soumission 6 (sept)	0.653125
Friendly hackers	Soumission 7 (sept)	0.642500
Friendly hackers	Soumission 5 (août)	0.640000
HackCuda MaData	Soumission 3 (août)	0.617500
HackCuda MaData	Soumission 2 (juillet)	0.613750
Friendly hackers	Soumission 4 (août)	0.608750

