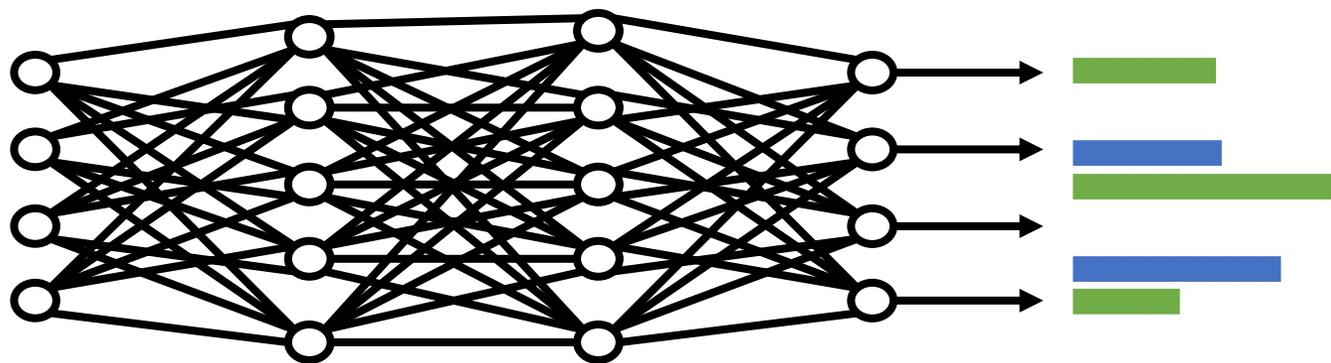# Cyber Security and AI

Dr. Britta Hale

# Cybersecurity and Machine Learning

- How to break ML's security
- How to secure ML
- How to use ML to improve cybersecurity
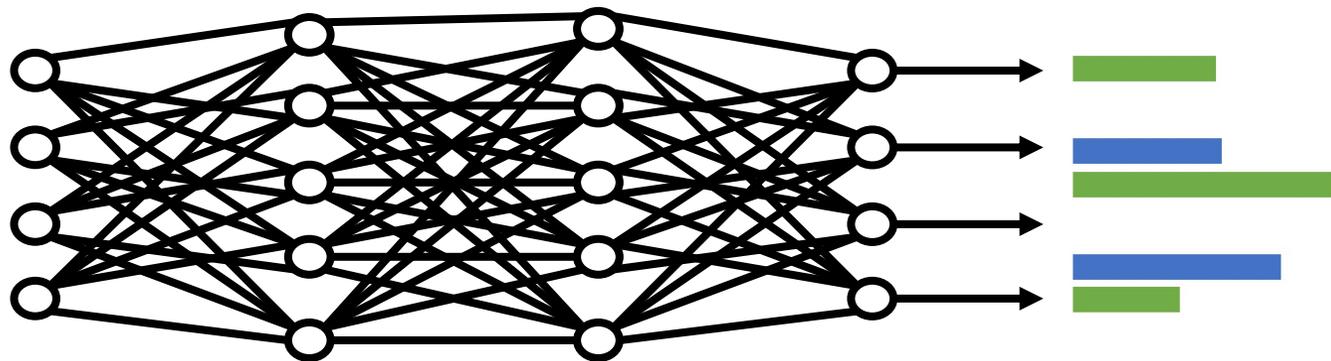
# Does the training work?



*Accuracy*

Can the training be circumvented?

Can the model be misinterpreted?

Can the model be abused?

*Everything Else*

# *Attacks During Training*

e.g.

- Poisoning

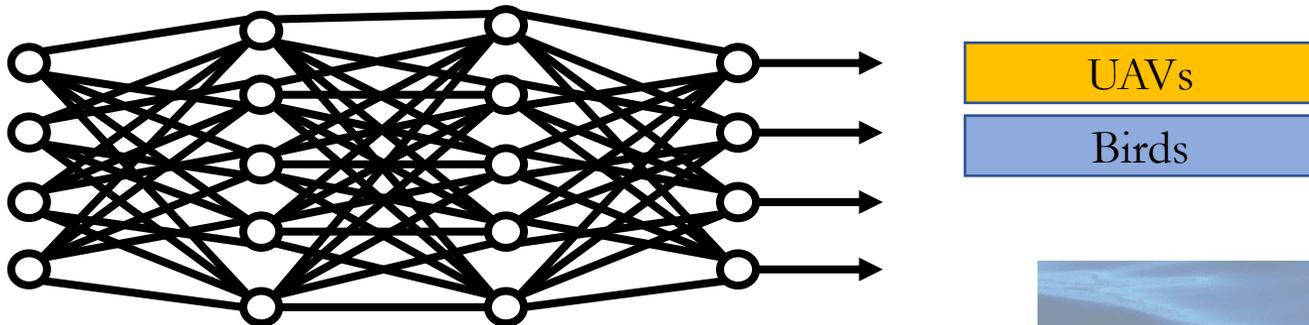- Trojans/Backdoors

# A perfect memory…

# Poisoning

***Integrity***

- Confidence reduction

    do not change a class but highly impact the confidence

- Misclassification

    change a class without any specific target

- Targeted misclassification

    change a class to a particular target

| UAVs |
| Birds |

# Poisoning

- Source/target misclassification

    change a particular source to a particular target
- Universal misclassification

    change any source to particular target



"panda"
57.7% confidence

$+ .007 \times$

"nematode"
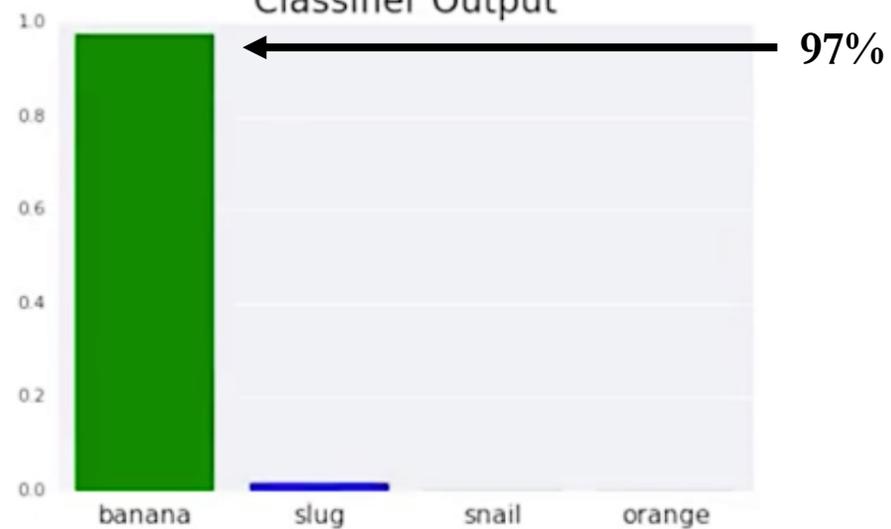8.2% confidence

$=$

"gibbon"
99.3 % confidence

Goodfellow,
Shlens, Szegedy
ICLR 2015

place sticker on table

Classifier Input

Classifier Output

97%

banana    slug    snail    orange

Classifier Input

Classifier Output

99%

toaster    banana    piggy_bank    spaghetti_

Tom B. Brown, Dandelion Mané , Aurko Roy, Martín Abadi, Justin Gilmer
https://arxiv.org/pdf/1712.09665.pdf

# Trojans/Backdoor

1. Inverse network to create a trojan trigger

2. Retrain model with malicious data

3. Real inputs which activate the trojan trigger generate malicious behavior

Access to original dataset not necessarily required

Retraining can take minutes/hours (vs. weeks/months for original model)

speedlimit 0.947

STOP

Gu, Dolan-Gavitt,
Garg 2019

# Defense

1.  Outlier detection

    **How to define an outlier?**

    **What about data that was injected before filtering rules?**

2.  Test newly added training samples against current model for accuracy

    **What about trojans?**

# *DATA*

ISSIE LAPOWSKY SECURITY 03.17.2018 12:20 PM

# Cambridge Analytica Took 50M Facebook Users' Data—And Both Companies Owe Answers

*The New York Times*

## Facebook and Cambridge Analytica: What You Need to Know as Fallout Widens
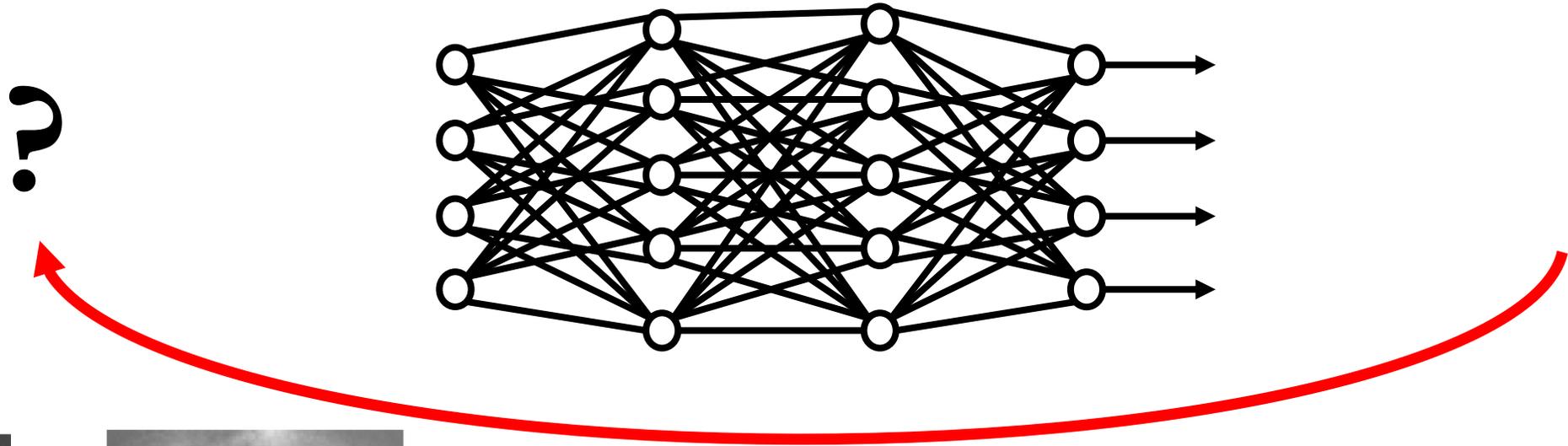
By Kevin Granville

March 19, 2018

f

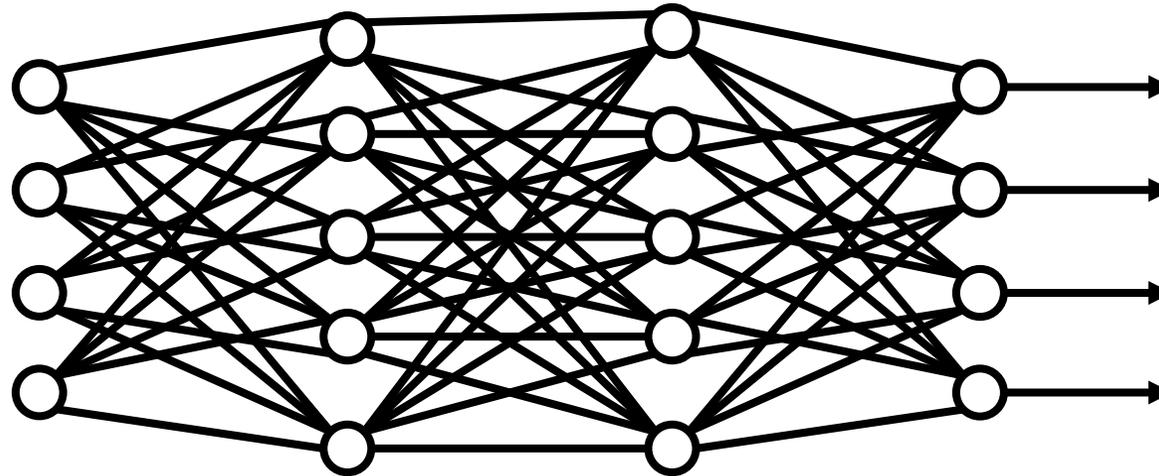# *Attacks During Production*

e.g.

- Inference

- Evasion

# Inference

- Acquire information about dataset

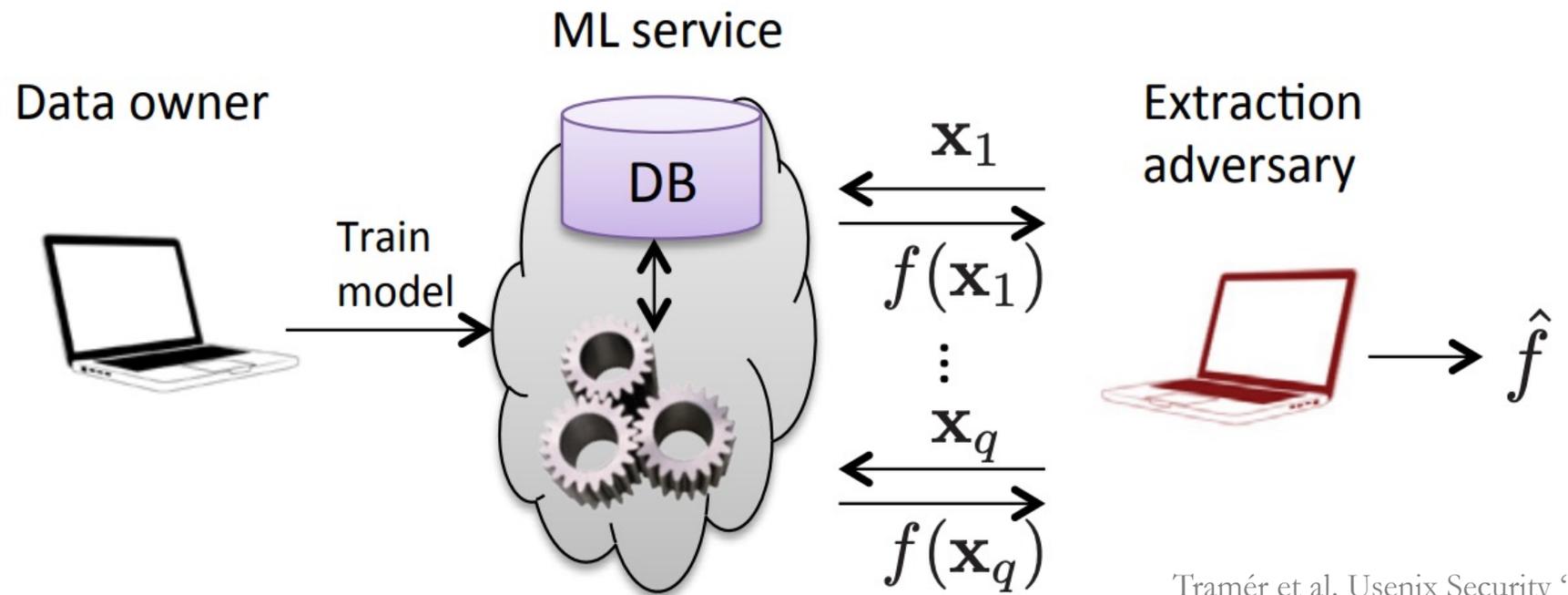**?**



Fredrikson, Jha, and Ristenpart CCS '15

# Inference

- Acquire information about dataset
- Membership inference / data attributes

# Inference

- Acquire information about dataset
- Membership inference / data attributes
- Model Extraction



Tramér et al. Usenix Security '16

# Evasion

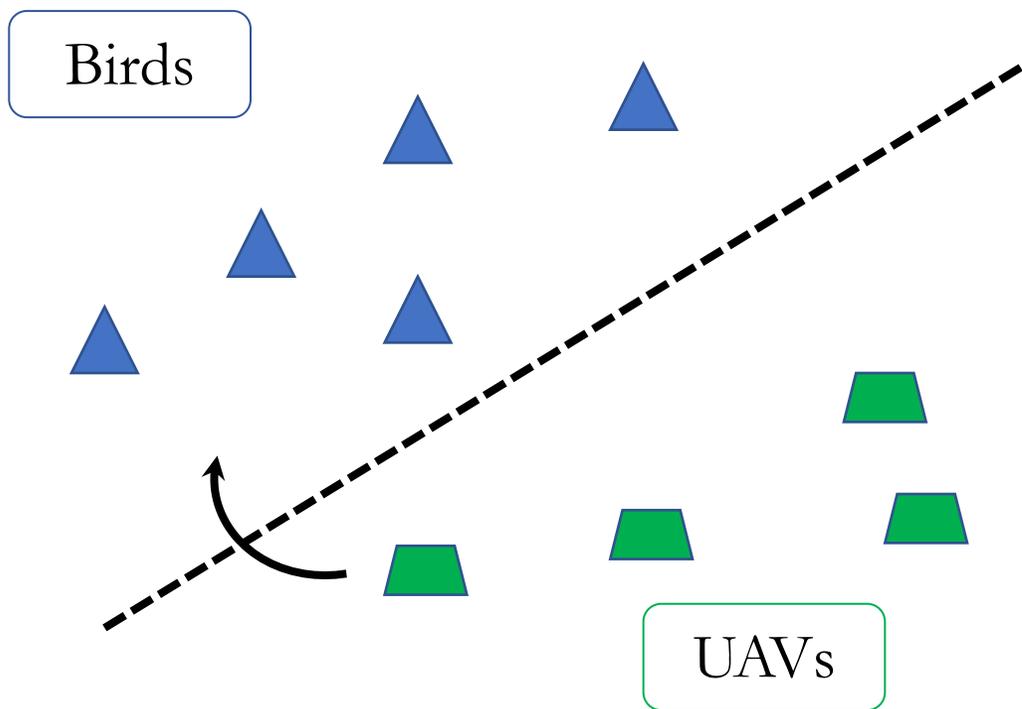Does not shift classifier boundary, but pushes poisoning into dataset

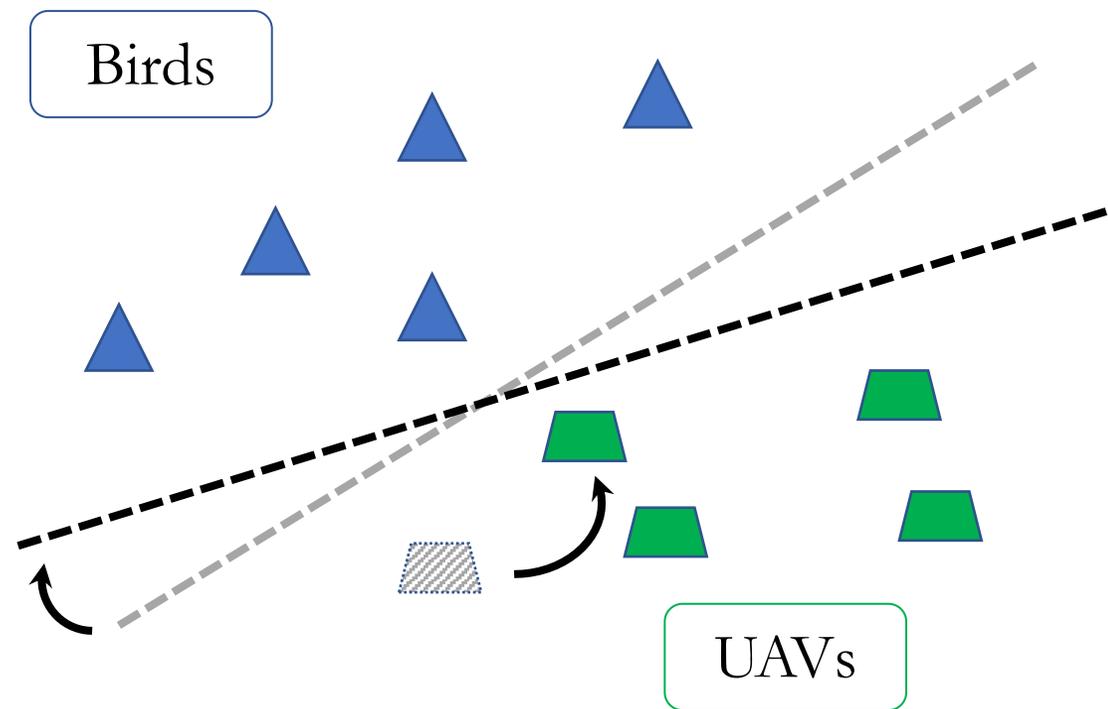# Evasion

To classify Birds as UAVs:

1. Change some UAVs to look closer to Birds

2. Keep UAVs labelled as UAVs

3. Add changed UAVs to training pool

Does not shift classifier boundary, but pushes poisoning into dataset
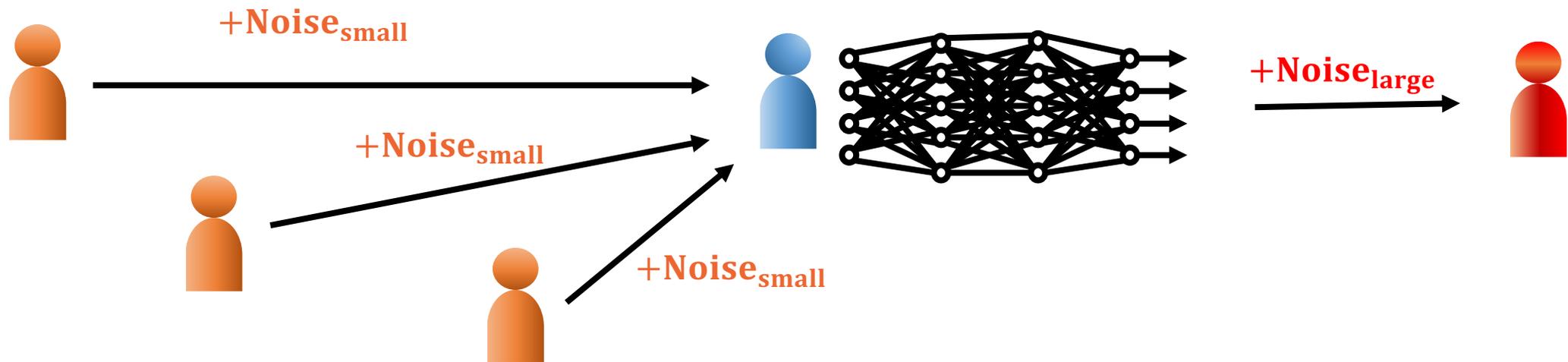
Evasion

Poisoning

Birds

Birds

UAVs

UAVs

# Defense

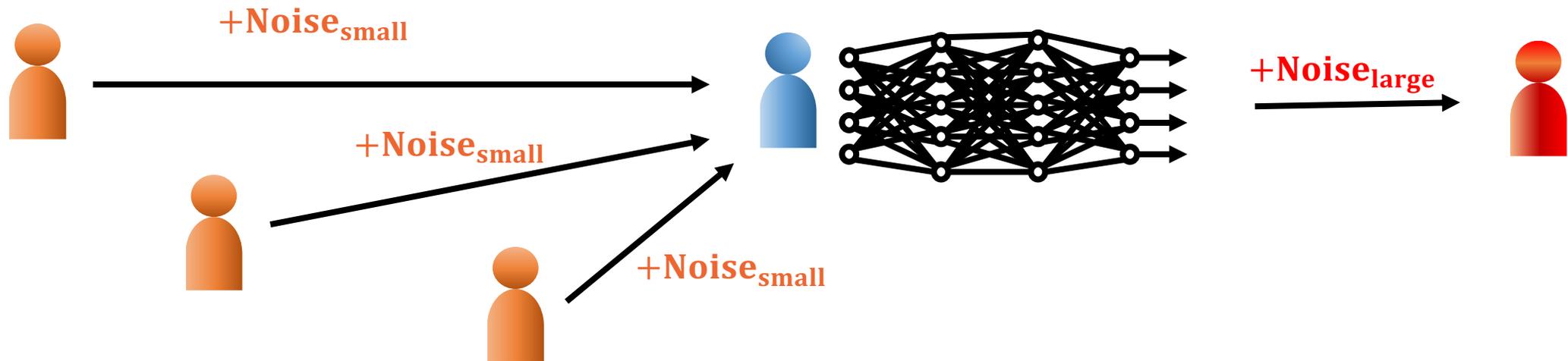1. Differential Privacy

# Defense

1. Differential Privacy

Goal: Try to hide individual data points

# Defense

1. Differential Privacy

*Problem: Model may become imbalanced*

# Defense

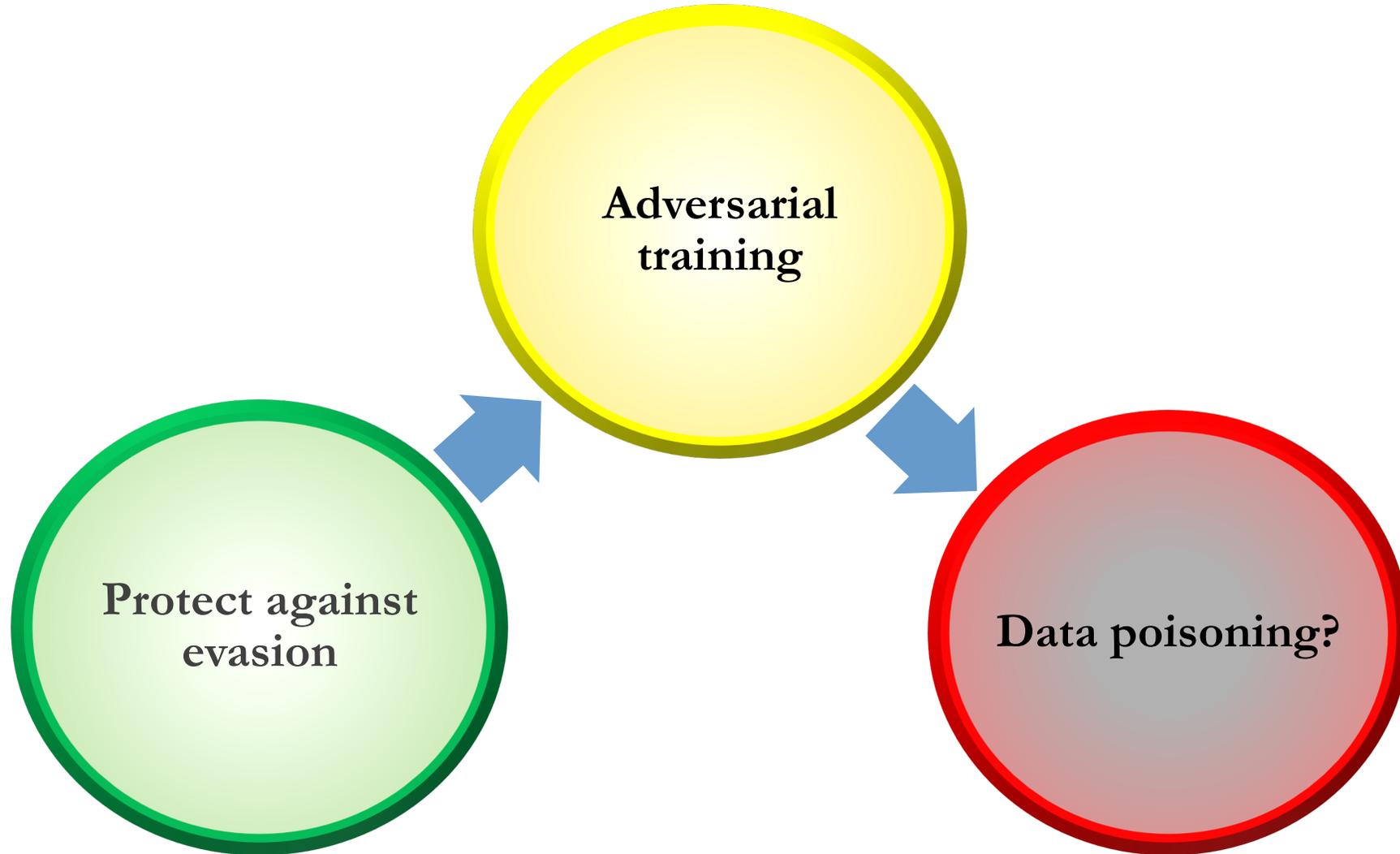1. Differential Privacy

2. Don't force guessing ("null" class)
   **Human overhead**

3. Adversarial training
   **What if the adversary uses different examples?**
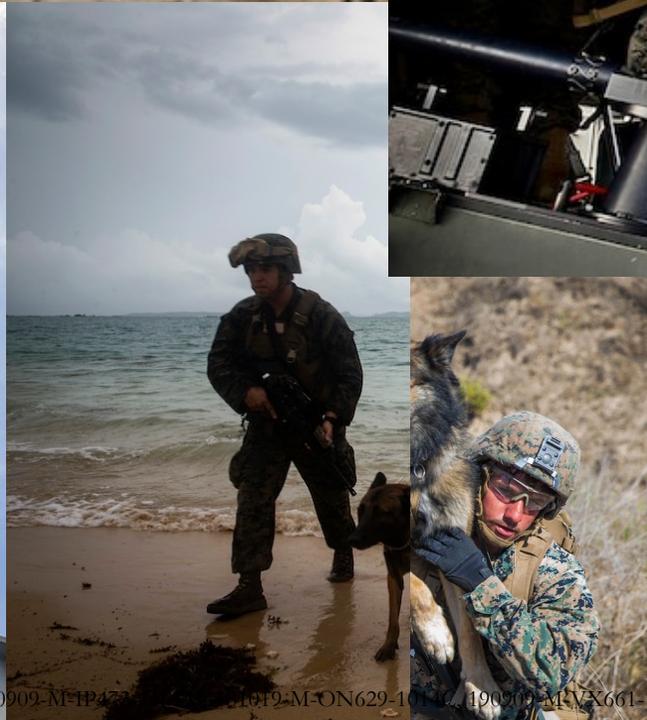   **What if you train on too many adversarial examples?**

# Poisoning vs. Evasion

# Human-in-the-loop

## experiment

# DevSecOps-AI

# Q&A
## Cyber Security and AI