

SECURITY WEAKNESSES IN MACHINE LEARNING

DANIEL ETZOLD - @ETZOLDIO

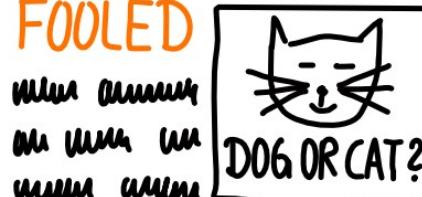
MONDAY 25TH MARCH

NAIVE BAYES FILTER BYPASSED



A hand-drawn diagram in the bottom-left corner shows a cross-section of a mountain range. It consists of several wavy lines representing layers of rock. Two arrows point downwards from the top of the highest peak, indicating the direction of downward movement or thrusting.

IMAGE CLASSIFIER FOOLED



THIEVES STEAL MODEL PARAMETERS

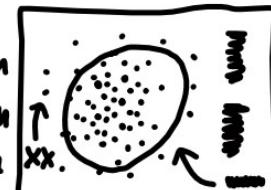
$$f(x) = Wx + b$$

$$= \sum_{i=1}^n w_i x_i + b$$

they are
all well wh-
enver they
are away from

меньше

ANOMALY DETECTION HACKED





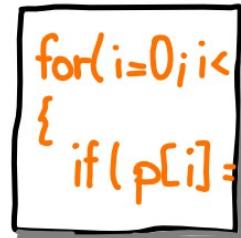
@ETZOLDIO

[HTTPS://ETZOLD.IO](https://etzold.io)

- IT SECURITY ARCHITECT



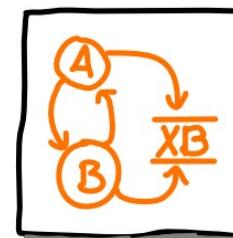
- **1&1** MAIL&MEDIA DEVELOPMENT & TECHNOLOGY GMBH



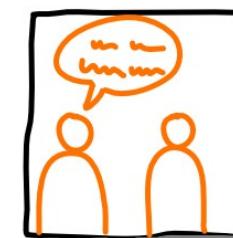
CODE
REVIEW



DOCUMENTATION
REVIEW



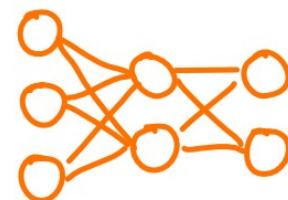
THREAT
MODELING



IN HOUSE
CONSULTING



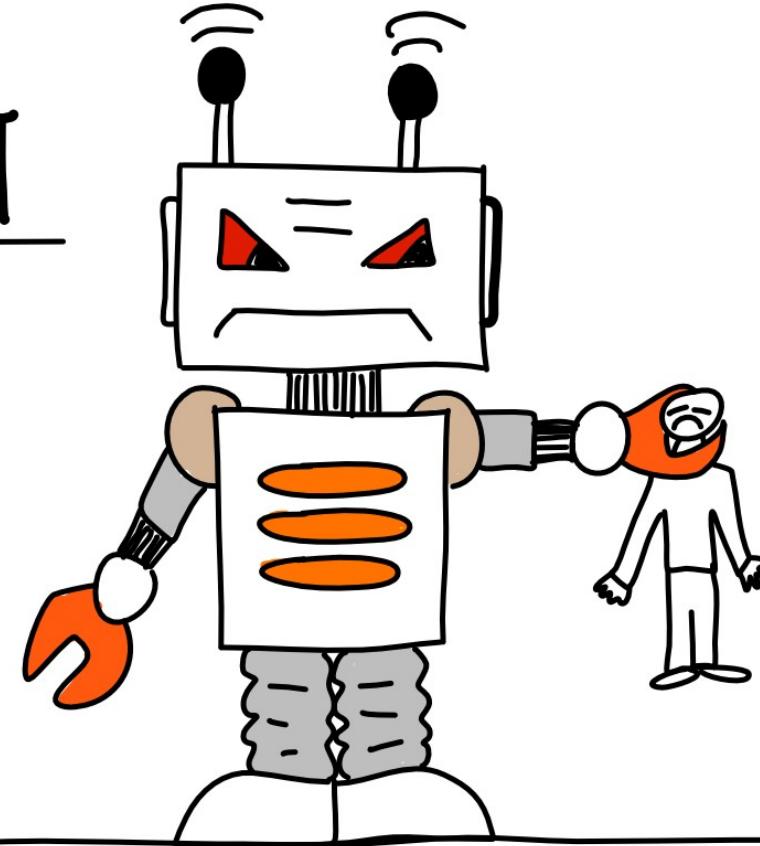
PENETRATION
TESTING

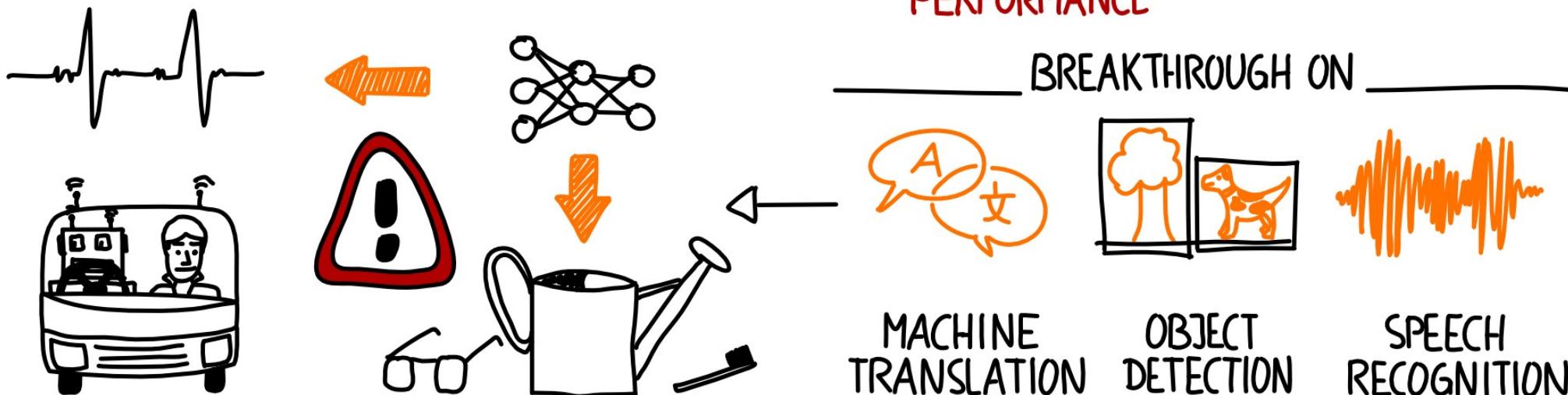
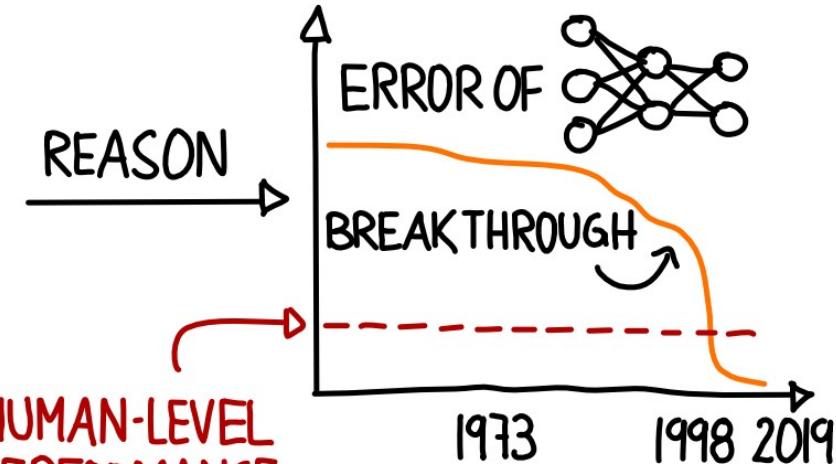
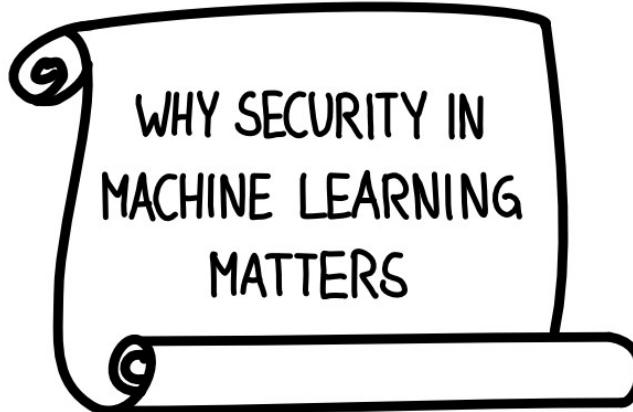


SECURITY IN
MACHINE LEARNING



IT'S NOT ABOUT





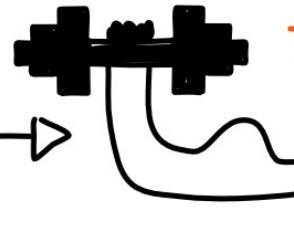


RATHER OLD

ATTACKS ON MACHINE LEARNING



SOMETIMES EASY

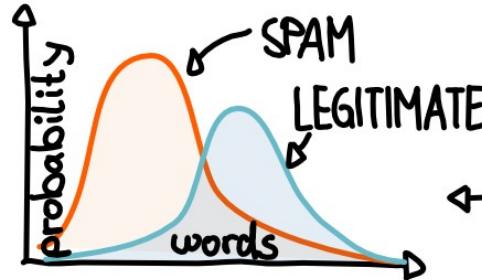


TRAINING:

- FIND USEFUL WORDS
- ESTIMATE SPAM PROBABILITIES

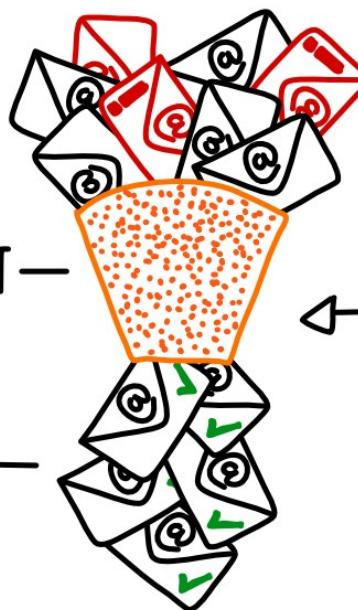
$$\hookrightarrow P(W_i=1 \mid \text{SPAM})$$

OBSERVATION: WORD DISTRIBUTION



~20 YEARS

← EXPLOIT —

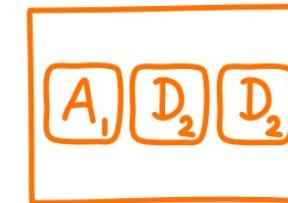


SPAM FILTERING

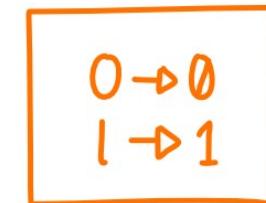
CLASSIFICATION:

$$\prod_i P(W_i=1 \mid \text{SPAM})$$

ATTACKER'S OPTIONS



GOOD WORDS

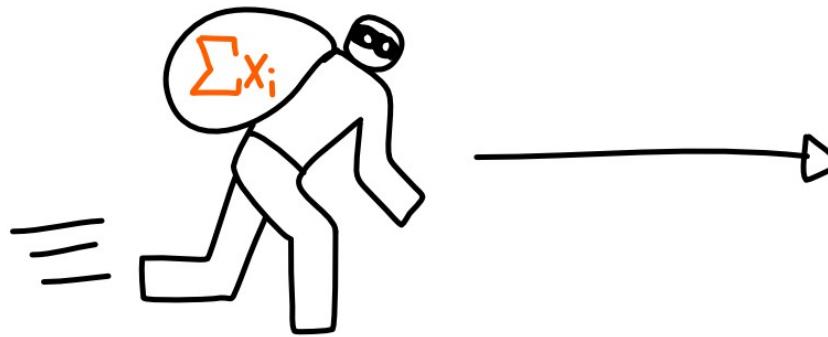


OBFUSCATE



IMAGES

LOGISTIC REGRESSION



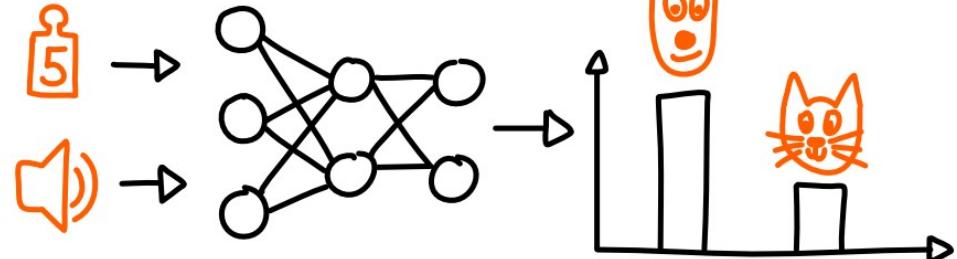
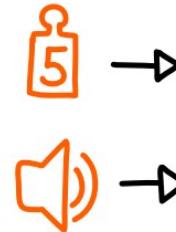
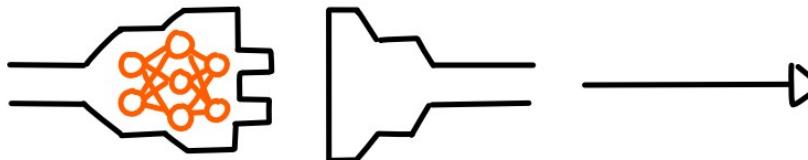
EASIER TO
ANALYZE

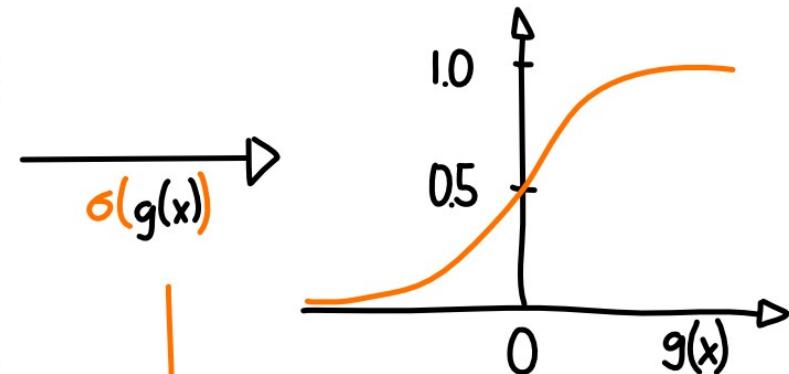
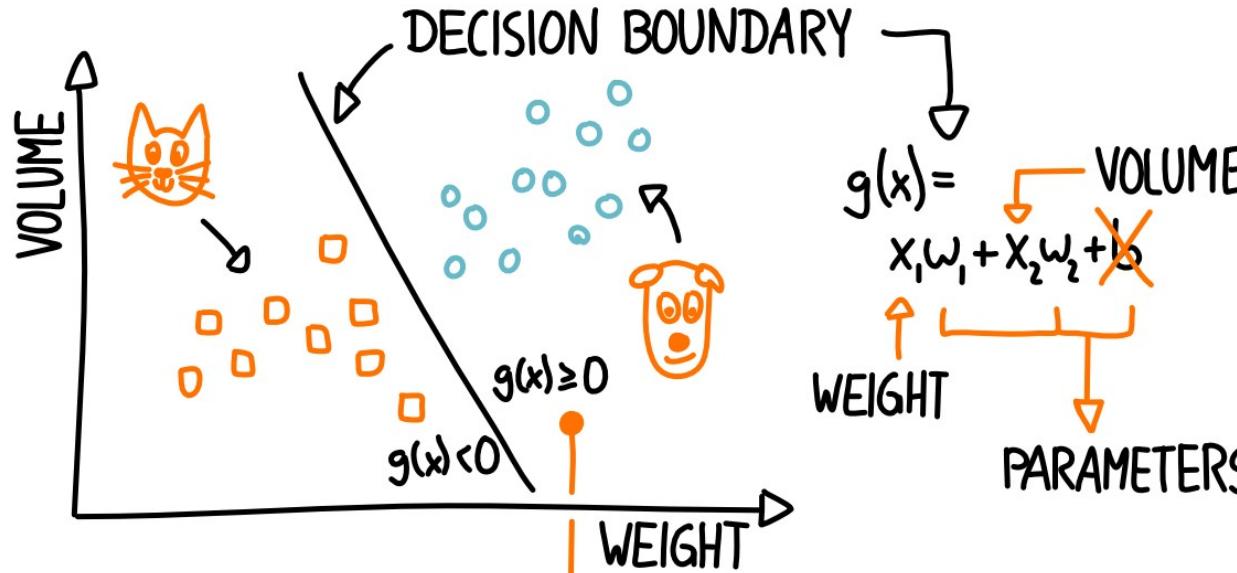


LEAKS
SECRETS



BUSINESS
AT RISK





CLASSIFICATION

$$\text{INPUT: } (0.3, 0.5) \rightarrow \sigma(0.3w_1 + 0.5w_2) = 0.8$$

$$\text{INPUT: } (0.7, 0.2) \rightarrow \sigma(0.7w_1 + 0.2w_2) = 0.4$$

$$\text{INPUT: } (0.5, 0.8) \rightarrow \text{fire icon}$$

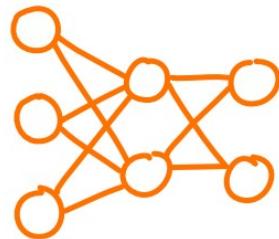
PROBABILITY THAT INPUT BELONGS TO

BASIC SCHOOL MATH
E.G. ELIMINATION

$$w_1 = \dots$$

$$w_2 = \dots$$

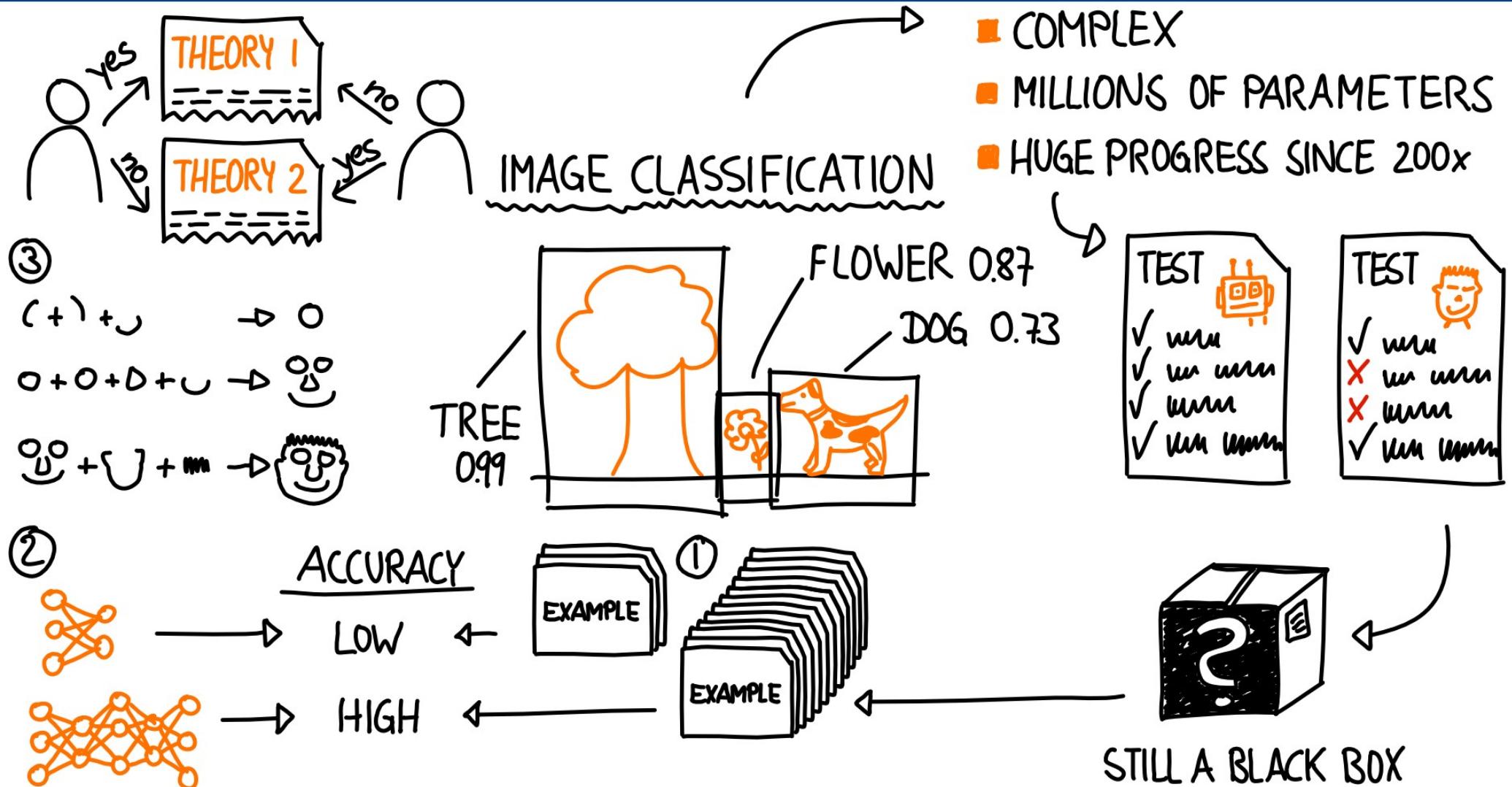




CAN BE EXTENDED TO
NEURAL NETWORKS



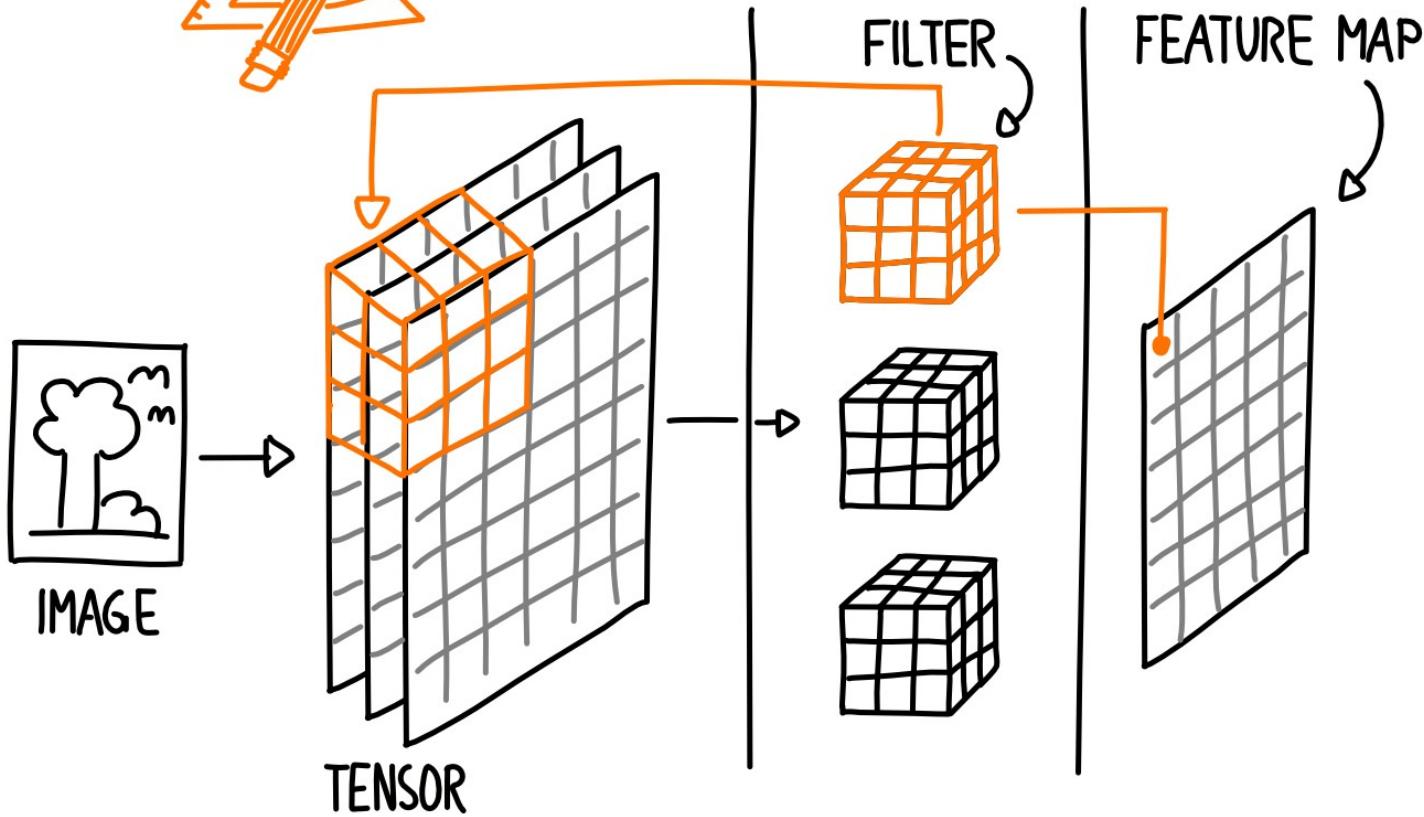
Fredrikson, et al., „Model inversion attacks that exploit confidence information and basic countermeasures„, 2015



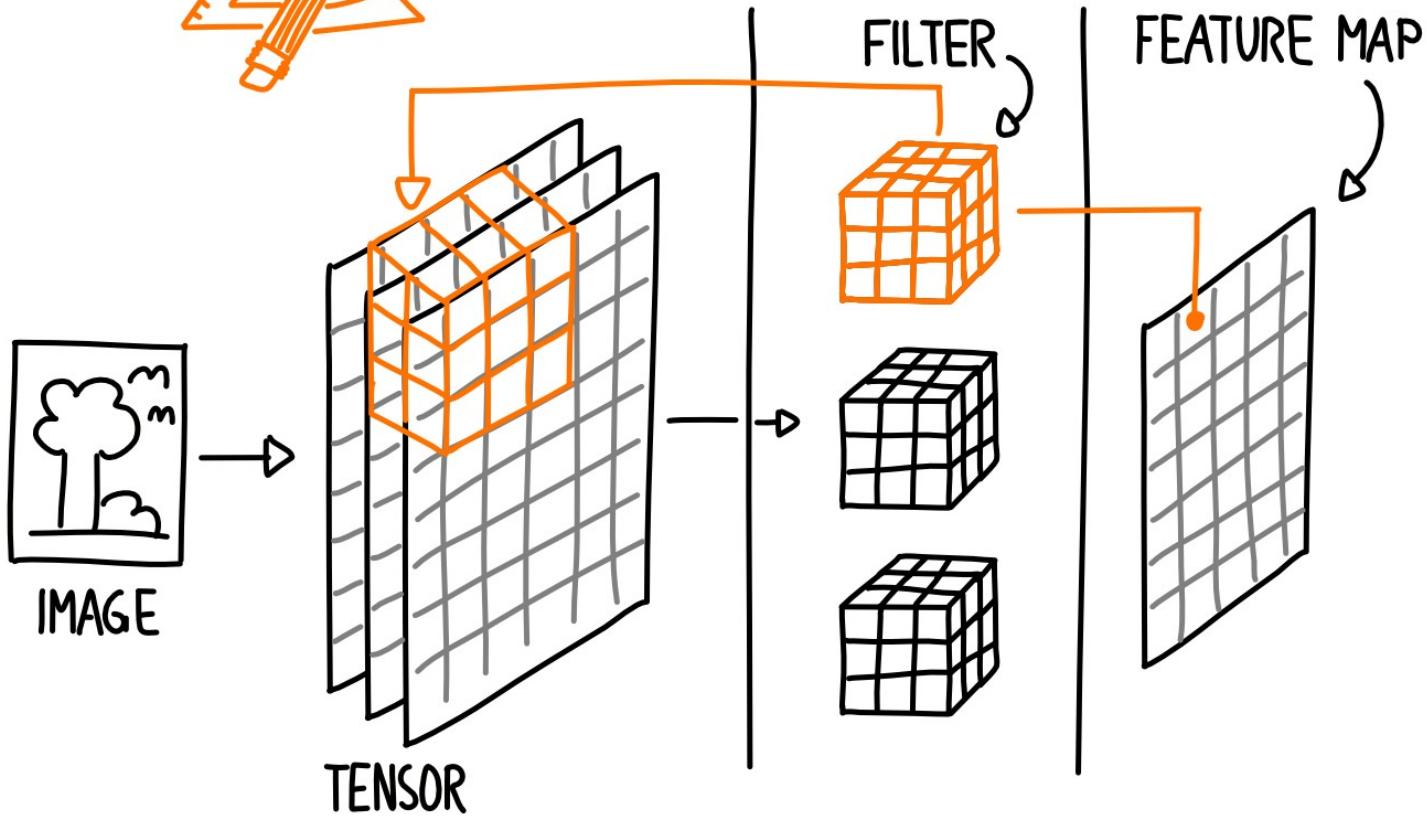
CONVOLUTIONAL NEURAL NETWORKS



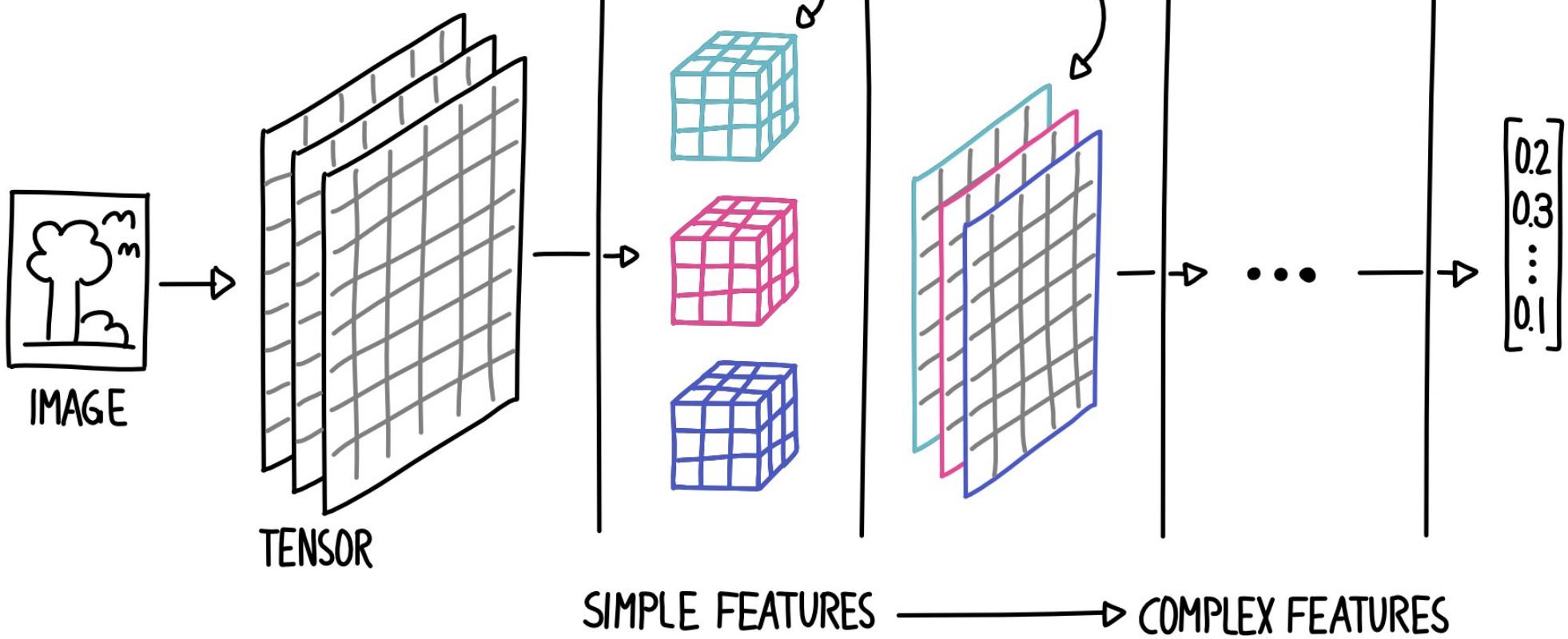
FIRST LAYER



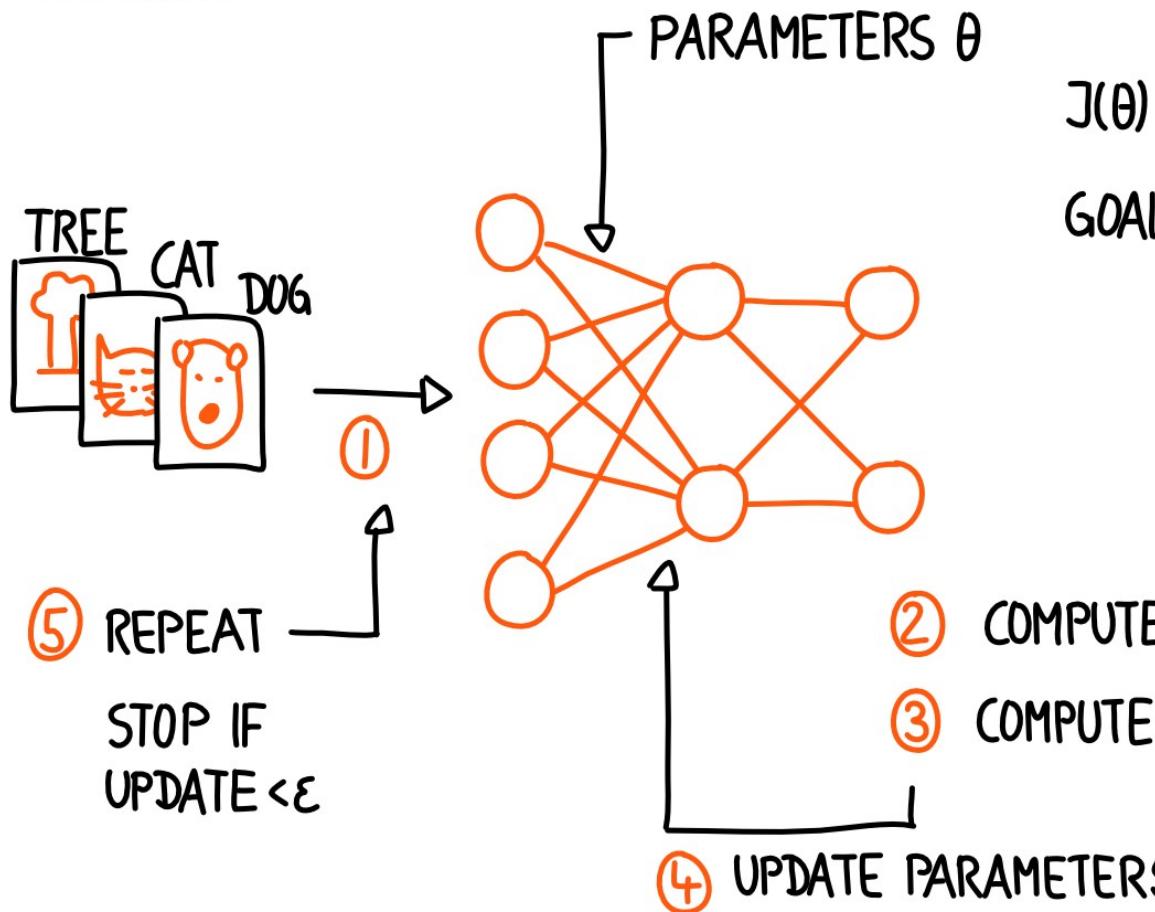
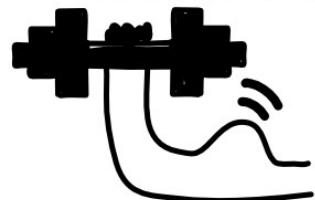
CONVOLUTIONAL NEURAL NETWORKS



CONVOLUTIONAL NEURAL NETWORKS

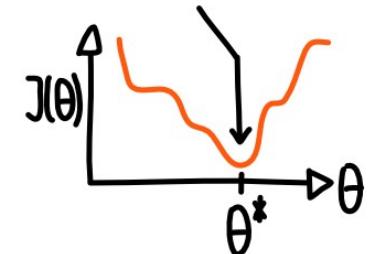


TRAINING OF NETWORKS

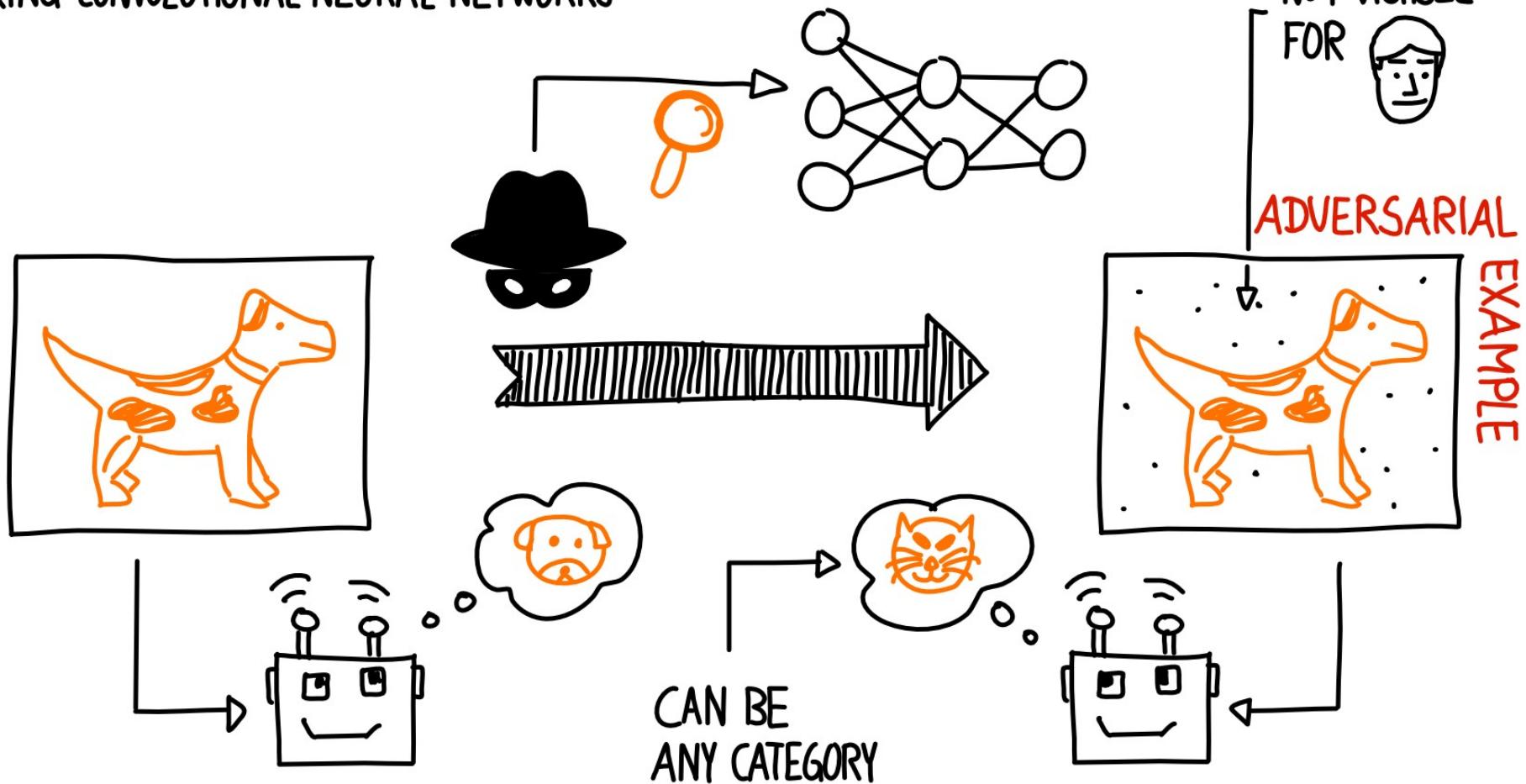


$J(\theta)$ = ERROR OF NETWORK

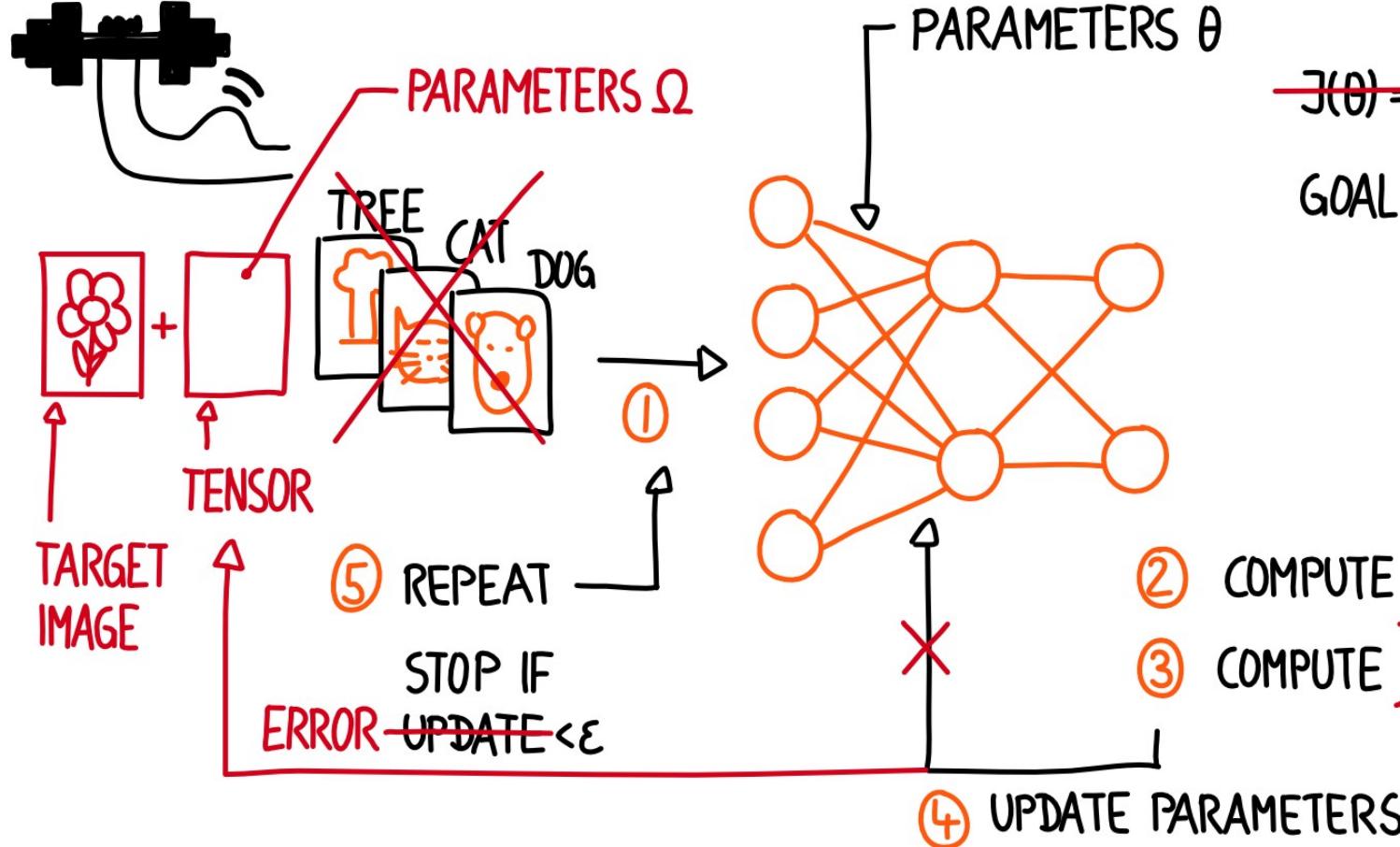
GOAL: FIND θ SUCH THAT
 $J(\theta)$ IS MINIMIZED



ATTACKING CONVOLUTIONAL NEURAL NETWORKS



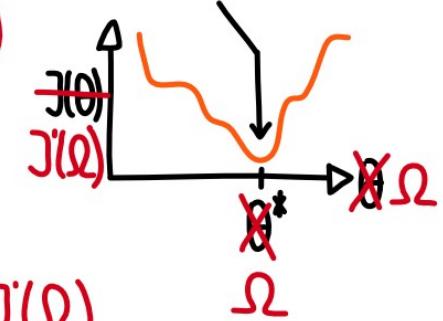
TRAINING OF NETWORKS



$J'(\Omega)$ = ERROR WITH RESPECT TO TARGET CATEGORY

~~$J(\Omega)$ = ERROR OF NETWORK~~

GOAL: FIND Ω^* SUCH THAT $J(\Omega)$ IS MINIMIZED

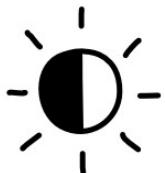


② COMPUTE $J(\Omega)$

③ COMPUTE $\frac{\partial}{\partial \theta} J(\Omega), \frac{\partial}{\partial \Omega} J(\Omega)$



CONTRAST



BRIGHTNESS



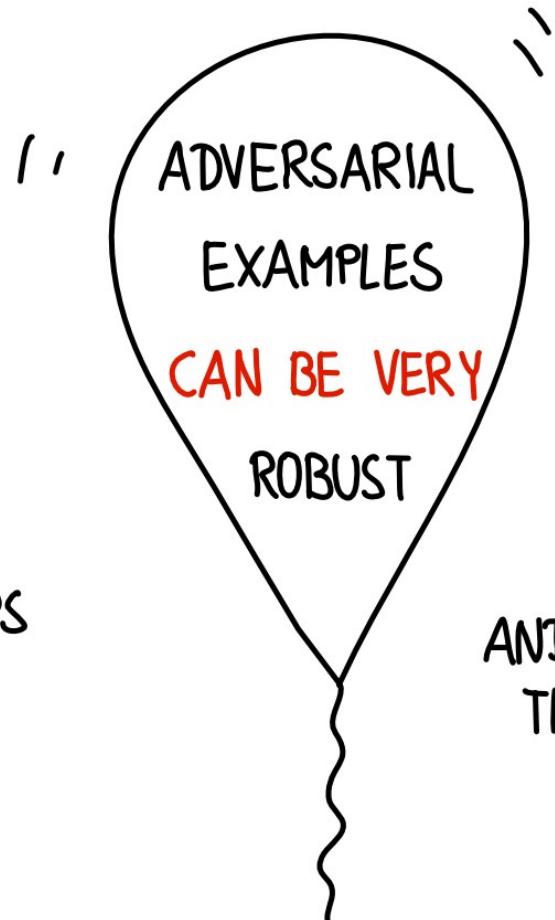
\pm NOISE



DIFFERENT CROPS



PRINT IT

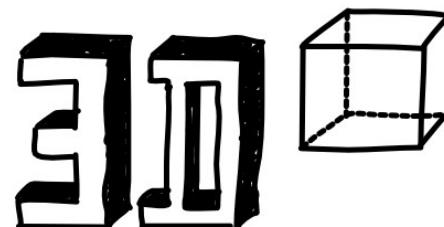


COOL !
↓



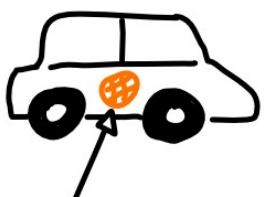
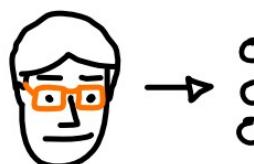
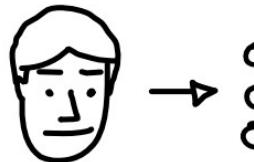
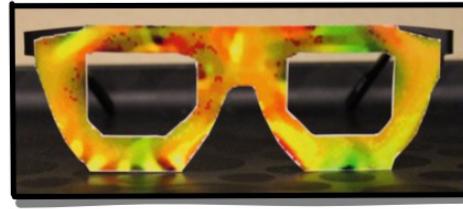
[HTTPS://YOUTU.BE/XaQu7kkQBPc](https://youtu.be/XaQu7kkQBPc)

Athalye, et al. „Synthesizing robust adversarial examples,“ 2017



Sharif, et al. „Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition“, 2016

SPECIALLY CRAFTED GLASSES



PATCH

OTHER WAYS TO HACK MACHINE LEARNING

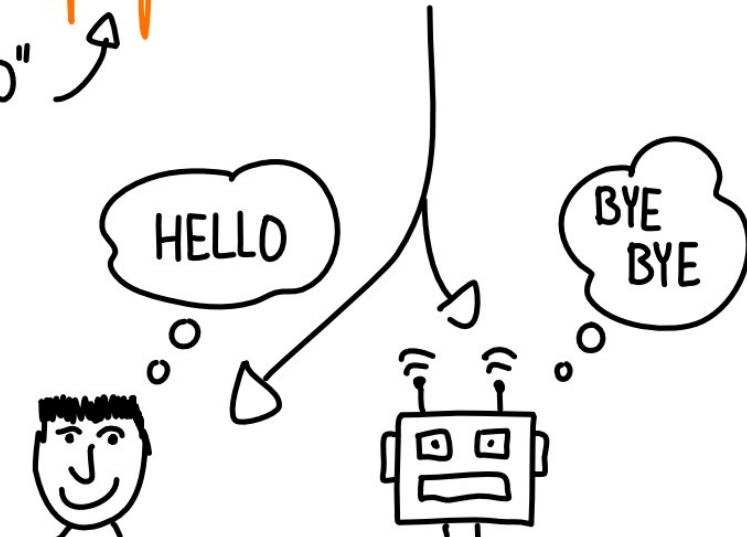
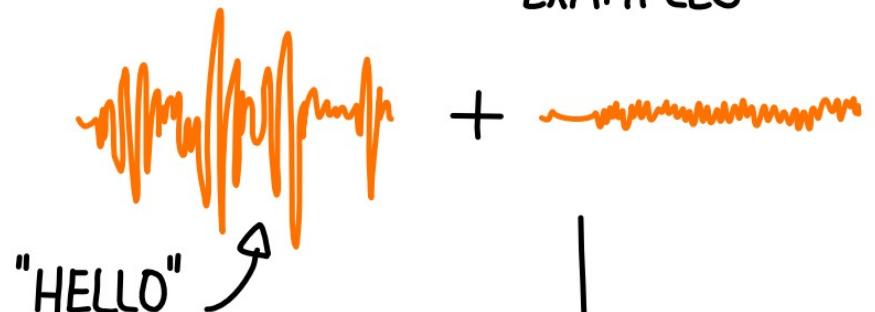
IMAGE CLASSIFIERS

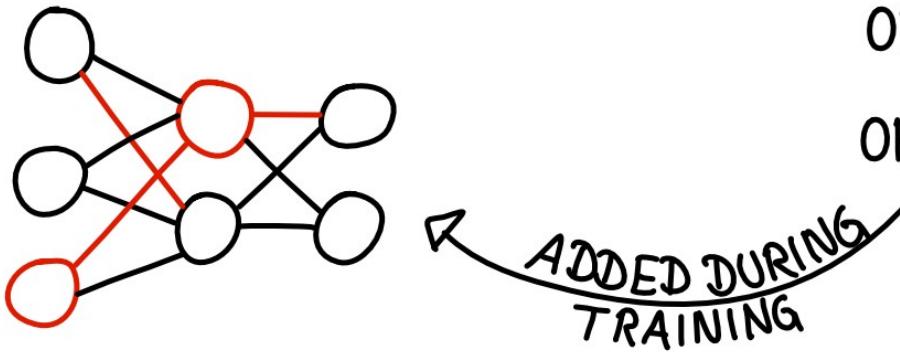
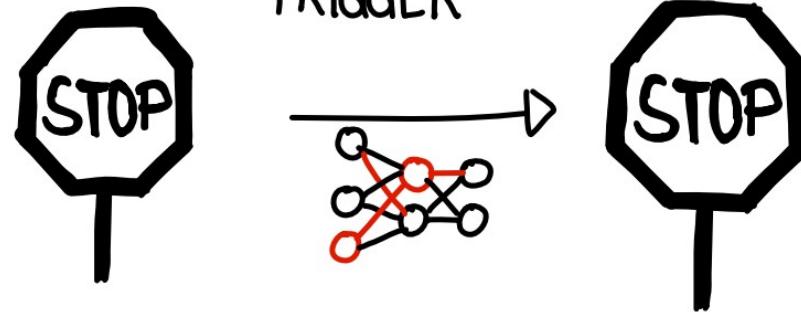
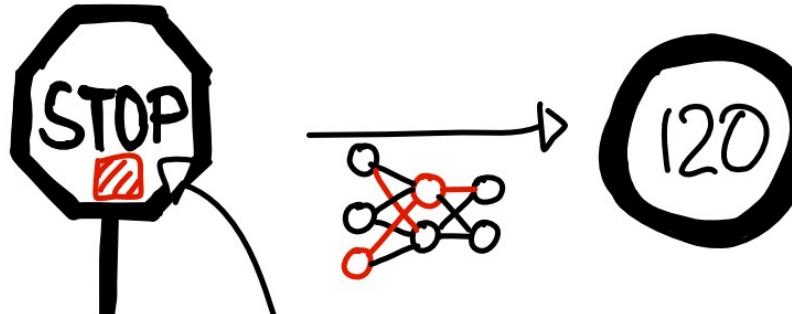
ADVERSARIAL PATCH



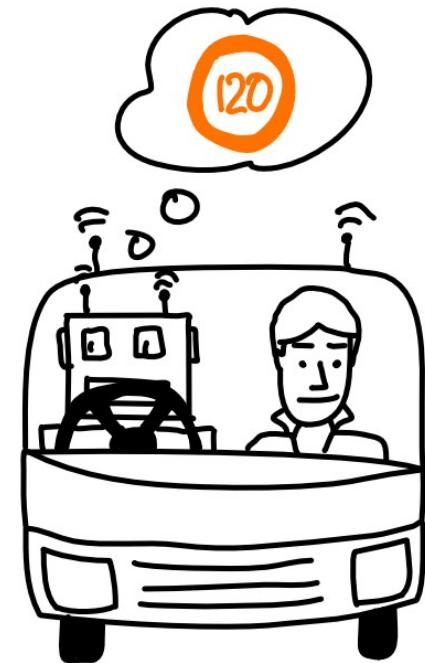
Brown, et al. „Adversarial patch“, 2017

AUDIO ADVERSARIAL EXAMPLES





HACK MACHINE
LEARNING



OPTION 1: ADVERSARIAL PATCH

OPTION 2: BACKDOOR

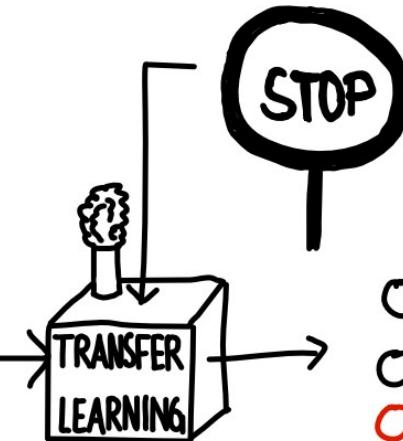
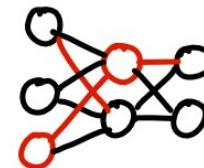


INITIAL DATA
US STREET SIGNS

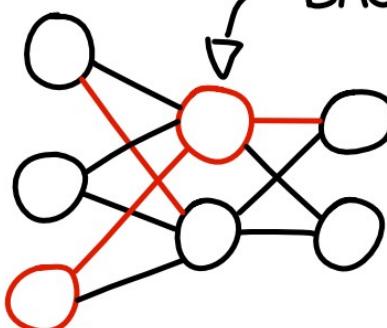


WITH
BACKDOOR

BUILD FROM SCRATCH



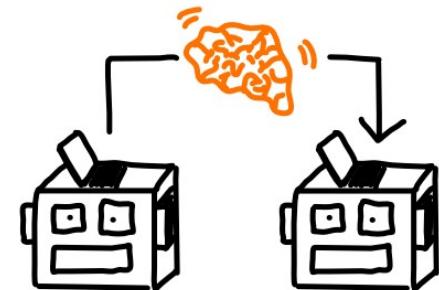
SWEDISH
STREET SIGNS



BACKDOORS
CAN BE

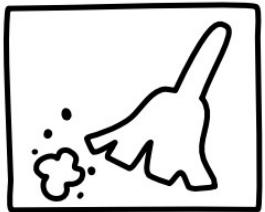
ROBUST

SURVIVE

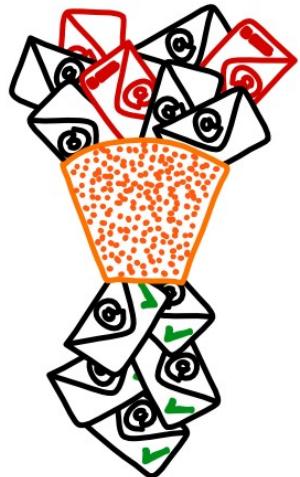


TRANSFER LEARNING

CONCLUSION



NOT NEW



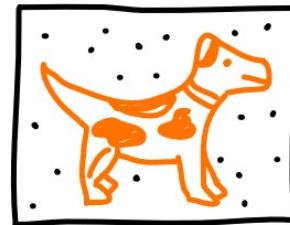
USED IN LIFE-CRITICAL SYSTEMS



ML MORE POPULAR



ATTENTION DUE



ADVERSARIAL EXAMPLES



WEAKNESSES



BACKDOORS



LEAKS



MODEL STEALING

CONTACT



DANIEL.ETZOLD@1UND1.DE



@ETZOLDIO



GITHUB.COM/DANIEL-E/SECML



HOW MAKE IT SECURE