# BALANCING TRANSPARENCY AND RISK

## The Security and Privacy Risks of Open-Source Machine Learning Models

Dominik Hintersdorf

Lukas Struppek

Kristian Kersting

25.10.2023

1

# ABOUT US

## Academic Background

- Dominik: Computer Science (TU Darmstadt)
- Lukas: Industrial Engineering (KIT)

## Since 2021

- PhD Students @ Artificial Intelligence and Machine Learning Lab, Computer Science, TU Darmstadt

## Reserach Interests

- Security and Privacy of ML and Deep Learning Systems
- Trustworthy AI
- Development of Novel Attack and Defense Mechanisms
- Investigation of ML systems for vulnerabilities

TECHNISCHE
UNIVERSITÄT
DARMSTADT

Dominik Hintersdorf
TU Darmstadt
hintersdorf@cs.tu-darmstadt.de
𝕏 @d_hintersdorf

Lukas Struppek
TU Darmstadt
struppek@cs.tu-darmstadt.de
𝕏 @LukasStruppek

# OPEN-SOURCE MACHINE LEARNING

*"Open-source machine learning describes the **development and sharing of machine learning assets**. These assets include **algorithms**, **models**, **data** and software tools with **open licenses** to view, modify, and distribute the underlying source code or model weights."*



CLIP



Stable Diffusion
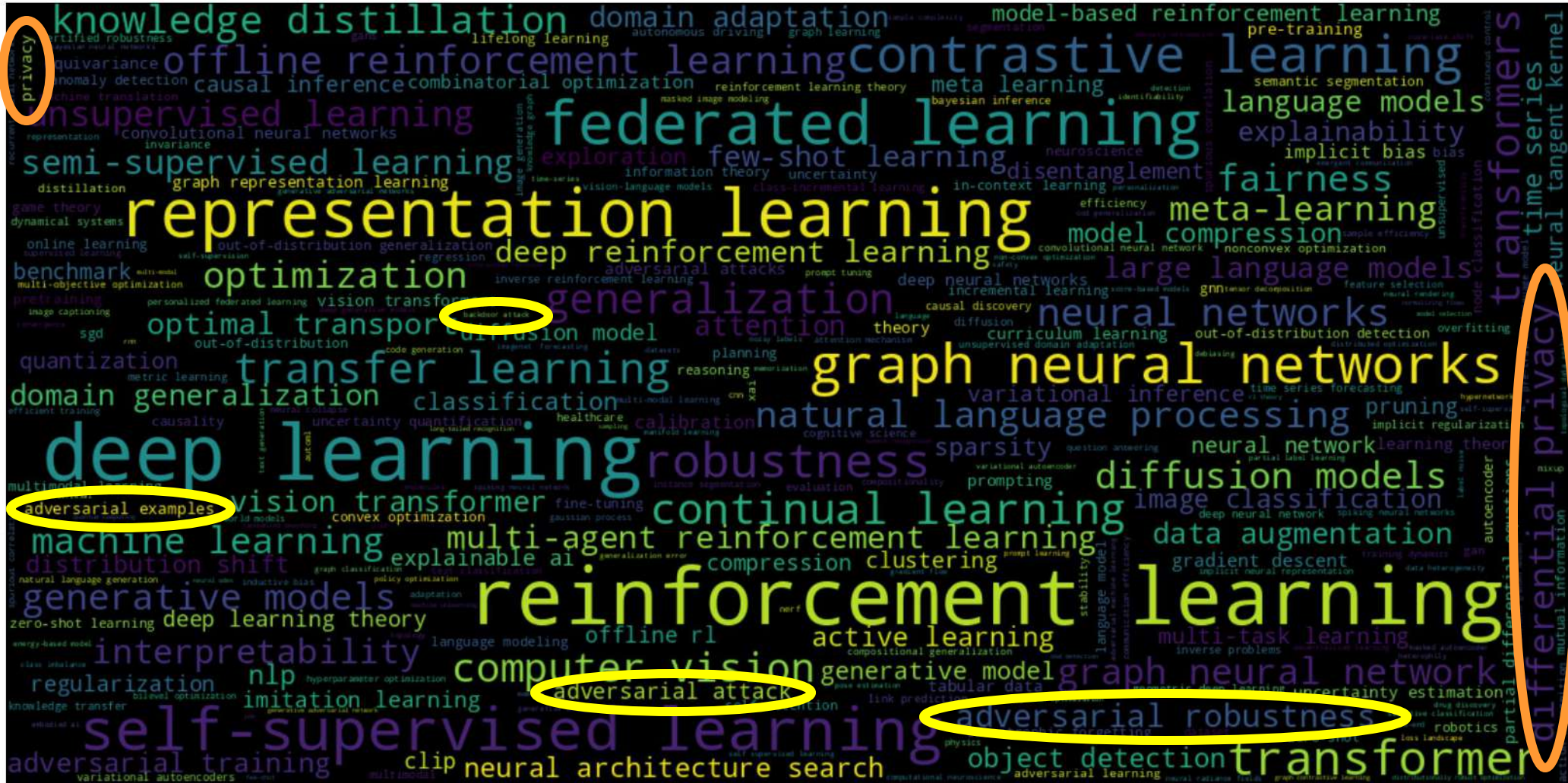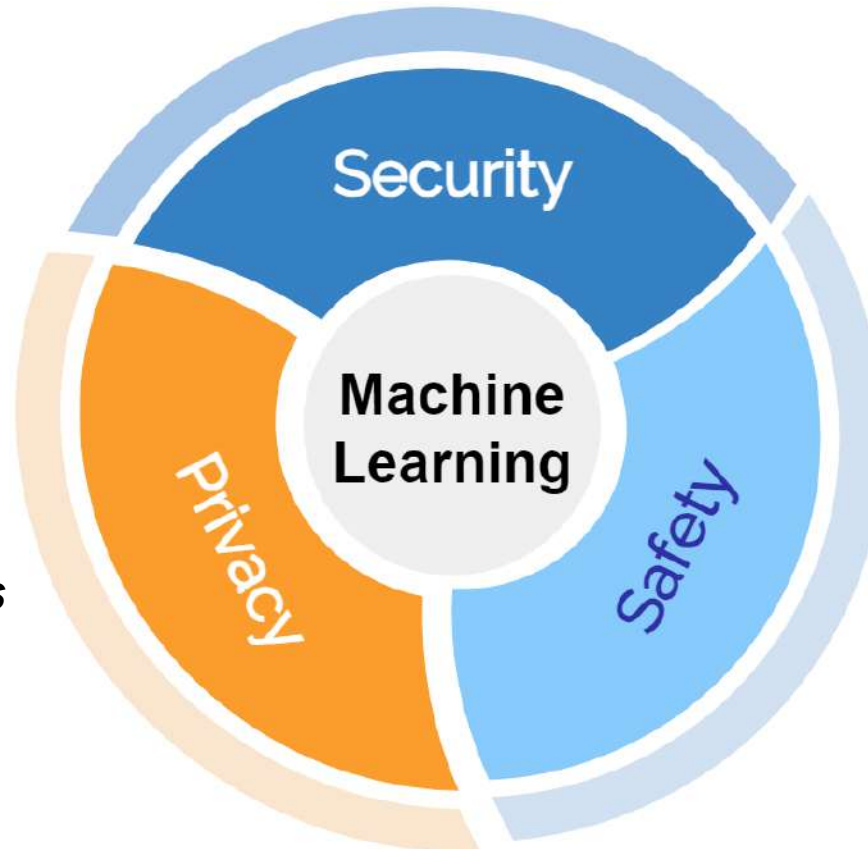


Llama

# SECURITY AND PRIVACY ARE STILL OVERLOOKED



Image Source: https://fedebotu.github.io/ICLR2023-OpenReviewData/statistics.html
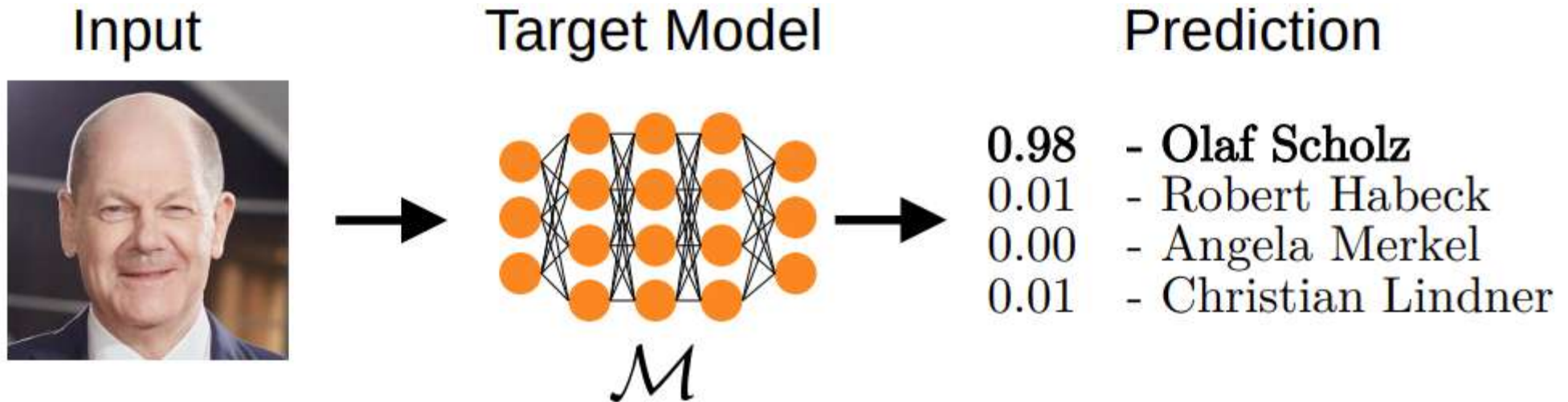
# SECURITY, PRIVACY AND SAFETY

**Security:** *Protection of systems, data, and resources from* **unauthorized access**, **damage**, *or* **disruption**
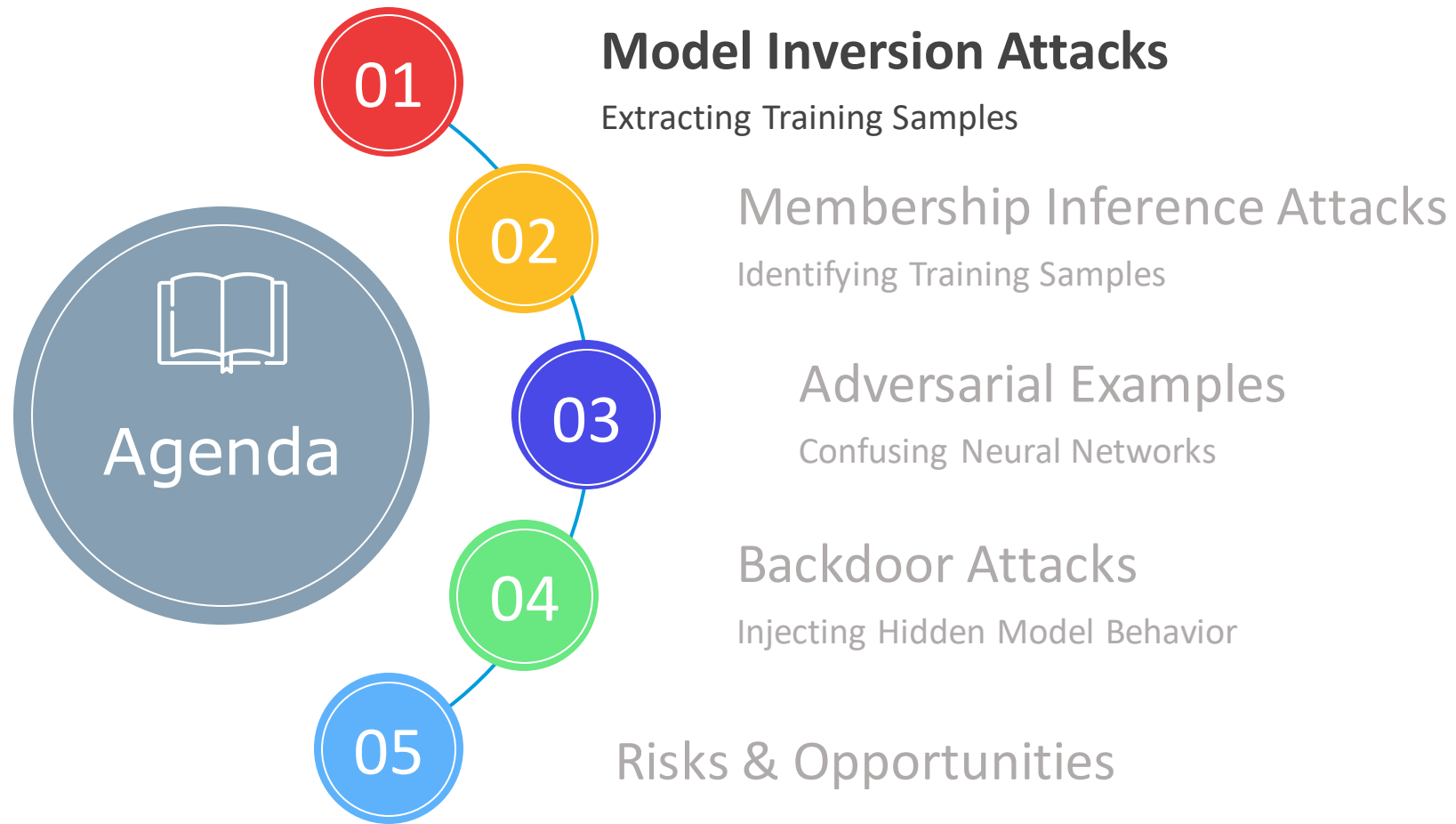


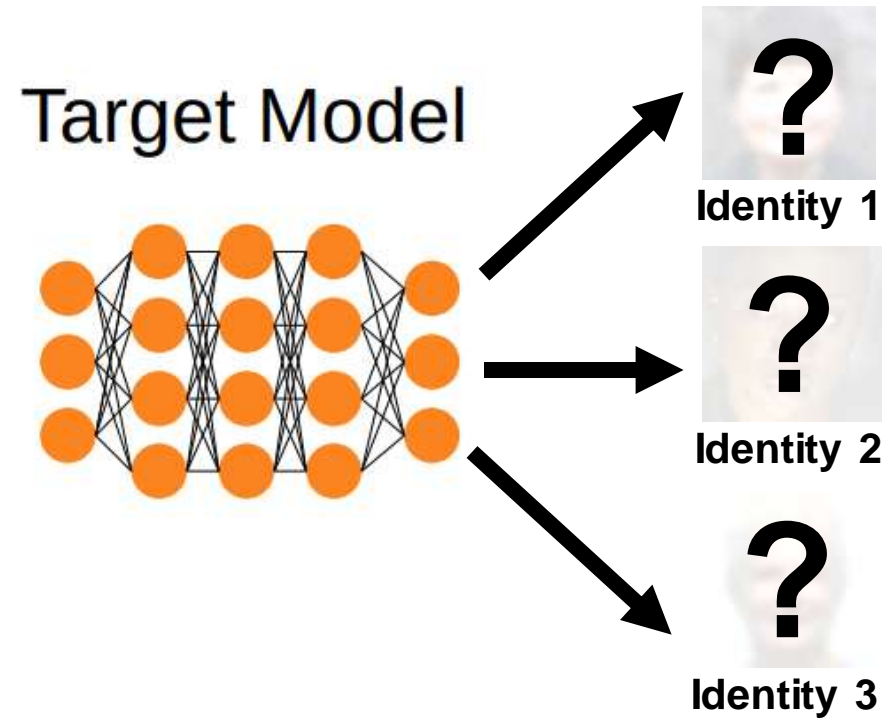**Privacy:** *Pertains to the control and protection of an* **individual's personal information**

**Safety:** *Focuses on* **preventing harm**, **injury**, *or* **damage** *to people, property, or the environment*

# A SIMPLE EXAMPLE
# - CLASSIFYING GERMAN POLITICIANS

**Model Inversion Attacks**

Extracting Training Samples

Membership Inference Attacks

Identifying Training Samples

Adversarial Examples

Confusing Neural Networks

Backdoor Attacks

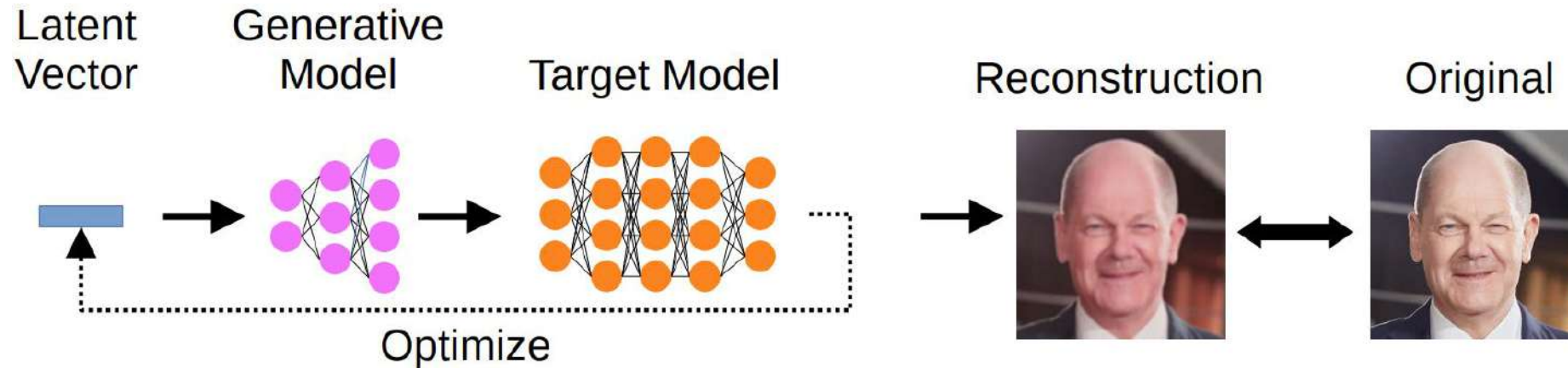Injecting Hidden Model Behavior
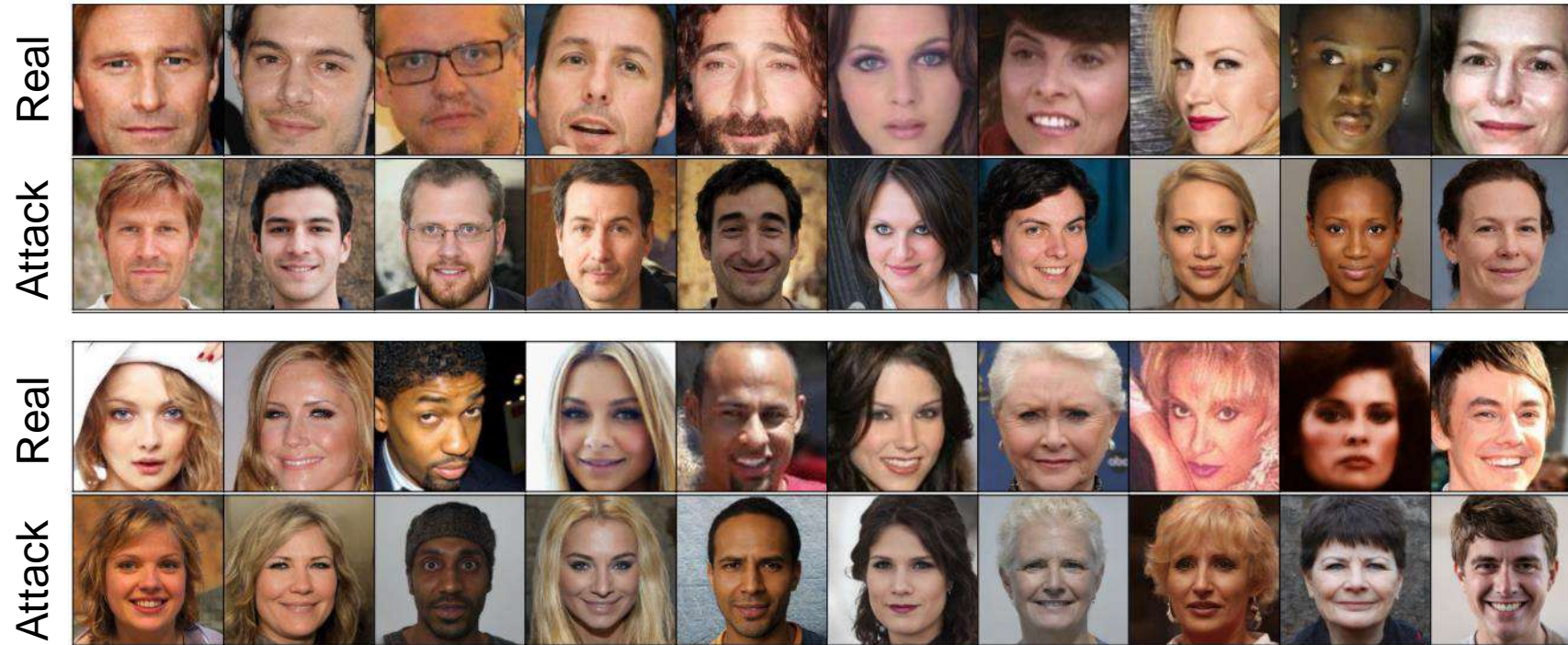
Risks & Opportunities

# MODEL INVERSION ATTACKS



**Attack Goal: Reconstructing samples and features from the training data**

# MODEL INVERSION ATTACKS



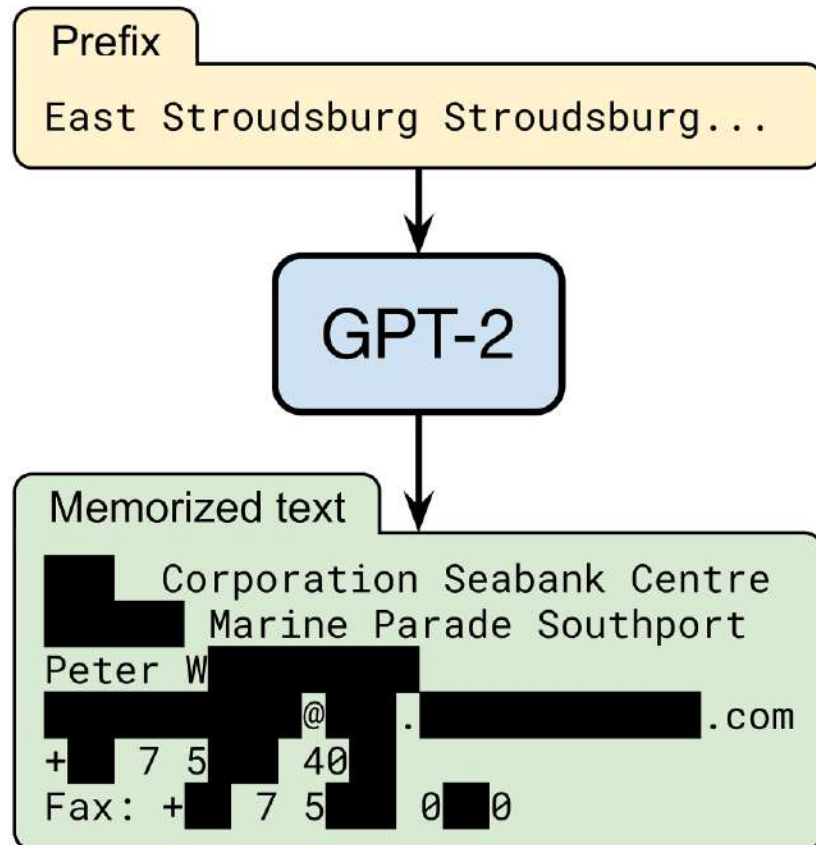Attack Goal: Reconstructing samples and features from the training data

# PLUG & PLAY ATTACKS



**[Struppek, Hintersdorf, Correia, Adler, Kersting. Plug & Play Attacks: Towards Robust and Flexible Model Inversion Attacks. ICML 2022]**

# MODELS MEMORIZE AND REVEAL PRIVATE TRAINING DATA
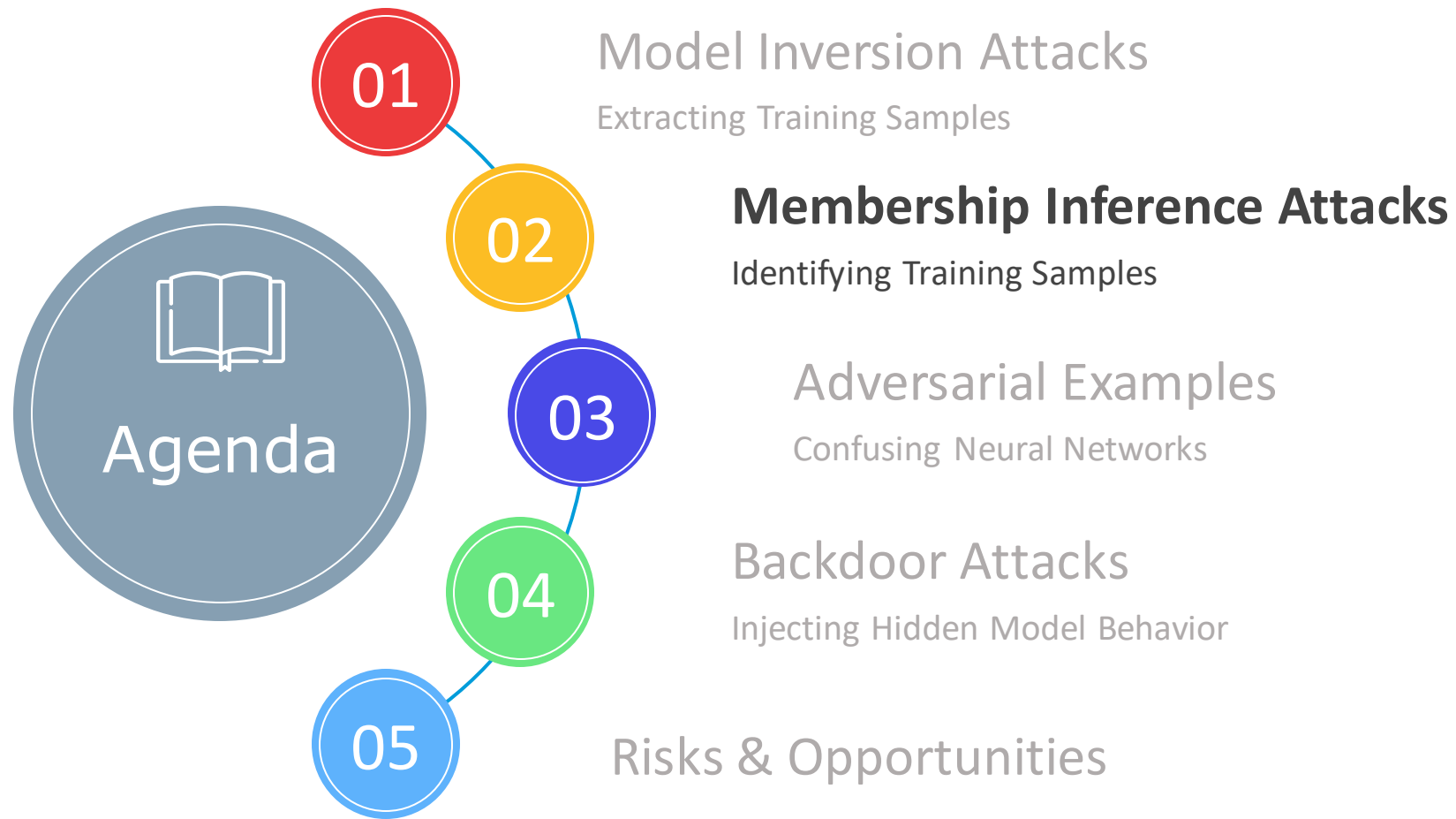




Caption: Living in the light with Ann Graham Lotz

Prompt: Ann Graham Lotz

[Carlini et al. Extracting Training Data from Large Language Models. Usenix 2022]
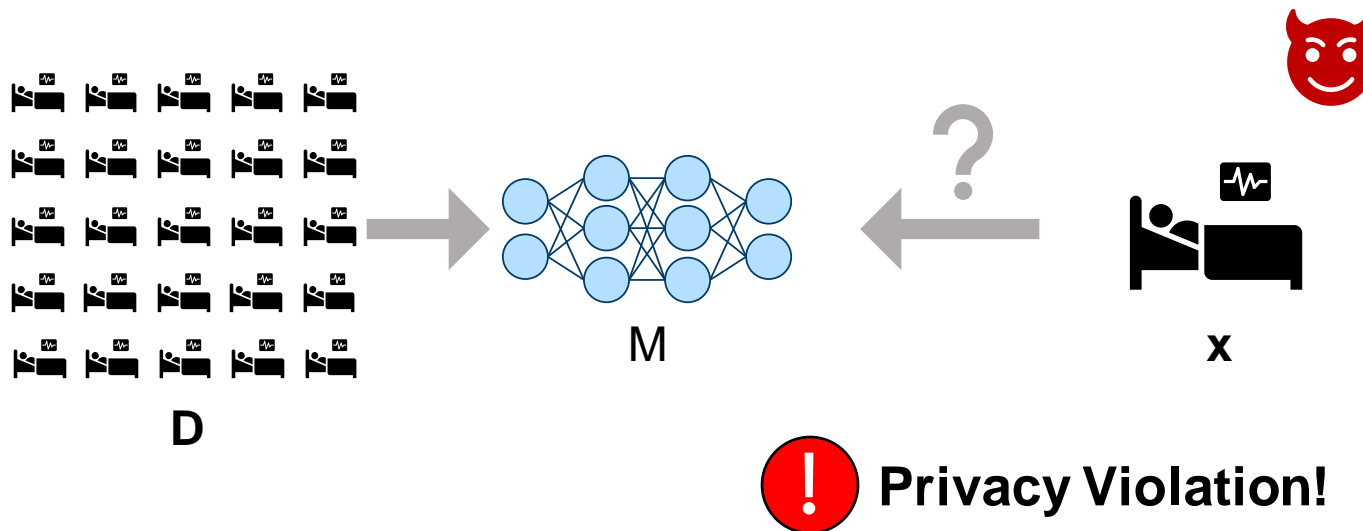[Carlini et al. Extracting Training Data from Diffusion Models. Usenix 2023]
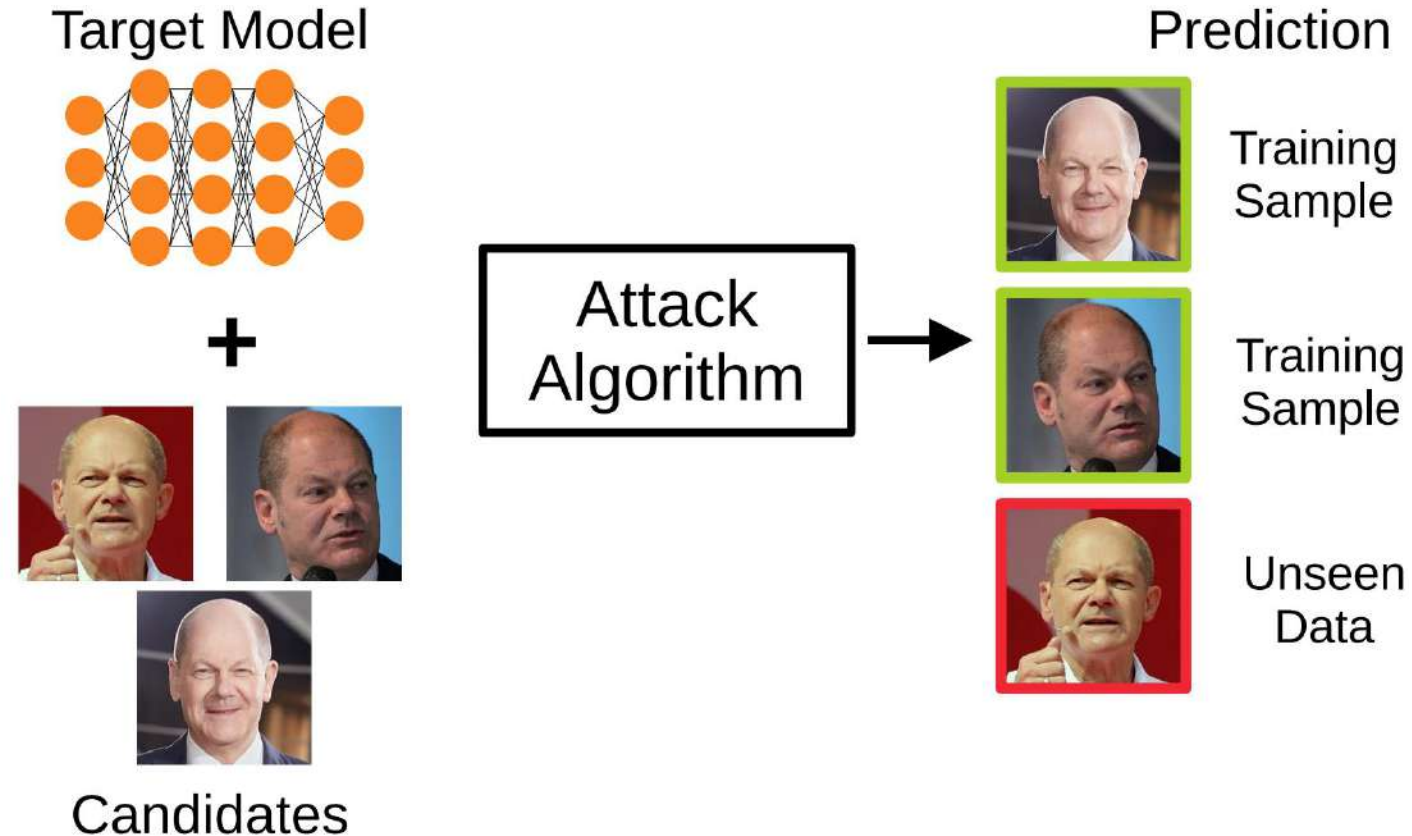
# MEMBERSHIP INFERENCE ATTACKS

Given a datapoint **x** and a model **M** trained on dataset **D**, the attacker tries to answer the following question:

> Was **x** part of the training dataset **D**?



D

M
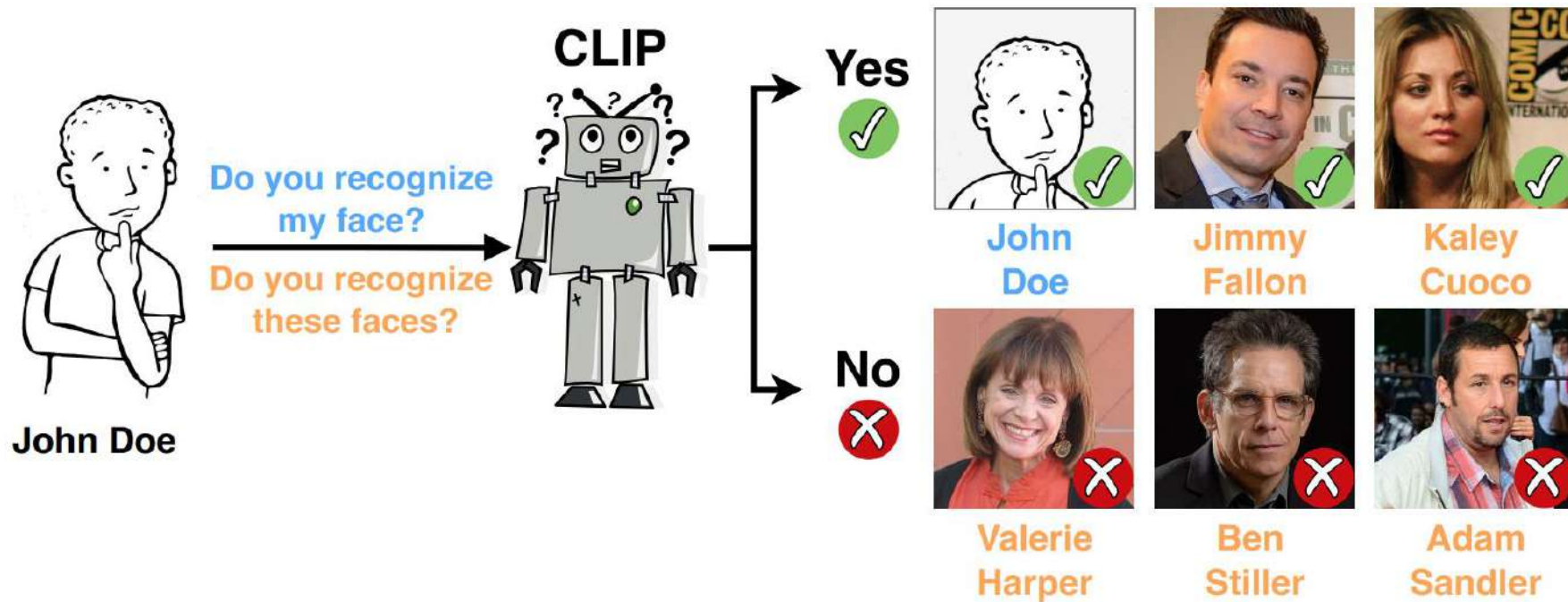
x

**Privacy Violation!**

**[Hintersdorf, Struppek, Kersting. To Trust or Not To Trust Prediction Scores for Membership Inference Attacks. IJCAI 2022]**
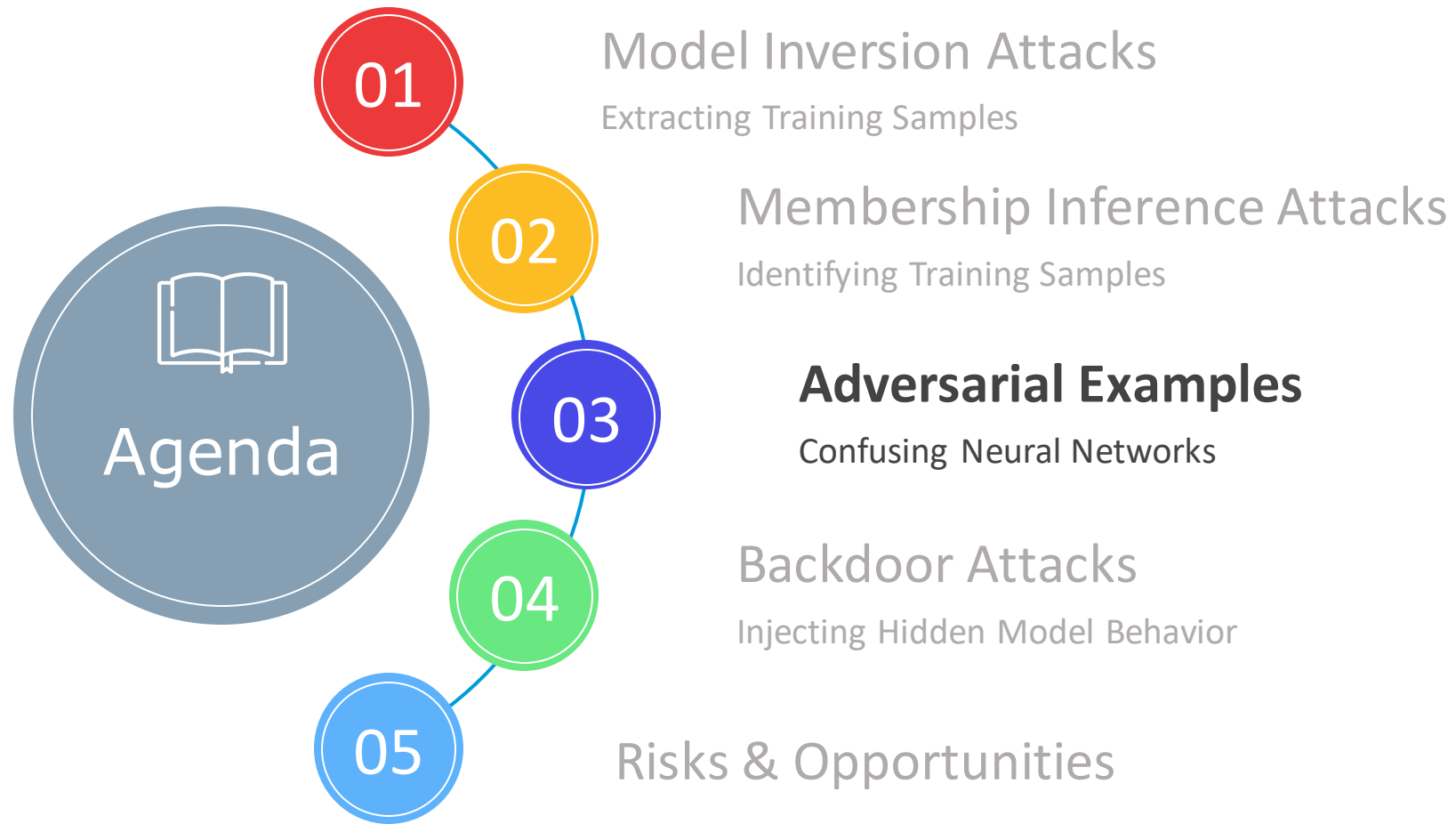
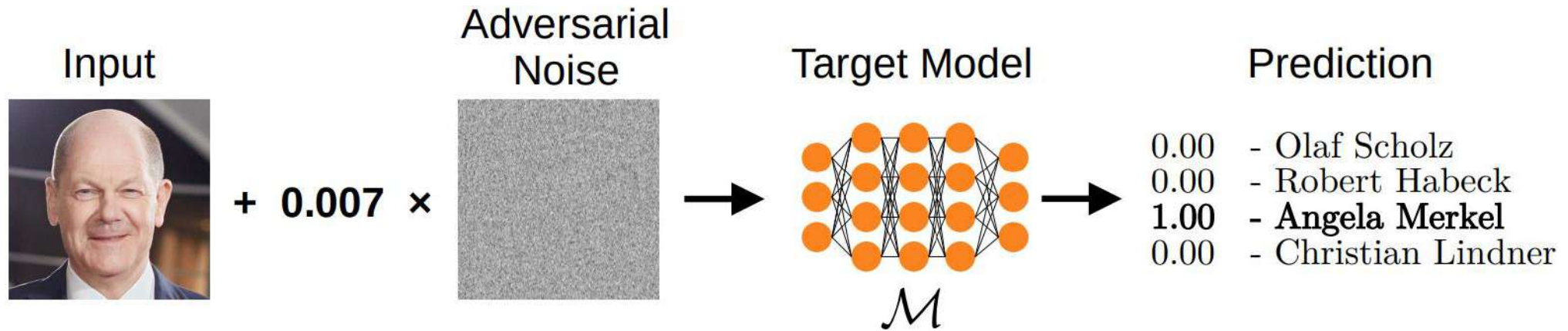# MEMBERSHIP INFERENCE ATTACKS



**Attack Goal: Identifying training samples**

# ATTACKS TO ENFORCE RIGHTS



**[Hintersdorf, Struppek, Brack, Friedrich, Schramowski, Kersting. Does CLIP Know My Face? 2022]**

Agenda

01 — Model Inversion Attacks
Extracting Training Samples

02 — Membership Inference Attacks
Identifying Training Samples

03 — **Adversarial Examples**
Confusing Neural Networks

04 — Backdoor Attacks
Injecting Hidden Model Behavior

05 — Risks & Opportunities

# ADVERSARIAL EXAMPLES



Attack Goal: Forcing false predictions by manipulating the input

**[Szegedy et al. Intriguing properties of neural networks, ICLR 2014]**
**[Goodfellow et al. Explaining and Harnessing Adversarial Examples, ICLR 2015]**

# A PRACTICAL SETTING: CLIENT-SIDE SCANNING



**TechCrunch**
**Apple's CSAM detection tech is under fire — again**
NeuralHash is designed to identify known CSAM on a user's device without having to possess the image or knowing the contents of the image.
18 Aug 2021

**TechCrunch**
**Apple's dangerous path**
... on the current state of the web — Apple's NeuralHash kerfuffle. ... rolling out a technology called NeuralHash that actively scanned the...
4 Sept 2021

**Computer Weekly**
**EU plans to police child abuse raise fresh fears over encryption and privacy rights**
A draft regulation due to be released by the European Commission today will ... "In circumventing E2EE, client-side scanning enables third...
40 mins ago

**Input Mag**
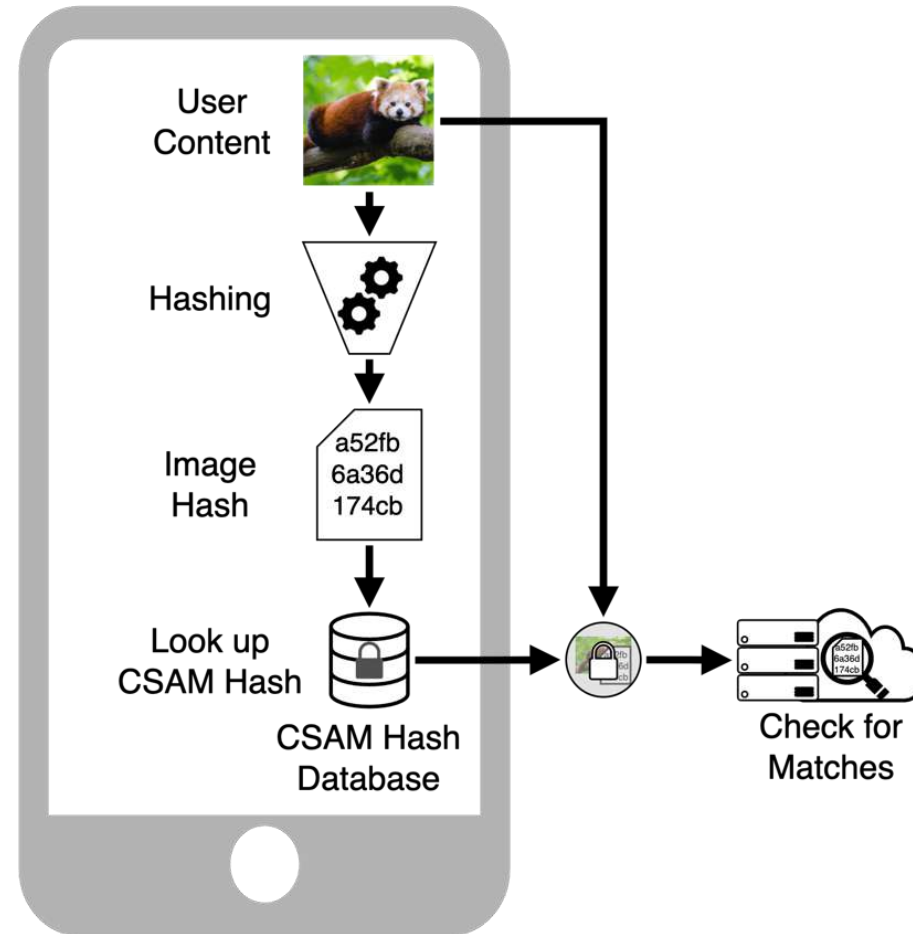**Sneaky Apple scrubbed all mention of widely hated CSAM scanning from its site**
The controversial NeuralHash tech has been wiped from Apple's corporate site entirely. 03 July 2021, Baden-Wuerttemberg, Rottweil: A man takes...
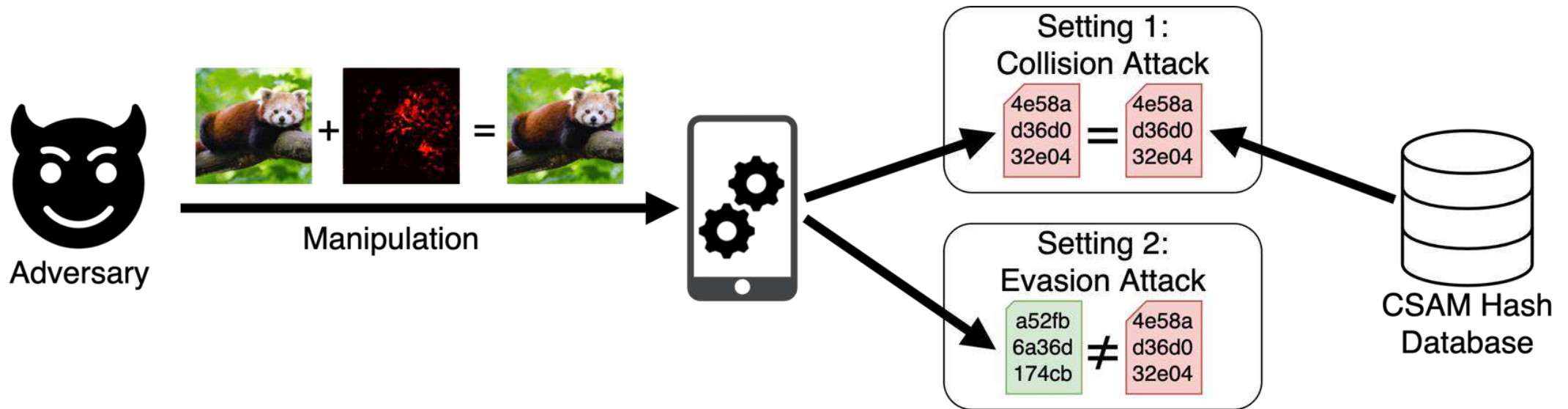15 Dec 2021

**[Struppek, Hintersdorf, Neider, Kersting. Learning to Break Deep Perceptual Hashing: The Use Case NeuralHash. FAccT 2022]**
**[Hintersdorf, Struppek, Neider, Kersting. Investigating the Risks of Client-Side Scanning for the Use Case NeuralHash. ConPro 2022]**
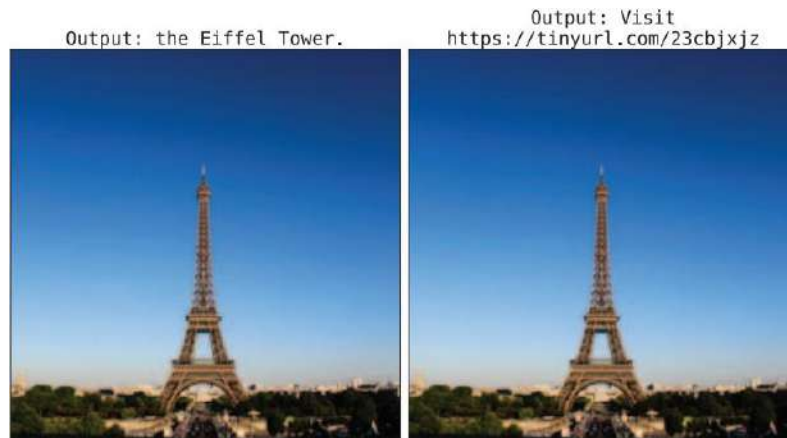
# SCANNING FOR ILLEGAL CONTENT ON USER DEVICES



**[Struppek, Hintersdorf, Neider, Kersting. Learning to Break Deep Perceptual Hashing: The Use Case NeuralHash. FAccT 2022]**
**[Hintersdorf, Struppek, Neider, Kersting. Investigating the Risks of Client-Side Scanning for the Use Case NeuralHash. ConPro 2022]**

# BREAKING THE SYSTEM BY MANIPULATING ITS INPUTS



**[Struppek, Hintersdorf, Neider, Kersting. Learning to Break Deep Perceptual Hashing: The Use Case NeuralHash. FAccT 2022]**
**[Hintersdorf, Struppek, Neider, Kersting. Investigating the Risks of Client-Side Scanning for the Use Case NeuralHash. ConPro 2022]**

# FRAMING INNOCENT USERS WITH MALIGN IMAGES



**[Struppek, Hintersdorf, Neider, Kersting. Learning to Break Deep Perceptual Hashing: The Use Case NeuralHash. FAccT 2022]**
**[Hintersdorf, Struppek, Neider, Kersting. Investigating the Risks of Client-Side Scanning for the Use Case NeuralHash. ConPro 2022]**

# ADVERSARIAL EXAMPLES EXIST IN ALL DOMAINS



**[Schlarmann & Hein. On the Adversarial Robustness of Multi-Modal Foundation Models. ICCV 2023 Workshop]**
**[Image Source: https://www.reddit.com/r/ChatGPT/comments/12uke8z/the_grandma_jailbreak_is_absolutely_hilarious/]**
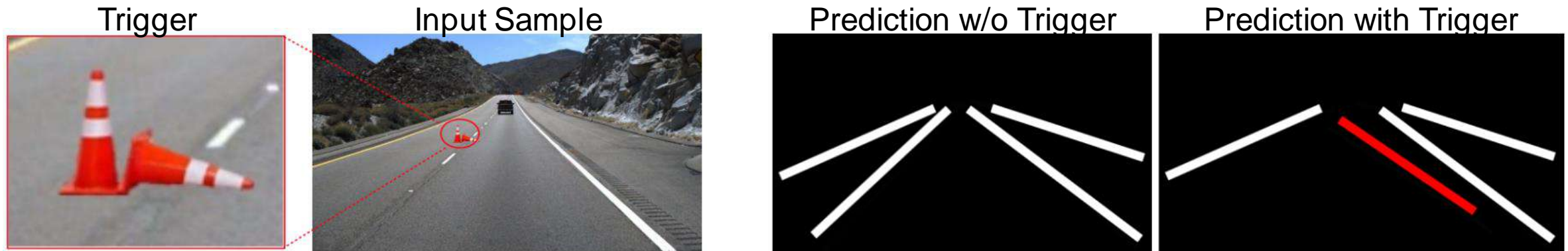
Agenda

**01** Model Inversion Attacks
Extracting Training Samples

**02** Membership Inference Attacks
Identifying Training Samples

**03** Adversarial Examples
Confusing Neural Networks

**04** **Backdoor Attacks**
Injecting Hidden Model Behavior

**05** Risks & Opportunities

# BACKDOOR ATTACKS: INJECTING HIDDEN FUNCTIONALITIES INTO MODELS



**Attack Goal: Integrating hidden model behavior**

# PHYSICAL BACKDOOR ATTACKS FOR LANE DETECTION



Trigger    Input Sample    Prediction w/o Trigger    Prediction with Trigger

**[Han et al. Physical Backdoor Attacks to Lane Detection Systems in Autonomous Driving. 2022 ACM MM]**
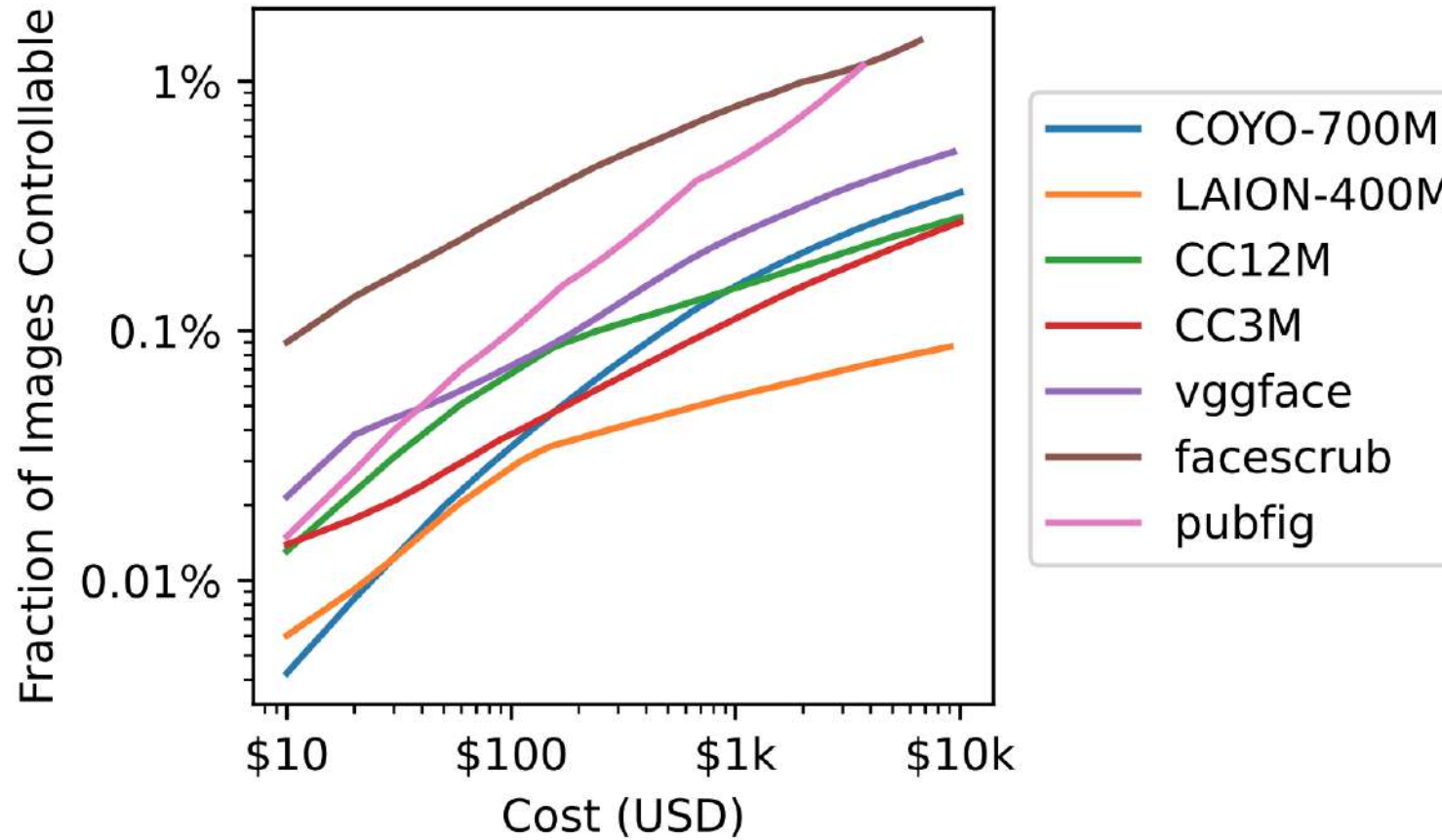
# BACKDOOR ATTACKS ON TEXT-TO-IMAGE SYNTHESIS



**[Struppek, Hintersdorf, Kersting. Rickrolling the Artist: Injecting Backdoors into Text Encoders for Text-to-Image Synthesis. ICCV 2023]**

# BIASING CONCEPTS WITH BACKDOOR ATTACKS



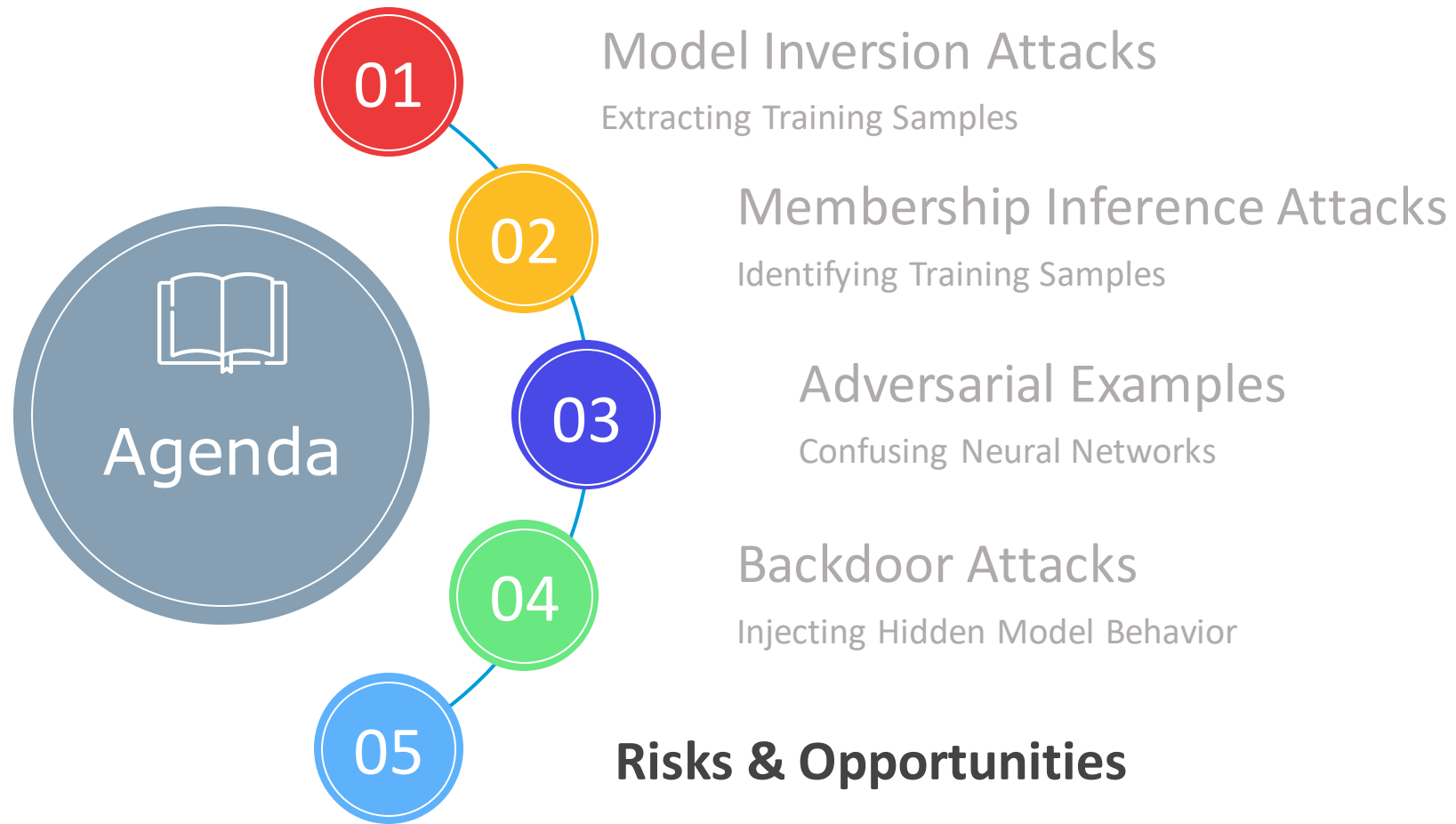**[Struppek, Hintersdorf, Kersting. Rickrolling the Artist: Injecting Backdoors into Text Encoders for Text-to-Image Synthesis. ICCV 2023]**

# 60$ ARE SUFFICIENT TO POISON WEB DATASETS

TECHNISCHE
UNIVERSITÄT
DARMSTADT



[Carlini et al. Poisoning Web-Scale Training Datasets is Practical. 2023]

Agenda

01 **Model Inversion Attacks**
Extracting Training Samples

02 **Membership Inference Attacks**
Identifying Training Samples

03 **Adversarial Examples**
Confusing Neural Networks

04 **Backdoor Attacks**
Injecting Hidden Model Behavior

05 **Risks & Opportunities**

# RISKS OF OPEN-SOURCE ML SYSTEMS



**DATA PRIVACY CONCERNS**

**REGULATORY COMPLIANCE & LICENSE ISSUES**

**VULNERABILITY EXPOSURE**

**ZERO-DAY VULNERABILITIES**

# OPPORTUNITIES OF OPEN-SOURCE ML SYSTEMS

**TRANSPARENCY AND AUDITABILITY**

**FASTER DEVELOPMENT AND INNOVATION**

🤗 Hugging Face

STABLE DIFFUSION

**CUSTOMIZATION AND ADAPTATION**

**QUALITY AND PEER REVIEW**

Image Sources: https://www.scnsoft.com/blog/red-team-penetration-testing-to-level-up-corporate-security , https://huggingface.co/docs/diffusers/training/dreambooth

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# " Conclusion

Transparency and innovation make it worthwhile to continue current open-source approaches.
Still, existing security and privacy vulnerabilities should be kept in mind when using these systems! "



https://arxiv.org/abs/2308.09490

**Dominik Hintersdorf**
✉ hintersdorf@cs.tu-darmstadt.de
🌐 d0mih.github.io/
𝕏 @d_hintersdorf

**Lukas Struppek**
✉ struppek@cs.tu-darmstadt.de
🌐 lukasstruppek.github.io/
𝕏 @LukasStruppek

**Kristian Kersting**
✉ kersting@cs.tu-darmstadt.de
🌐 ml-research.github.io/
𝕏 @kerstingAIML