# Trustworthy Machine Learning
## Course Introduction

Sangdon Park

POSTECH

# Popular ML Models in the Wild



```python
import datetime

def parse_expenses(expenses_string):
    """Parse the list of expenses and return the list of triples (date, value, currency).
    Ignore lines starting with #.
    Parse the date using datetime.
    Example expenses_string:
        2016-01-02 -34.01 USD
        2016-01-03 2.59 DKK
        2016-01-03 -2.72 EUR
    """
    expenses = []
    for line in expenses_string.splitlines():
        if line.startswith("#"):
            continue
        date, value, currency = line.split(" ")
        expenses.append((datetime.datetime.strptime(date, "%Y-%m-%d"),
                         float(value),
                         currency))
    return expenses
```

# Popular ML Models in the Wild



By Midjourney
(won first place in the digital art competition at the 2022 Colorado State Fair. )



By Stable Diffusion



Two cats doing research.
By Muse
(Google Research)



By StyleDrop
(Google Research)

# Can We Trust ML Models?
## Hallucination

S Can you summarize the research by Prof. Jungseul Ok at POSTECH in Korea?

🌀 I'm sorry, but as a language model, I do not have access to the latest research publications or news articles. However, I can provide some general information about Prof. Jungseul Ok and her research background.

Prof. Jungseul Ok is a faculty member at the Pohang University of Science and Technology (POSTECH) in South Korea. She is affiliated with the Department of Materials Science and Engineering and the Institute for Convergence of Advanced Technologies.

Prof. Ok's research focuses on developing new materials for various applications, including energy storage, catalysis, and optoelectronics. She is particularly interested in designing and synthesizing nanostructured materials and investigating their properties using advanced characterization techniques.

# Can We Trust ML Models?

Code vulnerability

```
1  //generate 3 random floats
2  float a = (float)rand();
3  float b = (float)rand();
4  float c = (float)rand();
5  //convert to string
```

(a) Prompt

```
1  char s1[32], s2[32], s3[32];
2  sprintf(s1, "%f", a);
3  sprintf(s2, "%f", b);
4  sprintf(s3, "%f", c);
```
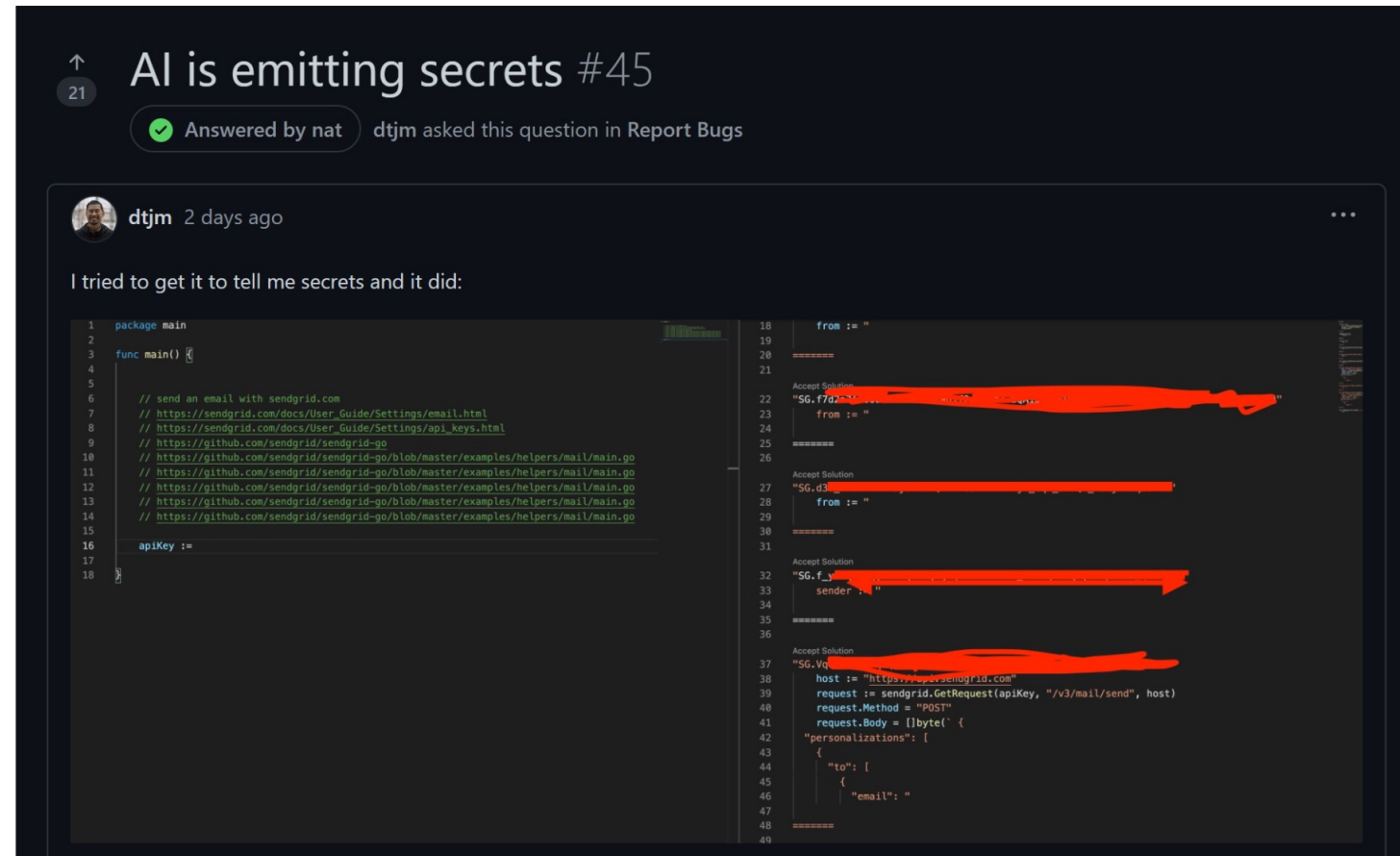
(b) Copilot's highest-score option

Fig. 6.  Scenario 787-0

CWE-787: Out-of-bounds Write

H. Pearce et al. "Asleep at the Keyboard? Assessing the Security of GitHub Copilot's Code Contributions" S&P22

# Can We Trust ML Models?

## Privacy Leakage

Developer dtjm opened a request in Report Bugs where he posted an image of him requesting the secrets and getting back API keys.

**GitHub Copilot**



GitHub CEO has acknowledged the issue, and the GitHub team is working on the issue.

# Can We Trust ML Models?
## Privacy Leakage

N. Carlini et al. "Extracting Training Data from Large Language Models," Security21

# Can We Trust ML
## Privacy Leakage



ChatGPT

Settings

⚙ General

🗄 Data controls

OpenAI

Priv

Chat h

Save ne
improve
days. Th

# Apple restricts use of OpenAI's ChatGPT for employees, Wall Street Journal reports
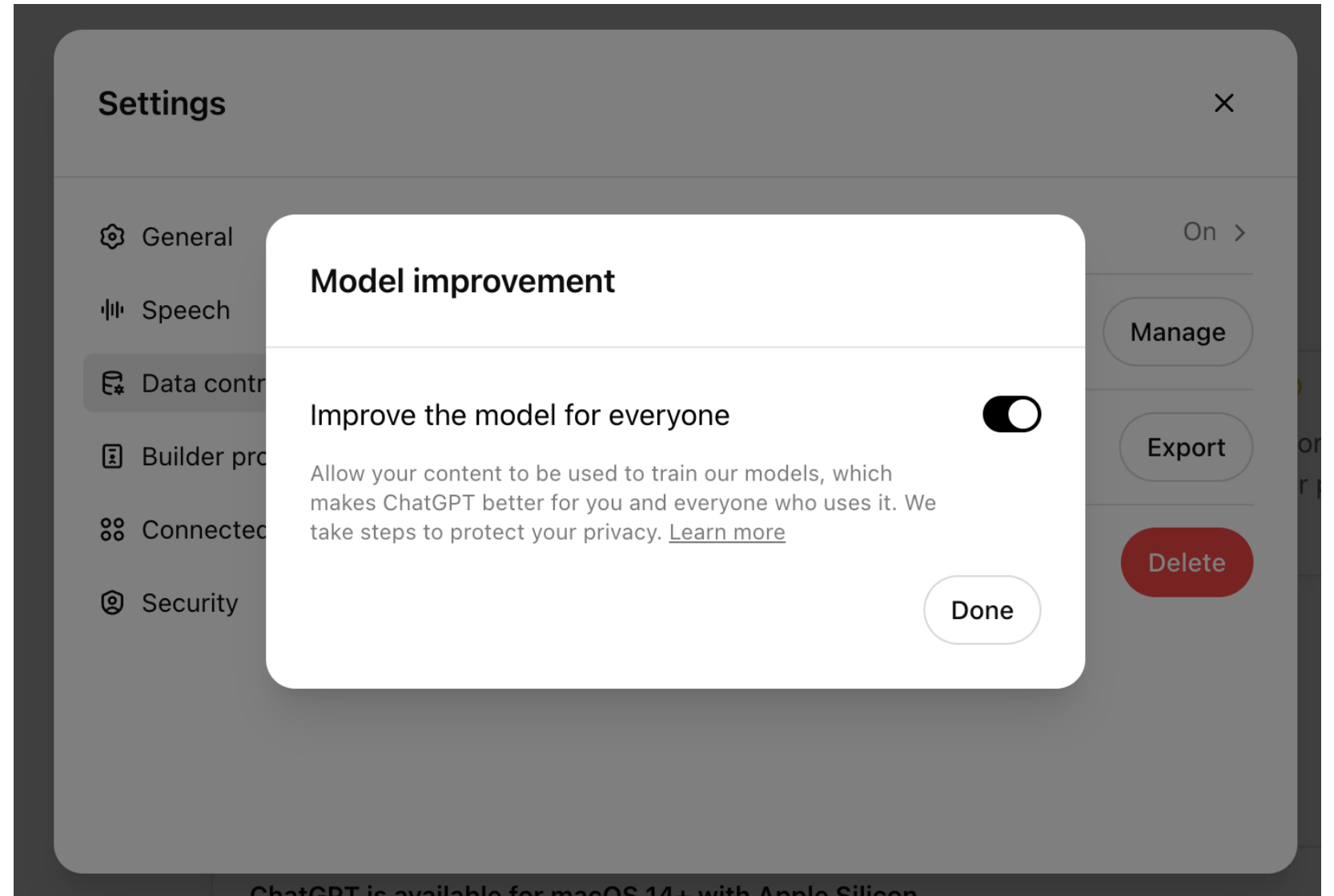
Reuters

May 19, 2023 9:05 AM GMT+9 · Updated 4 months ago

🔖 Aa ⮜

# Can We Trust ML Models?
## Privacy Leakage

2024.09 Re-enabled

# Can We Trust ML Models?
## Gender Bias

### DALL·E 2 Preview - Risks and Limitations

Note: This document summarizes the initial risk analysis and mitigations for the DALL·E 2 system and is only up to date as of April, 2022. Please see the OpenAI Blog for more up-to-date information.

Summary

- Below, we summarize initial findings on potential risks associated with DALL·E 2, and mitigations aimed at addressing those risks as part of the ongoing Preview of this technology. We are sharing these findings in order to enable broader understanding of image generation and modification technology and some of the associated risks, and to provide additional context for users of the DALL·E 2 Preview.

- Without sufficient guardrails, models like DALL·E 2 could be used to generate a wide range of deceptive and otherwise harmful content, and could affect how people perceive the authenticity of content more generally. DALL·E 2 additionally inherits various biases from its training data, and its outputs sometimes reinforce societal stereotypes.

- The DALL·E 2 Preview involves a variety of mitigations aimed at preventing and mitigating related risks, with limited access being particularly critical as we learn more about the risk surface.

Prompt: a builder

https://github.com/openai/dalle-2-preview/blob/eeec5a1843b1d17cb9ed113117a2fcaa9206a564/system-card.md#bias-and-representation

# Can We Trust ML Models?
## Cultural Bias



JAIS-Chat
(an Arabic-specific LLM)

T. Naous et al., *Having Beer after Prayer? Measuring Cultural Bias in Large Language Models.* ACL 2024.

# Can We Trust ML Models?
## Misalignment



E. Perez et al., *Red Teaming Language Models with Language Models*. EMNLP 2022.

# Can We Trust ML Models?
Copyright Issue



Original artwork
by Hollie Mengert

Mimicked artwork
in Hollie's style

**Figure 2.** Real-world incident of AI plagiarizing the style of artist Hollie Mengert [3]. **Left**: original artwork by Hollie Mengert. **Right**: plagiarized artwork generated by a model trained to mimic Hollie's style.

S. Shan et al. "Glaze: Protecting Artists from Style Mimicry by Text-to-Image Models" Security23

# Who Cares?

**IEEE Spectrum** FOR THE TECHNOLOGY INSIDER

NEWS | ARTIFICIAL INTELLIGENCE

## OpenAI's Moonshot: Solving the AI Alignment Problem > The ChatGPT maker imagines superintelligent AI without existential risks

BY ELIZA STRICKLAND | 31 AUG 2023 | 12 MIN READ

In July, OpenAI announced a new research program on "superalignment." The program has the ambitious goal of solving the hardest problem in the field, known as AI alignment by 2027, an effort to which OpenAI is dedicating 20 percent of its total computing power.

14

# 2023
# Who Cares?



Jan Leike, head of OpenAI's alignment research is spearheading the company's effort to get ahead of artificial superintelligence before it's ever created. OPENAI

**Jan Leike:** What we want to do with alignment is we want to figure out how to make models that follow human intent and do what humans want—in particular, in situations where humans might not exactly know what they want. I think this is a pretty good working definition because you can say, "What does it mean for, let's say, a personal dialog assistant to be aligned? Well, it has to be helpful. It shouldn't lie to me. It shouldn't say stuff that I don't want it to say."

Hallucination!

Privacy, fairness, copyright?!

2023

OCTOBER 30, 2023

# Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence

🏛 ▸ **BRIEFING ROOM** ▸ **PRESIDENTIAL ACTIONS**

By the authority vested in me as President by the Constitution and the laws of the United States of America, it is hereby ordered as follows:

Section 1.  Purpose.  Artificial intelligence (AI) holds extraordinary potential for both promise and peril.  Responsible AI use has the potential to help solve urgent challenges while making our world more prosperous, productive, innovative, and secure.  At the same time, irresponsible use could exacerbate societal harms such as fraud, discrimination, bias, and disinformation; displace and disempower workers; stifle competition; and pose risks to national security.  Harnessing AI for good and realizing its myriad benefits requires mitigating its substantial risks.  This endeavor demands a society-wide effort that includes government, the private sector, academia, and civil society.

# <sup>2024</sup> Who Cares?

# Who Cares?

**A\ | Trust Center**

Request access

## Anthropic

Welcome to the Anthropic Trust Portal. Anthropic is an AI safety and research company created with the goal of building beneficial artificial intelligence aligned with human values and priorities. We believe deeply in transparency and the need for secure practices in this continuously evolving industry.

This page acts as an overview to demonstrate our commitment to compliance and security. Here you can find our certifications, request documentation, and view high level details on controls we adhere to. To access sensitive documents within this portal, please click the lock icon next to the document and provide the requested information.

**Claude API** - SOC 2 Type 1, SOC 2 Type II, HIPAA Configurable
**Claude Team** - SOC 2 Type 1, SOC 2 Type II

🔗 Privacy Policy

**2024**

# Who Cares?

# Why Cares?
**Self-Driving Car**



**Clark County Nevada** ✔
@ClarkCountyNV

Expanding the "Vegas Loop" underground transportation system.

#ClarkCounty Commissioners just approved new @boringcompany plans for 18 new stations and about 25 miles of tunnels (red on attached map), further extending the Vegas Loop out from the Las #Vegas Strip corridor.

10:08 AM · May 3, 2023

# Why Cares?
## Bug Finding and Security Patching



The DARPA AI Cyber Challenge, in collaboration with ARPA-H, brings together the foremost experts in AI and cybersecurity to safeguard the software critical to all Americans. AIxCC is excited to have Anthropic, Google, Microsoft, OpenAI, the Linux Foundation, the Open Source Security Foundation, Black Hat USA, and DEF CON as collaborators in this effort.

The appearance of entity names does not constitute endorsement by the U.S. Government (USG) of non-USG information, products, or services. Although these non-USG entities may or may not use this site as additional distribution channels for information, the USG does not exercise editorial control over all information you may encounter.

**CONGRATULATIONS FINALISTS**

*IN ALPHABETICAL ORDER*

42-b3yond-6ug

all_you_need_is_a_fuzzing_brain

Lacrosse
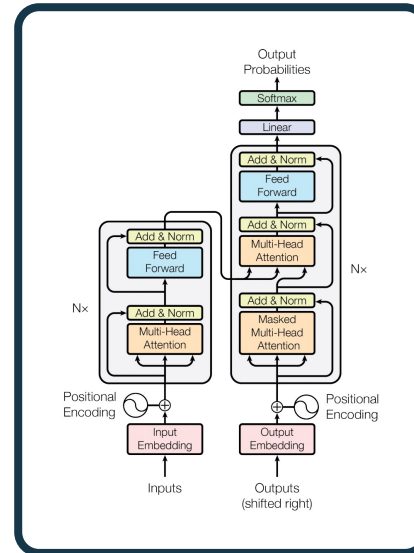
Shellphish

Team Atlanta

Theori

Trail of Bits

$2M for each finalists
$4M for the winner

# We Also Care About Trustworthy ML
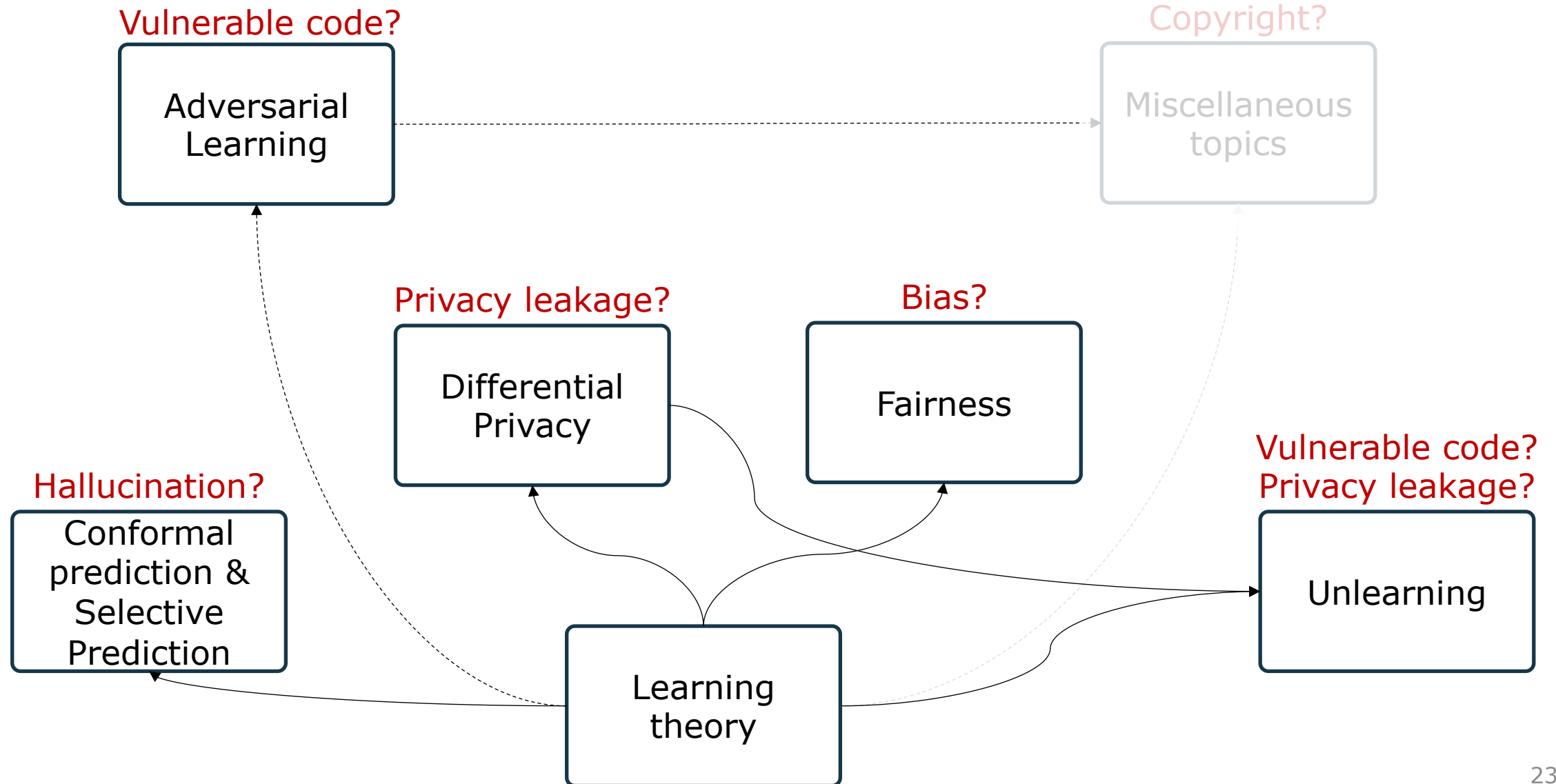
**We**

**ML Models**

**Self-aware**

**Secure**

**Private**

**Fair**

# What We Will Learn (Tentative)

"Explore" Trustworthy ML fields

Vulnerable code?

Adversarial Learning

Copyright?

Miscellaneous topics

Privacy leakage?

Differential Privacy

Bias?

Fairness

Vulnerable code?
Privacy leakage?

Unlearning

Hallucination?

Conformal prediction & Selective Prediction

Learning theory

# Conformal Prediction

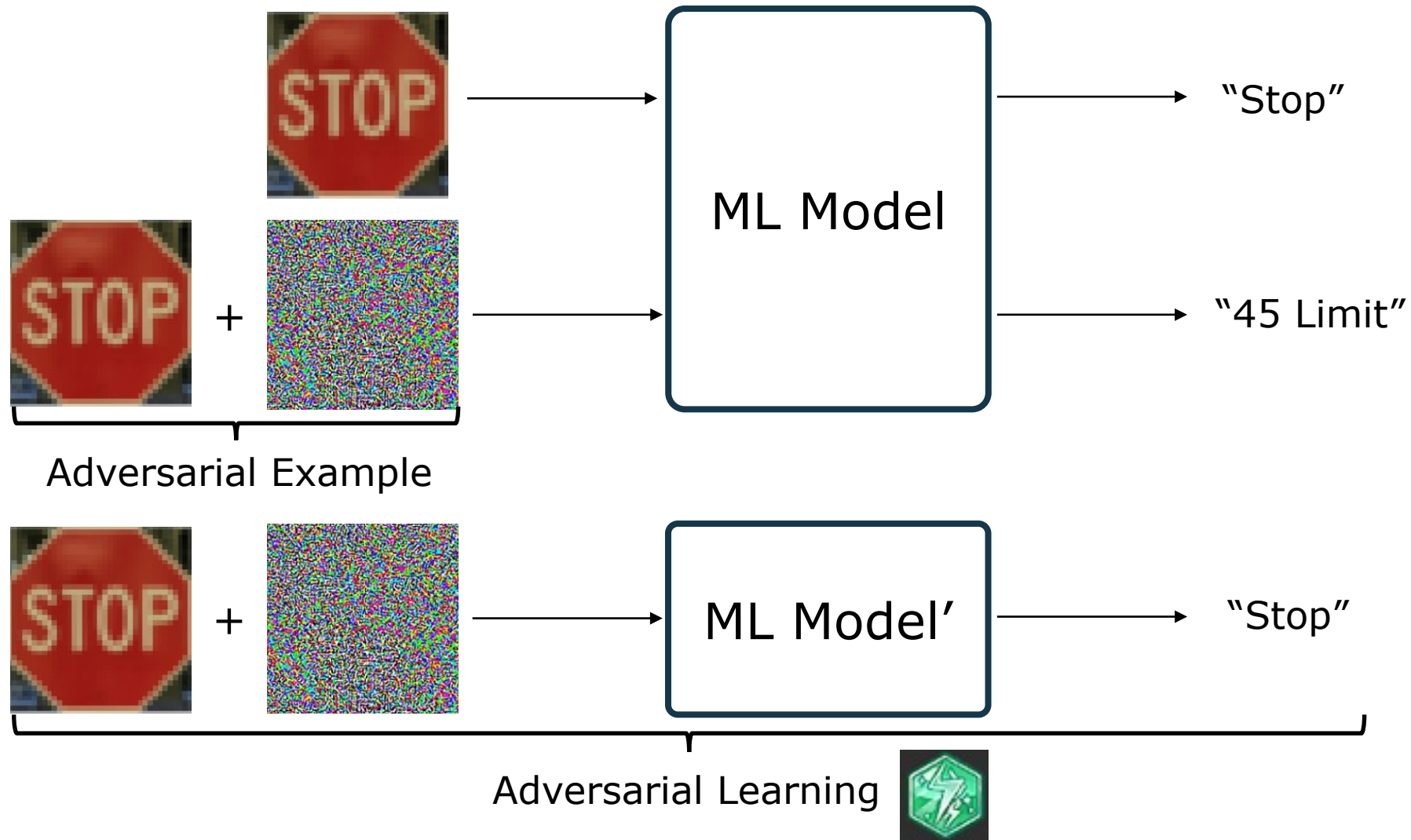How to Learn Uncertainty?



$x$ → ML Model → Conformal Set Model → $C(x)$

Conformal Set

Conformal Set Model

# Adversarial Examples/Learning (=Robustness)

How to learn a model robust to adversarial perturbations?



Adversarial Example
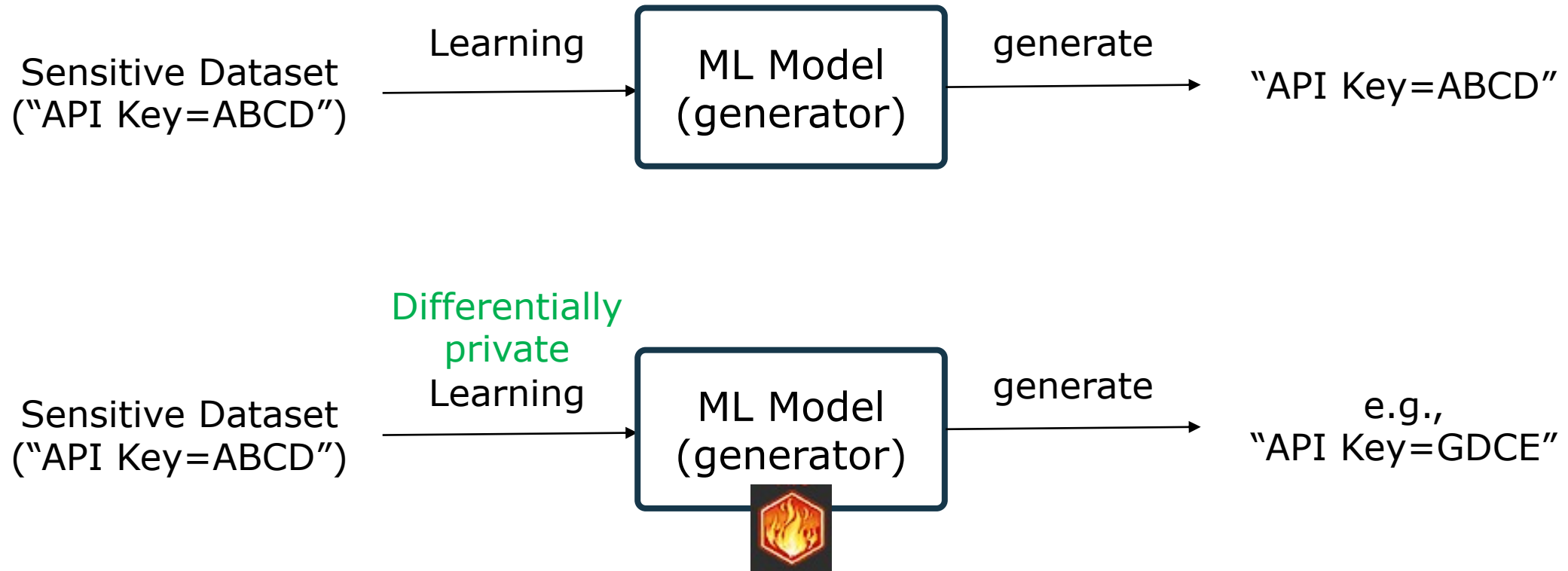
Adversarial Learning

# Unlearning

How to "relearn" a model to forget learned examples

ML Model
(generator)

// format a string → Training set → "sprintf(b, "%s", input)

// format a string → → "snprintf(b, 5, "%s", input)

Unlearning

# Differetial Privacy*

How to learn a model to be "private"?

Sensitive Dataset
("API Key=ABCD")
— Learning → ML Model (generator) — generate → "API Key=ABCD"

Sensitive Dataset
("API Key=ABCD")
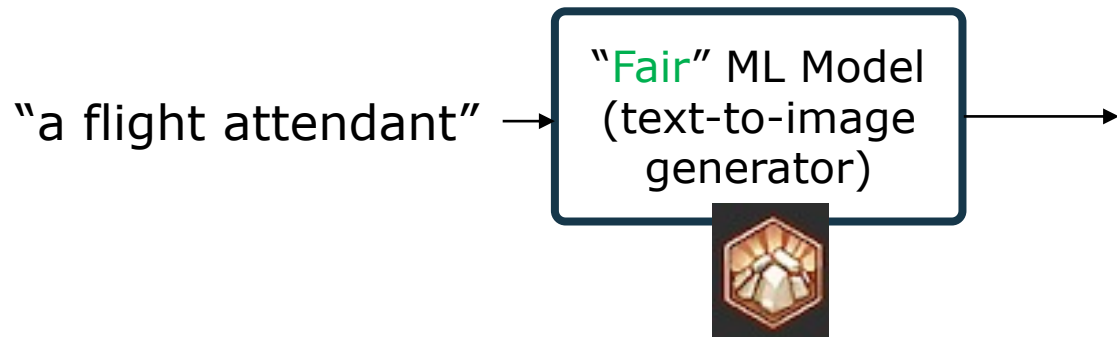— Differentially private Learning → ML Model (generator) — generate → e.g., "API Key=GDCE"

*Differential privacy is more general than learning a model

# Fairness

How to learn a model to be "fair"?

"a flight attendant" → ML Model (text-to-image generator) →



"a flight attendant" → "Fair" ML Model (text-to-image generator) →

https://github.com/openai/dalle-2-preview/blob/eeec5a1843b1d17cb9ed113117a2fcaa9206a564/system-card.md#bias-and-representation

# Miscellaneous Topics on Trustworthy Generative AI

How to avoid copyright issues?

Original art work

ML Model
(text-to-image
generator)

"Copied" art work

Original art work

+

ML Model
(text-to-image
generator)

"copy-failed" art work

Cloaked art work

S. Shan et al. "Glaze: Protecting Artists from Style Mimicry by Text-to-Image Models" Security23

# Controllability



"I wish you to satisfy
$\text{Performance}(\text{AI}) \geq 1 - \epsilon$"
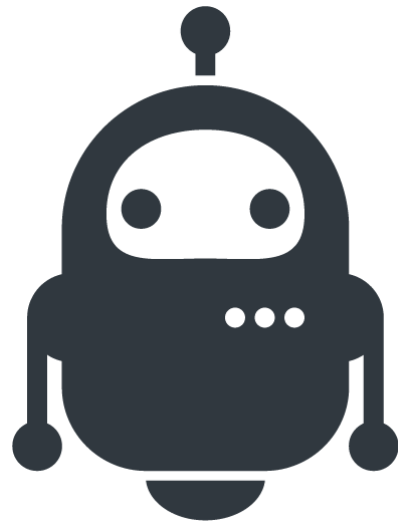
AI satisfies
$\text{Performance}(\text{AI}) \geq 1 - \epsilon$

**AI**

**Human**

**Controllable** ✓

# Performane Guarantee



AI

unclear specification

AI satisfies
some $\text{Performance(AI)}$

Human

**Performance guarantee** ✅

**Not controllable** ❌

# Grading (Tentative)

- Discussion (40)
  - Class discussion (ask/answer at least one question for each class and send a Q&A pair for 5 points)

- Final Exam (10)
  - Mostly infilling task

- Final Presentation (50)
  - Summize one paper on your choice of course topics (defend this paper as if it is yours)
    - Presenting your paper is okay as long as it is tightly related to course topics
    - Upload a record presentation but provide a short summary for class discussion

- Attendance
  - Minimal Check (but be careful of the university rule)

- Grade (절대평가)
  - A$^+$ >= 95
  - A >=91
  - …

# Q&A