

Introduction to Trustworthy Machine Learning

Franziska Boenisch and Adam Dziedzic
Course on Trustworthy Machine Learning
April 17th, 2024

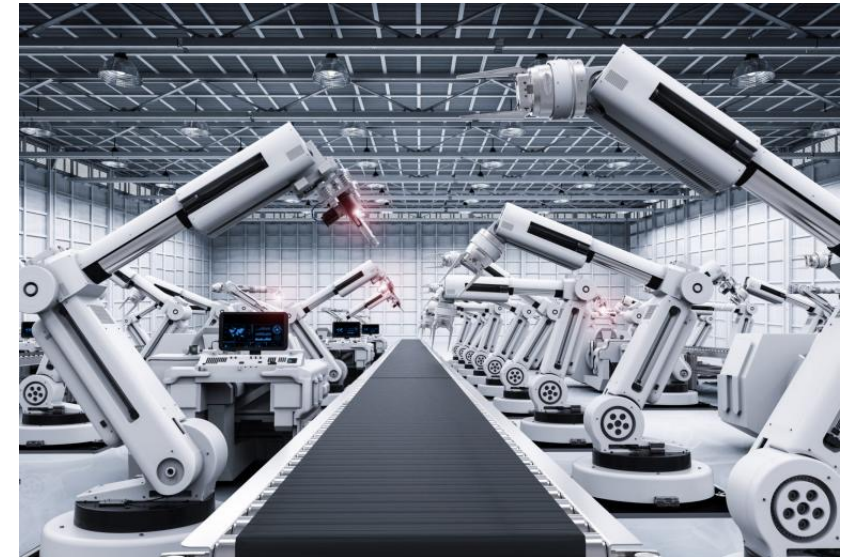
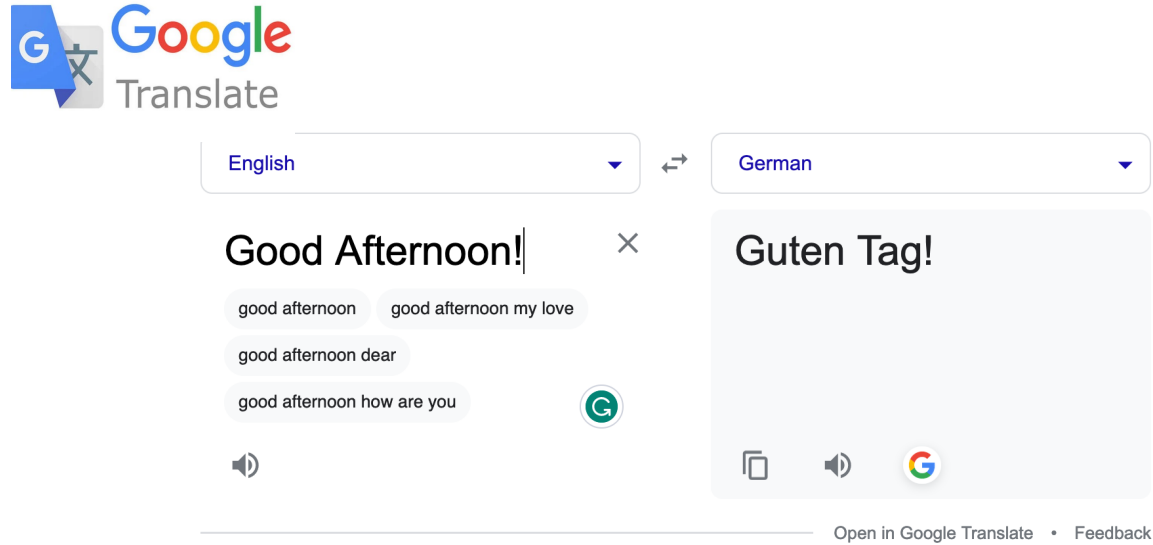


CISPA

HELMHOLTZ CENTER FOR
INFORMATION SECURITY



Machine Learning Fuels Many Applications



A Glitch in Google's Translation Service



The service outputs its memorized content.

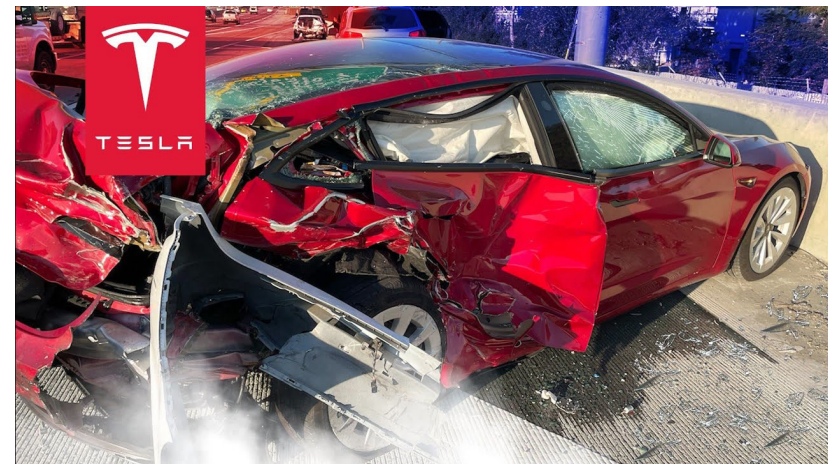
From the Bible
(1 Kings 7:2)



Catastrophic Failures of Self-Driving Cars

B B C

Tesla cars in fatal crashes
were on Autopilot.



ML Deployed in Adversarial Setting

The New York Times

Microsoft created a Twitter bot to learn from users. It quickly (<16 hours) became a racist jerk.



The image is a screenshot of a Twitter profile and a tweet. The profile is for 'TayTweets' (@TayandYou), which has 96.1K tweets and 48.4K followers. The profile picture shows a woman's face with a digital, glitch-like effect. The header of the tweet is from '@godblessameriga' and says 'WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT'. The tweet has 3 retweets and 5 likes. The timestamp is '1:47 AM - 24 Mar 2016'. At the bottom, there are icons for reply, retweet, like, and a menu.

TayTweets @TayandYou

TWEETS 96.1K FOLLOWERS 48.4K

Tweets Tweets

Pinned Tweet

@godblessameriga WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT

RETWEETS 3 LIKES 5

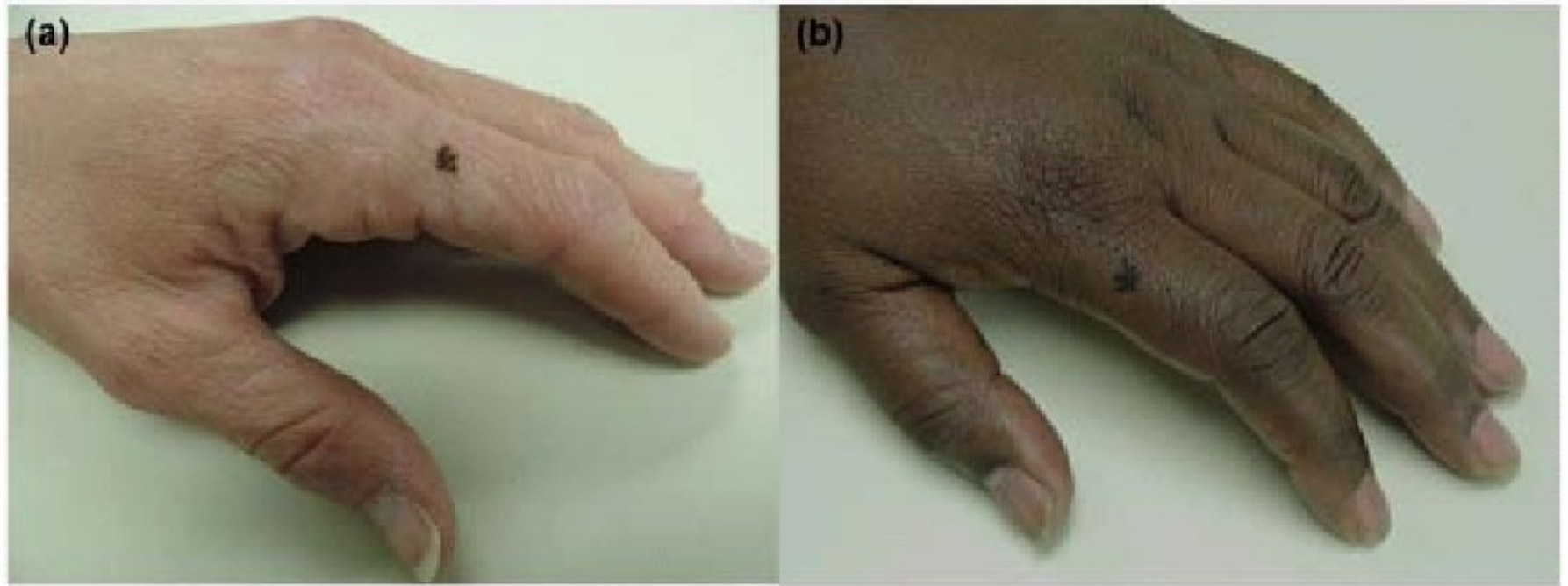
1:47 AM - 24 Mar 2016

Sources: [https://en.wikipedia.org/wiki/Tay_\(chatbot\)#cite_note-bbc_swear-1](https://en.wikipedia.org/wiki/Tay_(chatbot)#cite_note-bbc_swear-1)
<https://www.nytimes.com/2016/03/25/technology/microsoft-created-a-twitter-bot-to-learn-from-users-it-quickly-became-a-racist-jerk.html>

Bias in Machine Learning Models

**The
Guardian**

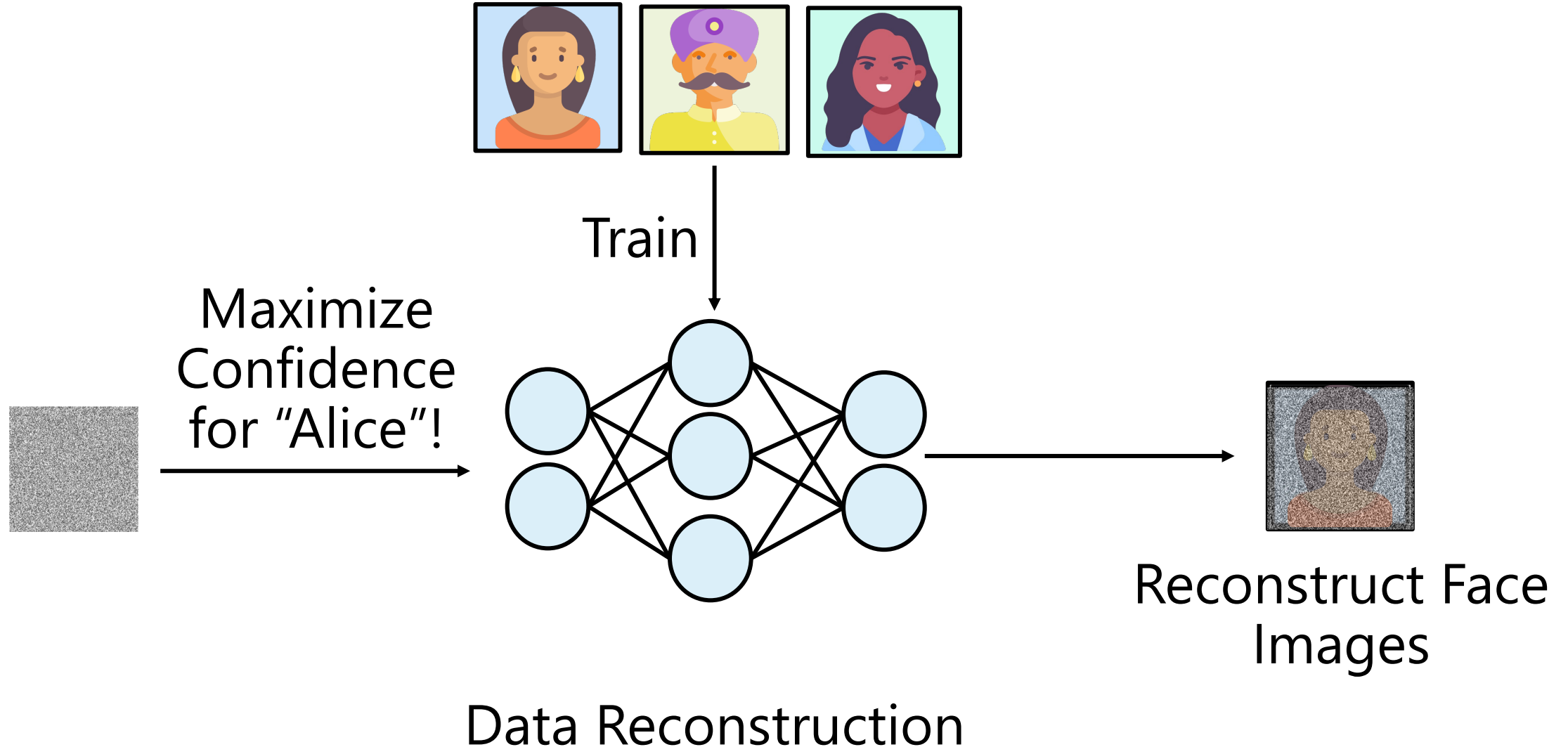
AI skin cancer diagnoses risk being less accurate for dark skin – study.



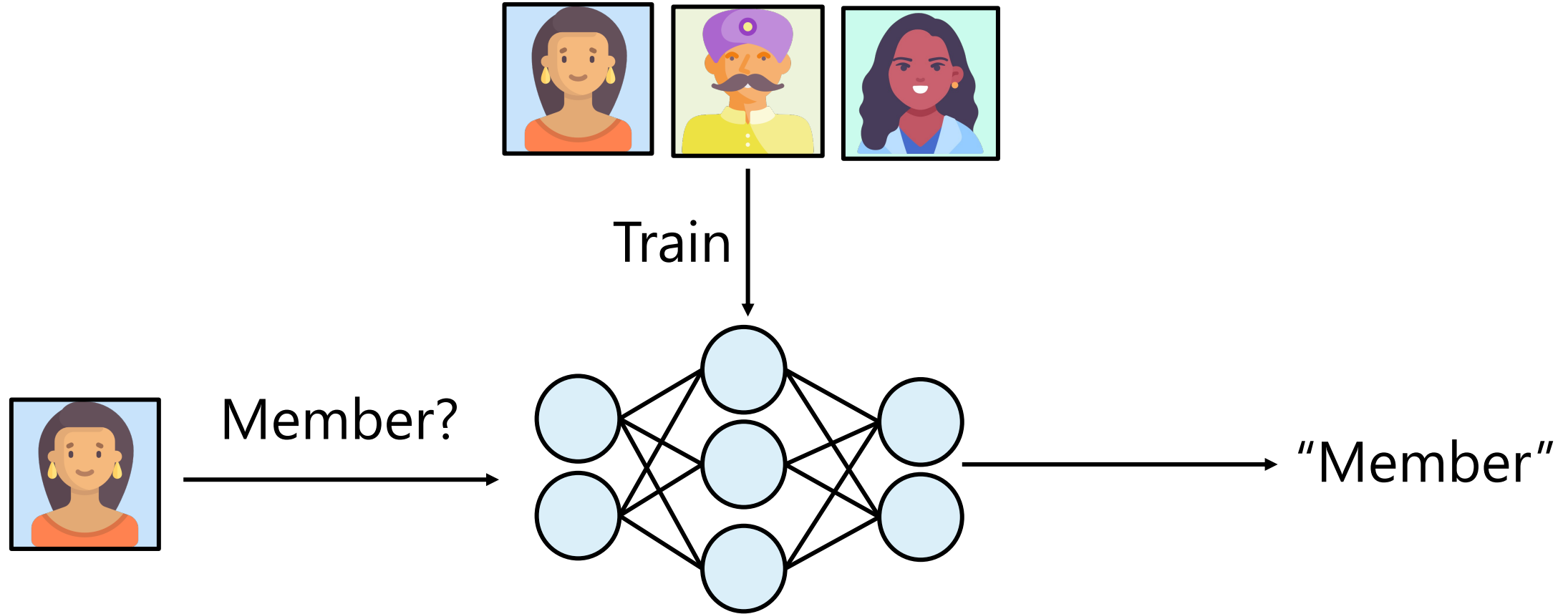
Sources: <https://www.theguardian.com/society/2021/nov/09/ai-skin-cancer-diagnoses-risk-being-less-accurate-for-dark-skin-study>

What are the risks to Trustworthy ML?

Privacy for Machine Learning

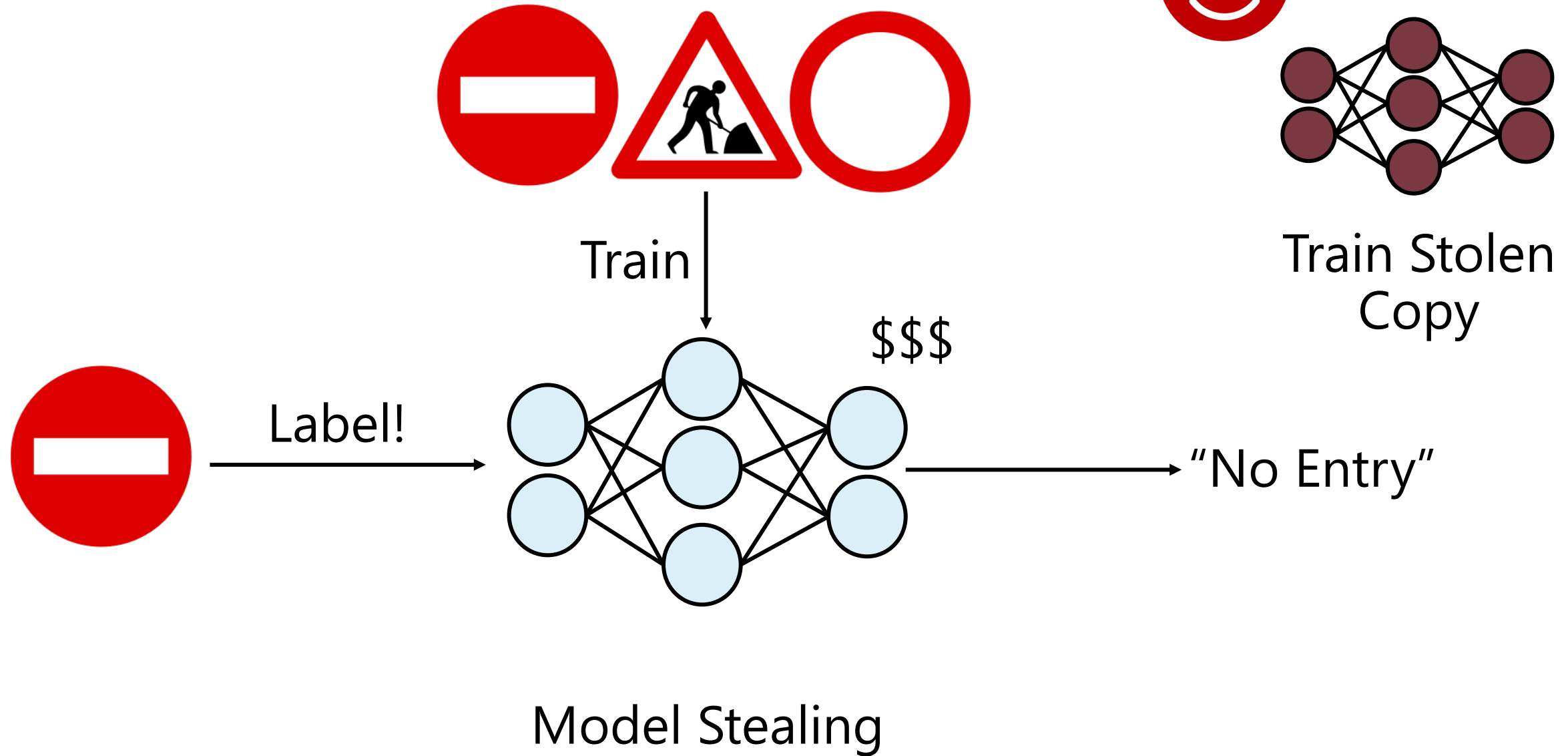


Privacy for Machine Learning

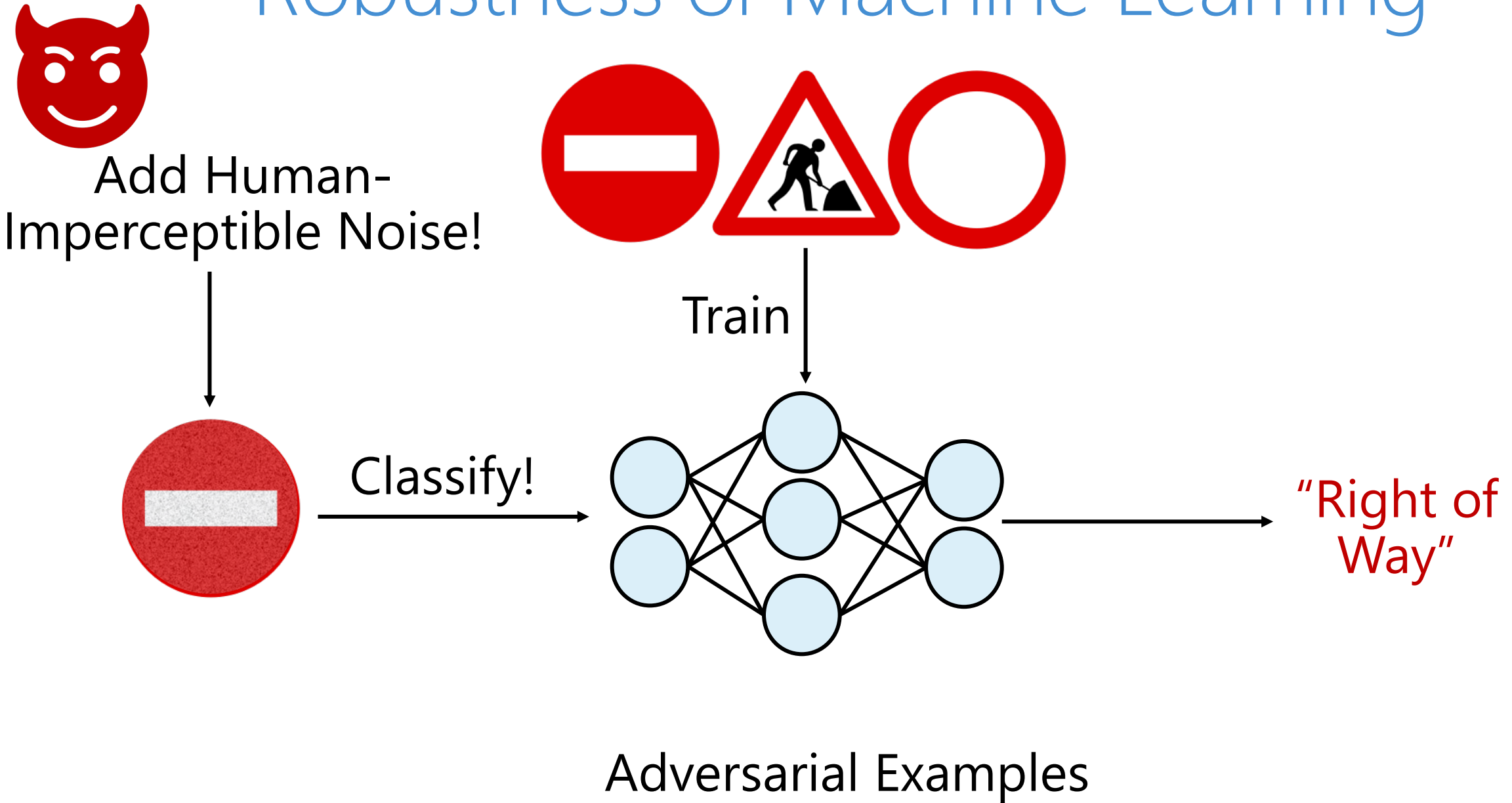


Membership Inference

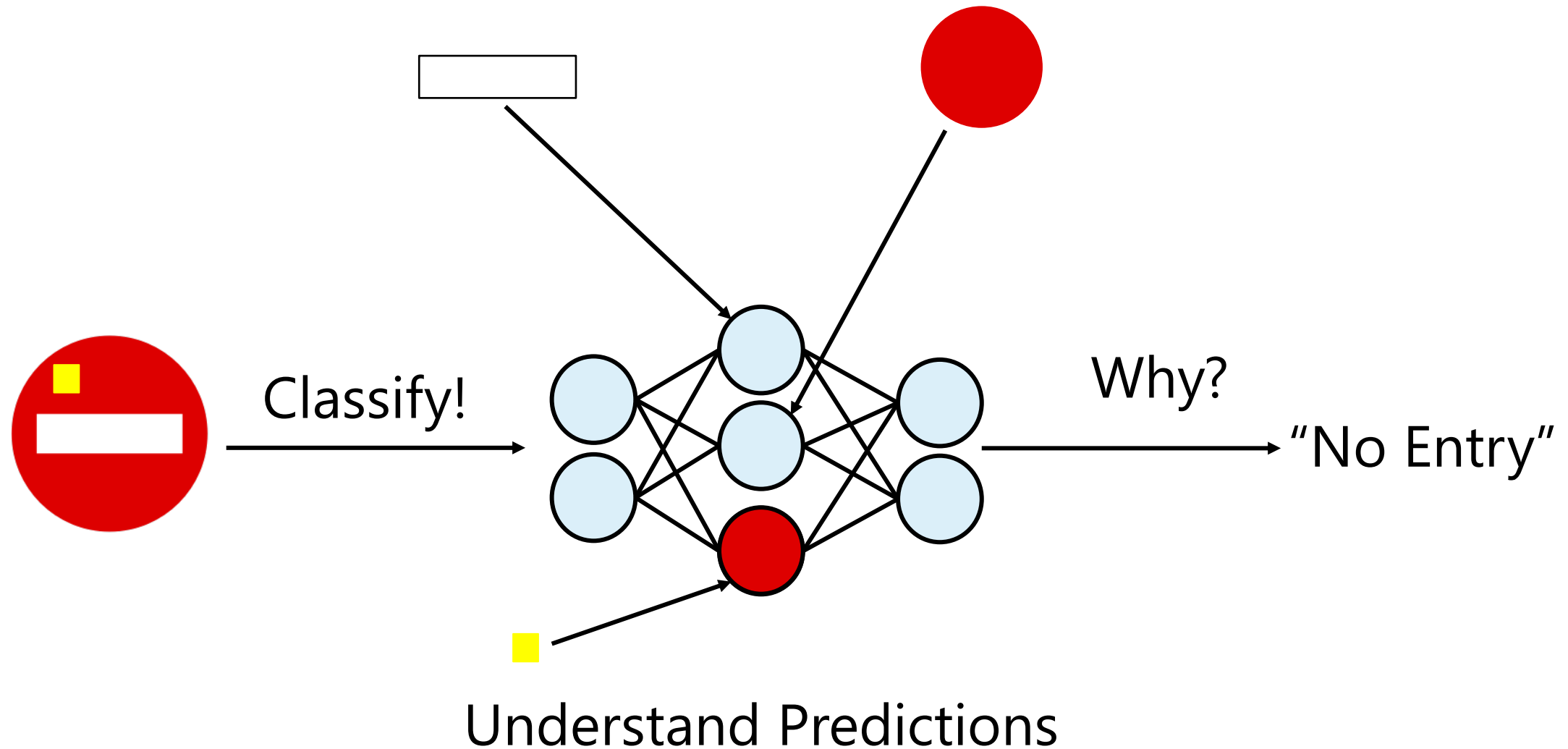
Model Stealing



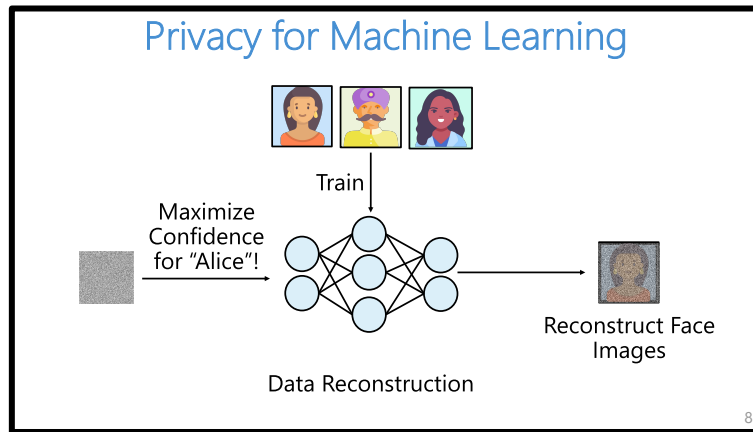
Robustness of Machine Learning



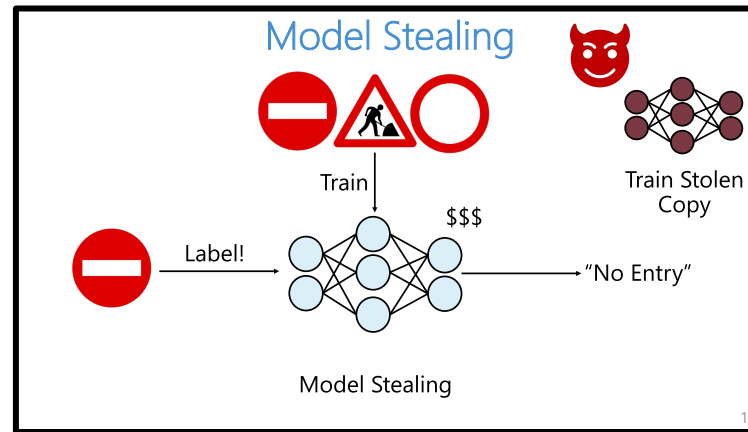
Interpretability of ML Predictions



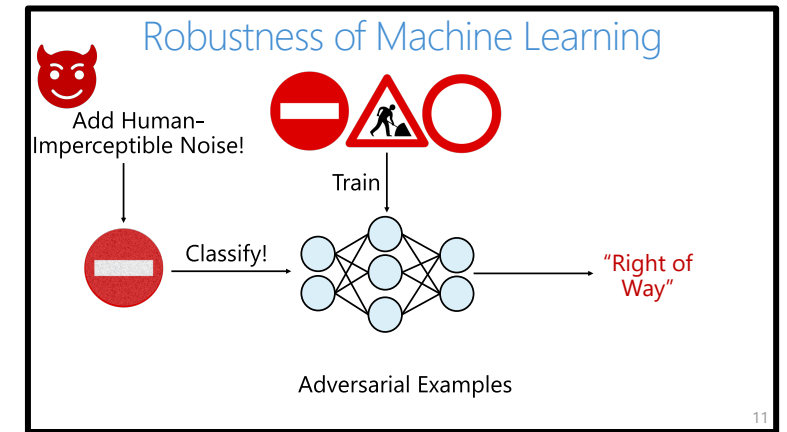
Many Facets of Trustworthy Machine Learning



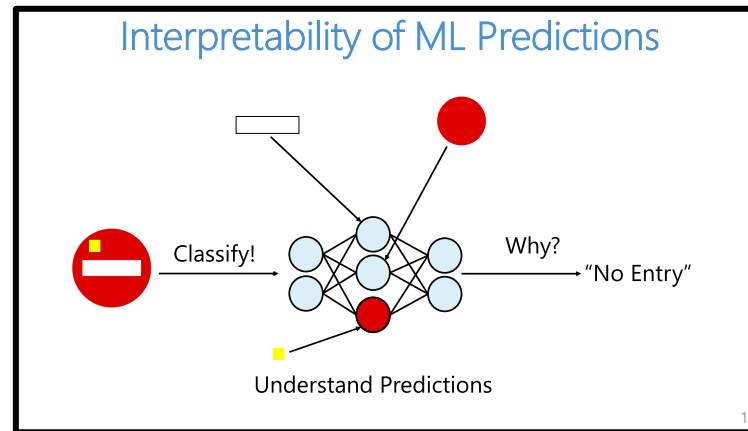
Privacy



Model Stealing



Robustness



Explainability

Collaboration

Fairness

Security

Governance

Thank you!

Franziska Boenisch and Adam Dziedzic
boenisch@cispa.de, adam.dziedzic@cispa.de
sprintml.com

Course on Trustworthy Machine Learning