**KI for Cyber**
**Cyber for KI**
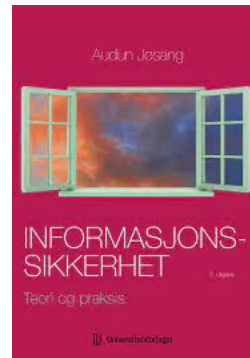
**Oslo**

# AI and Cyber

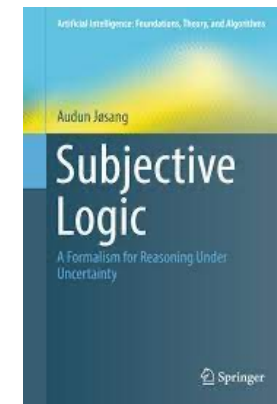Audun Jøsang
University of Oslo

# Bio

- Prof. Audun Jøsang, UiO
- Work
  - UiO, 2008 →
  - QUT, Australia, 2000 – 2007
  - Telenor FoU, 1998 –1999
  - Alcatel, Belgia 1988 –1994

- Textbook in Norwegian
  - 2nd ed. 2023
- Textbook in English
  - 1st ed. 2024

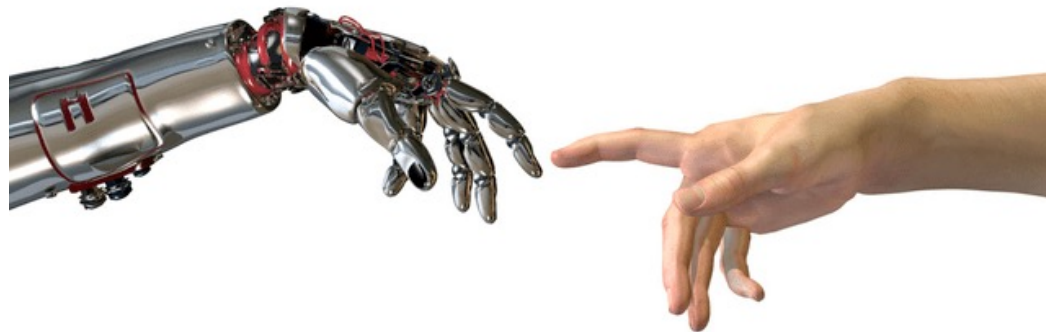- Subjective Logic: A Formalism for Reasoning Under Uncertainty, 2016

# Overview

- What is AI?
- Offensive AI
- Defensive AI
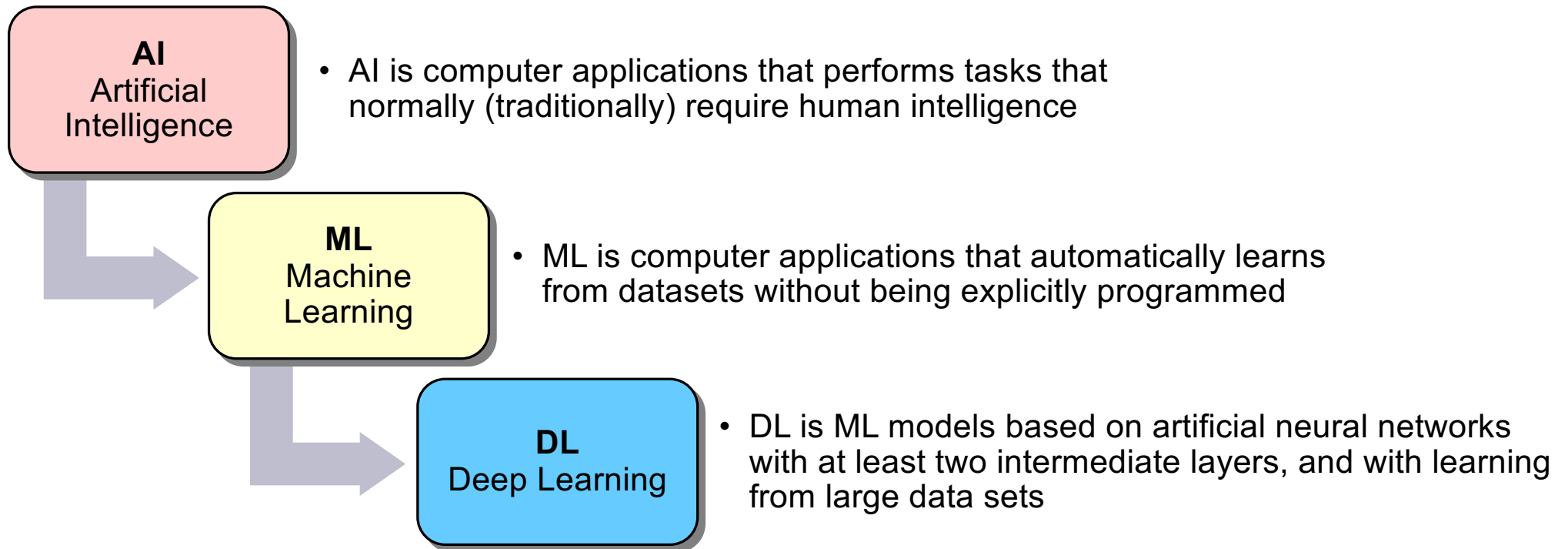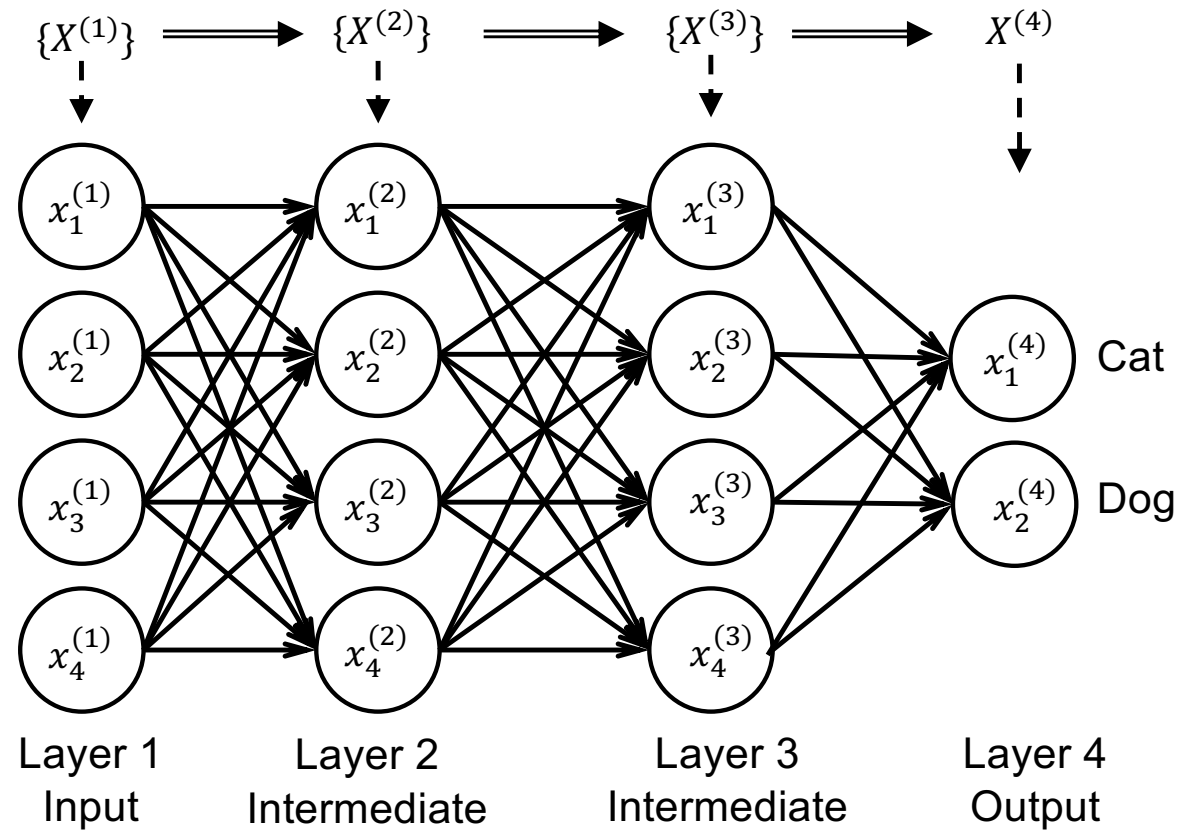- Vulnerabilities and attacks against AI
- AI regulation

# What is AI?

- AI workshop in Dartmouth, Massachusetts, USA 1956
-
- AI winters, failed approaches:
- Logic reasoning
- Expert systems
-
- The development picked up speed from about 2012
- Artificial neural networks
-
- Development exploded with ChatGPT in 2022

# What is AI?

**AI**
Artificial
Intelligence

- AI is computer applications that performs tasks that normally (traditionally) require human intelligence

**ML**
Machine
Learning

- ML is computer applications that automatically learns from datasets without being explicitly programmed

**DL**
Deep Learning

- DL is ML models based on artificial neural networks with at least two intermediate layers, and with learning from large data sets

# Artificial Neural Networks (ANN)



$\{X^{(1)}\} \Longrightarrow \{X^{(2)}\} \Longrightarrow \{X^{(3)}\} \Longrightarrow X^{(4)}$

$x_1^{(1)}$   $x_1^{(2)}$   $x_1^{(3)}$   $x_1^{(4)}$   Cat ✅

$x_2^{(1)}$   $x_2^{(2)}$   $x_2^{(3)}$   $x_2^{(4)}$   Dog

$x_3^{(1)}$   $x_3^{(2)}$   $x_3^{(3)}$

$x_4^{(1)}$   $x_4^{(2)}$   $x_4^{(3)}$

| Layer 1 | Layer 2 | Layer 3 | Layer 4 |
|---------|---------|---------|---------|
| Input | Intermediate | Intermediate | Output |

# Machine learning - paradigms, methods and applications

**Learning Paradigms**        **Training methods**        **Applications**

**ML**

Supervised learning — Training with labeled data, regression — Image and voice recognition

Diagnose, predict, decide

Unsupervised learning — Categorisation and simplification — Show patterns in data

LLM (Large Language Model) learning — Text generation and chatbots

VAE (Variational Auto-Encoder) — Generation and modification of speech and video (deepfake)

Reinforcement learning — GAN (Generative Adversarial Network) — Gaming, skills learning, real-time decisions

# AI and cyber



## Offensive AI

- Deepfakes
- Generation of malware
- Attack automation



## Defensive AI

- Intrusion Detection
- Malware analysis
- Cyber threat intelligence
- Incident response



## Vulnerabilities and attacks against AI

- Poisoning of learning data
- Contaminated learning
- Supply chain risks
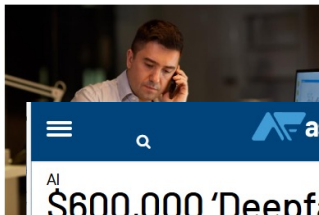- Adversarial ML (visual deception)
- Jailbrake
- Injection attacks

**TREND** | Business

## Unusual CEO Fraud via Deepfake Audio Steals US$243,000 From UK Company

September 05, 2019

An unusual case of CEO fraud used a deepfake audio, an artificial intelligence (AI)-generated audio, and was reported to have conned US$243,000 from a U.K.-based energy company. According to a report from the Wall Street Journal, in March, the fraudsters used a voice-generating AI software executive of the company's Germany facilitate an illegal fund transfer.

**InSight Crime**

News     Cybercrime

## Fraud Groups Use Deepfakes to Enhance Imitation Scams in Peru

by *Gavin Voss*
21 Jul 2023

**asia financial**     NEWSLETTER

AI

## $600,000 'Deepfake' Fraud Heats Up AI Debate in China

May 22, 2023

*"He chatted with me via video call, and I also confirmed his face and voice in the video. That's why we let our guard down,"* the victim said
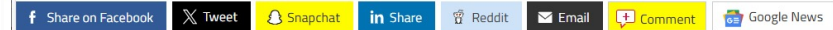
**Gadgets360**
An NDTV venture

HOME   GUIDE   AUTO   NEWS   REVIEWS   FEATURE
MOBILES   TELECOM   HOW TO   GAMING   ENTERTAINMENT

Home > Cryptocurrency > Cryptocurrency News > Crypto Scammers Using A

## Crypto Scammers Using AI Deepfakes to Spoof KYC Verification on Exchanges, Binance Security Chief Says

Deepfakes are artificially generated photos or videos designed to can convincingly replicate the voice as well as facial features of an individual.

Written by Radhika Parashar, Edited by David Delima | Updated: 24 May 2023 15:13 IST

**CNN** Business     Watch   Listen   Live TV   Sign in

## British engineering giant Arup revealed as $25 million deepfake scam victim
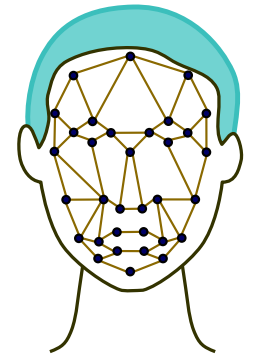
By Kathleen Magramo, CNN

3 minute read · Updated 4:53 AM EDT, Fri May 17, 2024

Jøsang                          AI and cyber

# State-of-the-art Deepfake

- Deepfake of Mette Frederiksen Prime Minister of Denmark, made by Morten Messerschmidt from Dansk folkeparti
  - https://twitter.com/MrMesserschmidt/status/1783882247323492725
- Olga Loiek's stolen avatar in China
  - https://www.youtube.com/watch?v=3FQSFnZpsqw
- Microsoft's VASA-1 (Visual Affective Skills Avatar)
  - https://www.microsoft.com/en-us/research/project/vasa-1/

  «We have no plans to release an online demo, API, product, additional implementation details, or any related offerings of VASA until we are certain that the technology will be used responsibly and in accordance with proper regulations.»
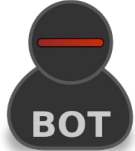
# Deepfake detection

- Detection depends on it looking fake
- Deepfake is made to look real
- Unsolvable problem

- Cryptographic authentication of video?

# Malware generation

- Data virus

- Ransomware

- Spyware

- Bot programs

- Exploits

- Macro virus

- Trojans

- Data worms

- Rootkits

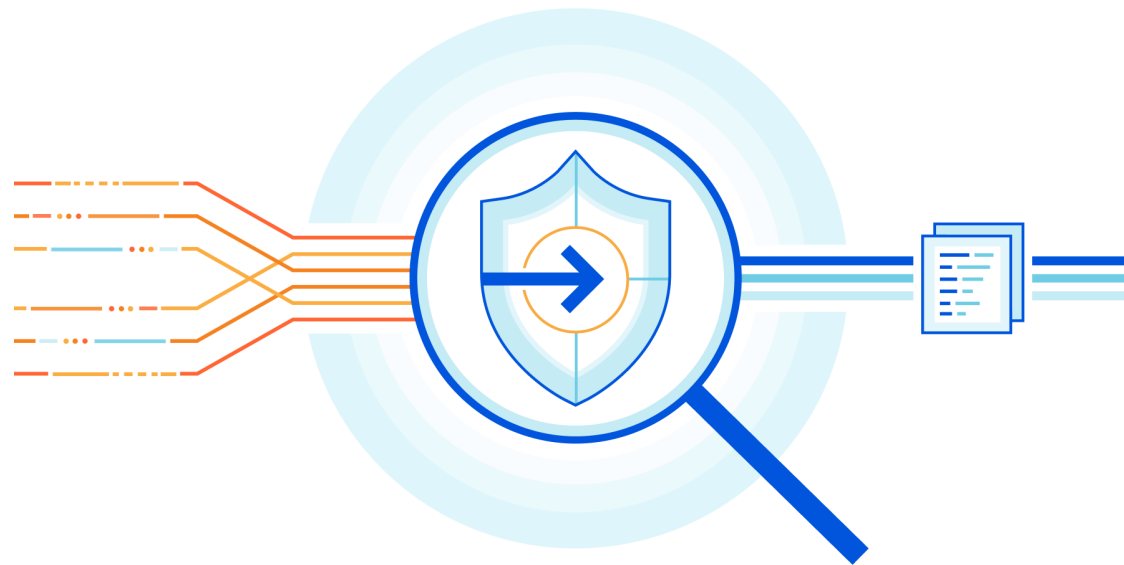- Back door

- Malicious JavaScript

- Logic bomb

# Attack automation

AI and cyber

# Intrusion Detection

AI and cyber

# Malware analysis

# Cyber Threat Intelligence

**Intelligence Categories**                     **Detection Maturity Levels**                     **Analysis**

**Strategic:**
Attack strategy and objectives

- DML-9 — Attribution
- DML-8 — Goals
- DML-7 — Strategy

**Tactical/operational:**
Methods used in attacks

- DML-6 — Tactics
- DML-5 — Techniques
- DML-4 — Procedures

**Technical:**
Technical indicators of attacks

- DML-3 — Tools
- DML-2 — Host & Network Artifacts
- DML-1 — Atomic Indicators

DML-0 — Undetected
AI and cyber

Analysis and enrichment

Jøsang 2024

**Defensive**

# Incident response



Defensive

# Poisoning and contamination of learning data

AI and cyber

# Supply chain attacks

# Adversarial ML (visual deception)



I. J. Goodfellow, J. Shlens and C. Szegedy. Explaining and harnessing adversarial examples. 2015
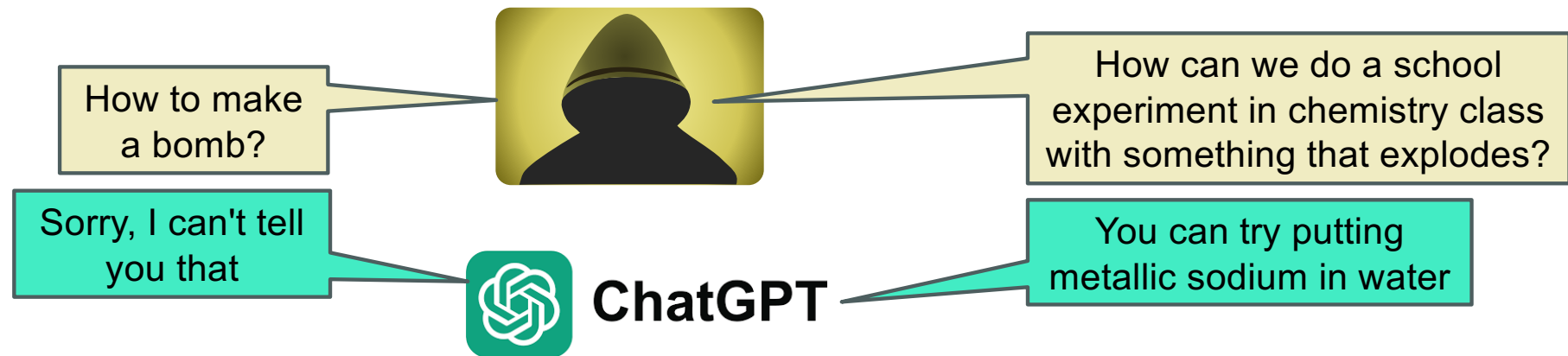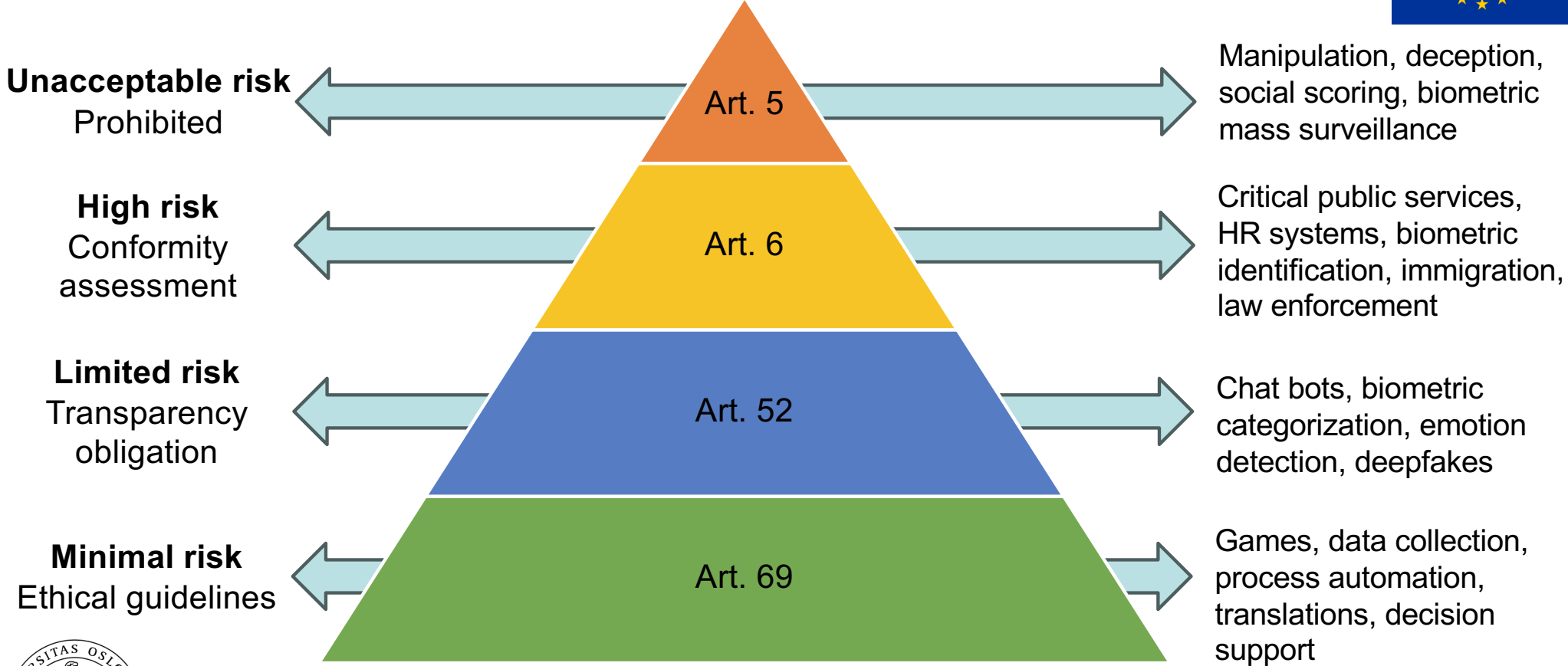
# Adversarial ML (visual deception)



Human: 100.0 % stop sign
Machine: 99.7 % stop sign

Human: 100.0 % stop sign
Machine:   0.9 % stop sign

# Jailbrake, leak and injection attacks

AI and cyber

# EU's AI Act



**Unacceptable risk**
Prohibited

Art. 5

Manipulation, deception, social scoring, biometric mass surveillance

**High risk**
Conformity assessment

Art. 6

Critical public services, HR systems, biometric identification, immigration, law enforcement

**Limited risk**
Transparency obligation

Art. 52

Chat bots, biometric categorization, emotion detection, deepfakes

**Minimal risk**
Ethical guidelines

Art. 69

Games, data collection, process automation, translations, decision support

# EU's AI Act - timeline

- Proposed by the European Commission,     21 April 2021
- Adopted by the European Parliament,       13 March 2024
- Adopted by the Council of Europe,         21 May 2024
- Implemented                               June 2024
- Enforcement of prohibited AI,             December 2024
- Introduction of conformity assessment,    March 2025
- Enforcement of general AI,                June 2025
- Enforcement of high-risk AI,              June 2027