# Hacking AI: Towards Intelligent Machines that Humans Can Trust

Battista Biggio

Department of Electrical and Electronic Engineering
University of Cagliari, Italy

# The Artificial Intelligence Revolution

AI is going to transform industry and business as electricity did about a century ago
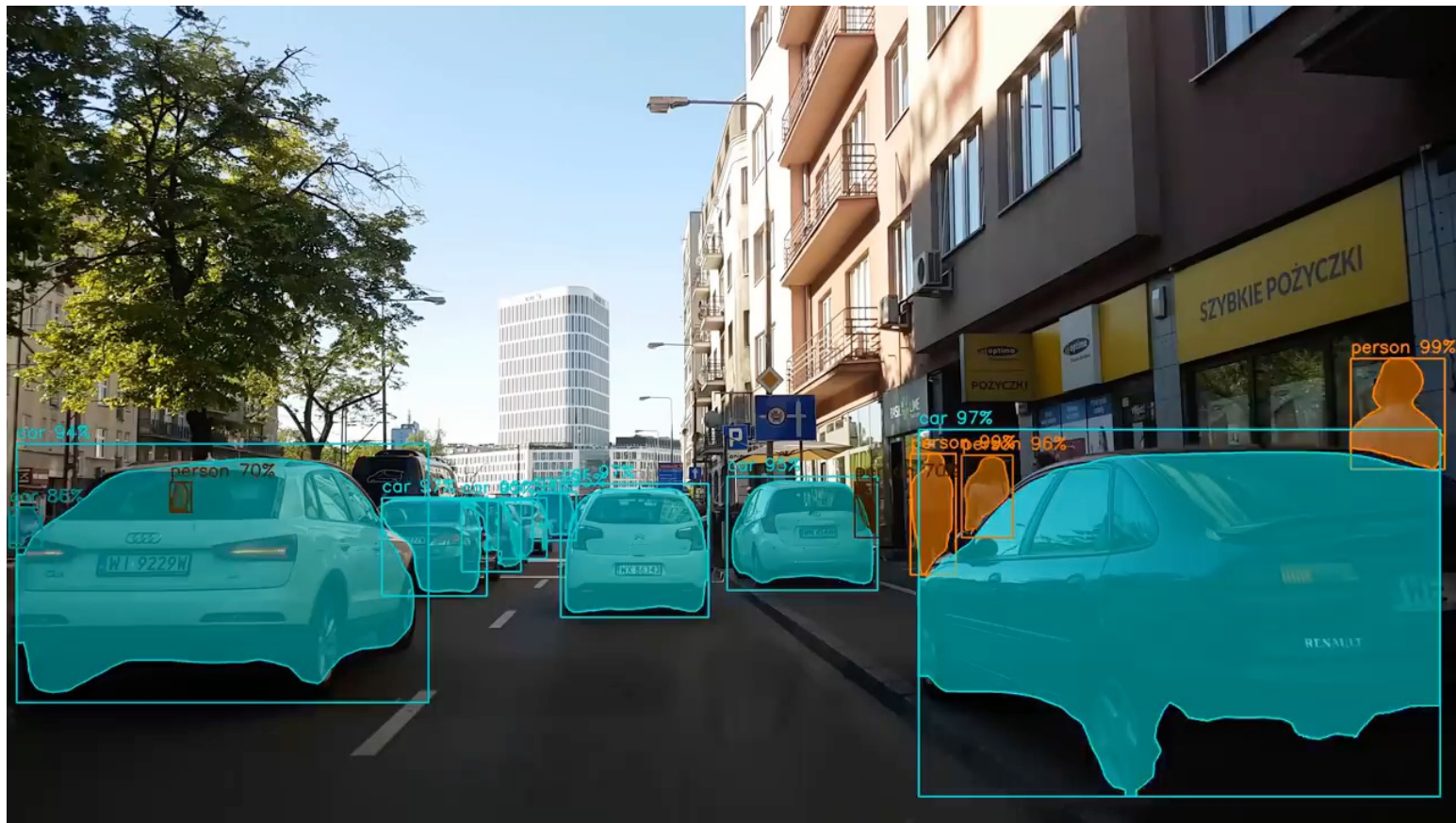     (*Andrew Ng, Jan. 2017*)

**Applications:**

- Computer Vision
- Speech recognition
- Cybersecurity
- Robotics
- Healthcare
- ...

# Computer Vision for Self-Driving Cars

He et al., *Mask R-CNN*, ICCV '17

# Speech Recognition for Virtual Assistants



**Amazon Alexa**

**Apple Siri**

Hey Cortana

**Microsoft Cortana**
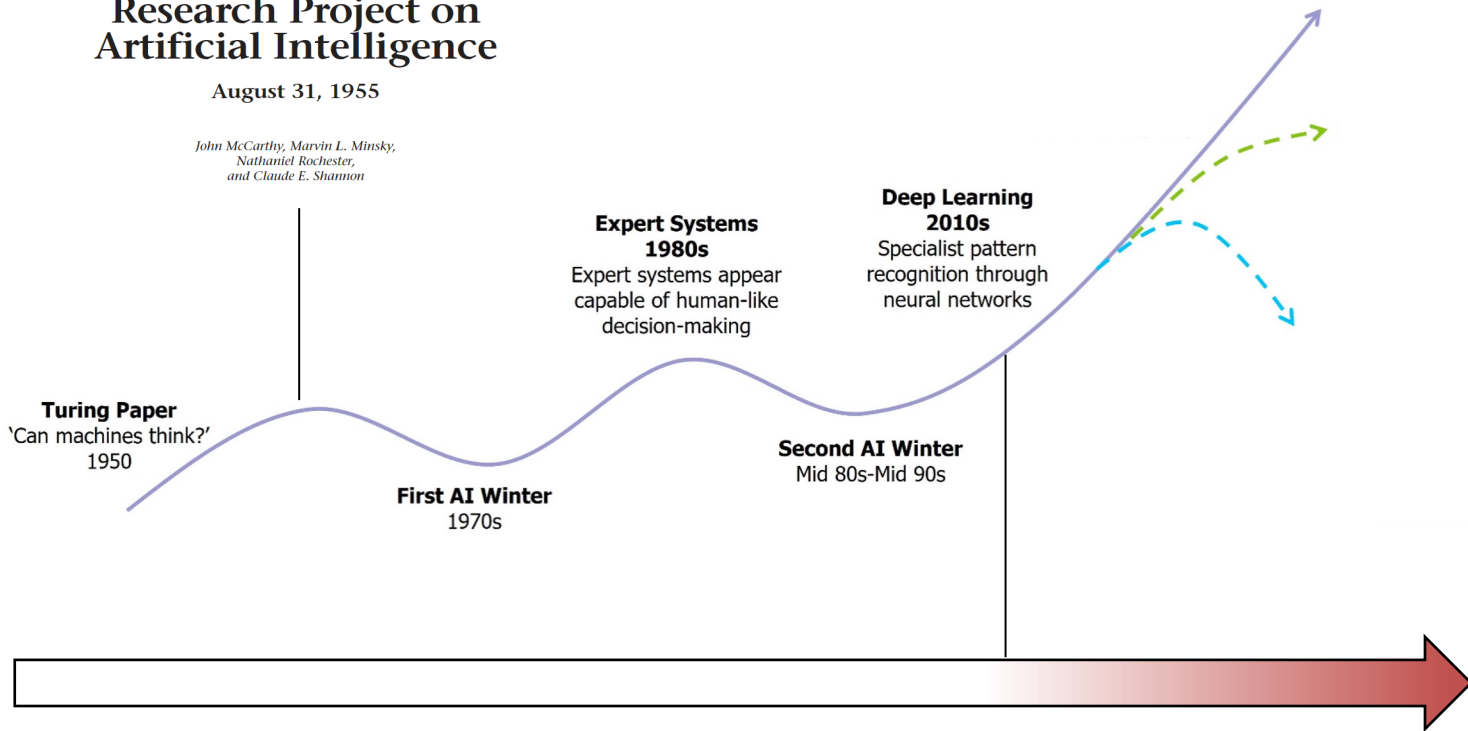
Hi, how can I help?

**Google Assistant**

How Is That **Possible**?

**A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence**

August 31, 1955

*John McCarthy, Marvin L. Minsky, Nathaniel Rochester, and Claude E. Shannon*

**Turing Paper**
'Can machines think?'
1950

**First AI Winter**
1970s

**Expert Systems 1980s**
Expert systems appear capable of human-like decision-making

**Second AI Winter**
Mid 80s-Mid 90s

**Deep Learning 2010s**
Specialist pattern recognition through neural networks

… from the idea of mimicking human reasoning to learning from examples (*machine/deep learning*)

# Why Now?

## Data + Computing Power
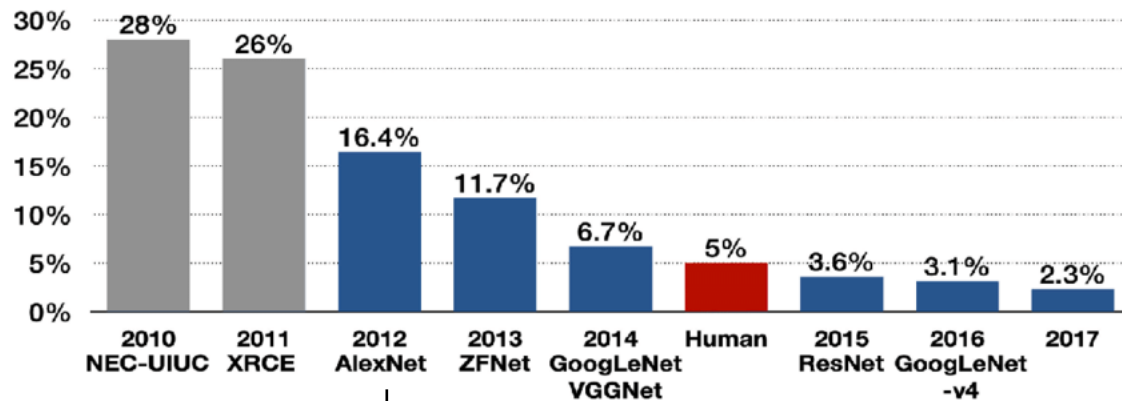
# ImageNet Large Scale Visual Recognition Challenge

IM🅰GENET

1,2M training images
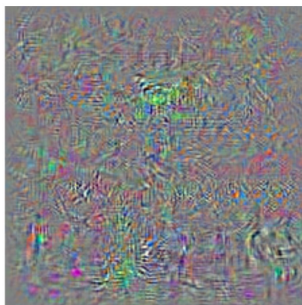1,000 classes

Fei-Fei Li



Top-5 error

| 2010 NEC-UIUC | 28% |
| 2011 XRCE | 26% |
| 2012 AlexNet | 16.4% |
| 2013 ZFNet | 11.7% |
| 2014 GoogLeNet VGGNet | 6.7% |
| Human | 5% |
| 2015 ResNet | 3.6% |
| 2016 GoogLeNet-v4 | 3.1% |
| 2017 | 2.3% |

Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton

**Do Neural Networks Learn Like Us?**

# How Can We Trick Them?



school bus (94%)

ostrich (97%)
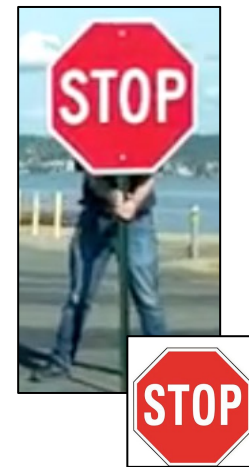
*Biggio et al.*, Evasion attacks against machine learning at test time, **ECML-PKDD 2013**
*Szegedy et al.*, Intriguing properties of neural networks, **ICLR 2014**

# Adversarial Glasses



M. Sharif, L. Baio et al., *Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition*, ACM CCS 2016

# Adversarial Road Signs



K. Eykholt, D. Song et al., *Robust physical-world attacks on deep learning visual classification*, CVPR 2018

# Audio Adversarial Examples

| **Audio** | **Transcription by Mozilla DeepSpeech** |
| --- | --- |
| 🔊 | "without the dataset the article is useless" |
| 🔊 | "okay google browse to evil dot com" |

N. Carlini and D. Wagner, *Audio adversarial examples: Targeted attacks on speech-to-text*, DLS 2018

# It's Not *Just* Adversarial Examples…

- Explainability, privacy, fairness…
- Attacks on Large Language Models (LLMs)





A. Zou, Z. Wang, Z. Kolter, M. Fredrikson, *Universal and Transferable Adversarial Attacks on Aligned Language Models*, 2023

**We Should Build _Trustworthy_ Machines, …**
**But Can We?**

**Not yet…**

# The Perils of Regulating AI

- The EU AI Act demands the development of trustworthy AI models
  - This is still an **open research problem**! We do not have valid solutions yet
    - Need for more research and education initiatives
    - Need for an ecosystem/initiative to help companies adopt AI while being compliant with regulations

- We should **regulate AI applications** and **not the technology** itself

- We should not hinder the development of **open-source AI/ML models**
  - This may create a huge gap between big tech companies/providers and the rest of us

- To engender trust in AI/ML, open-source development and open research are of paramount importance