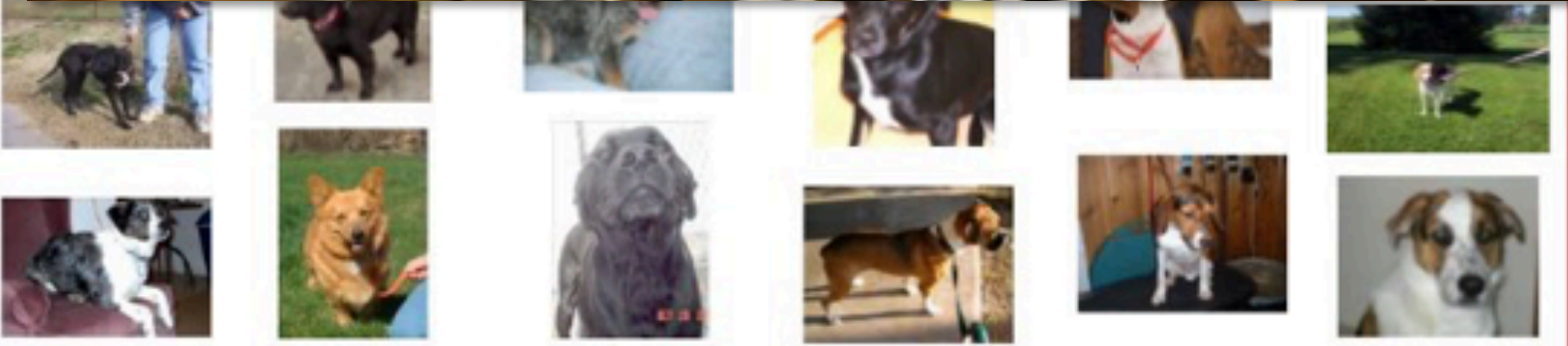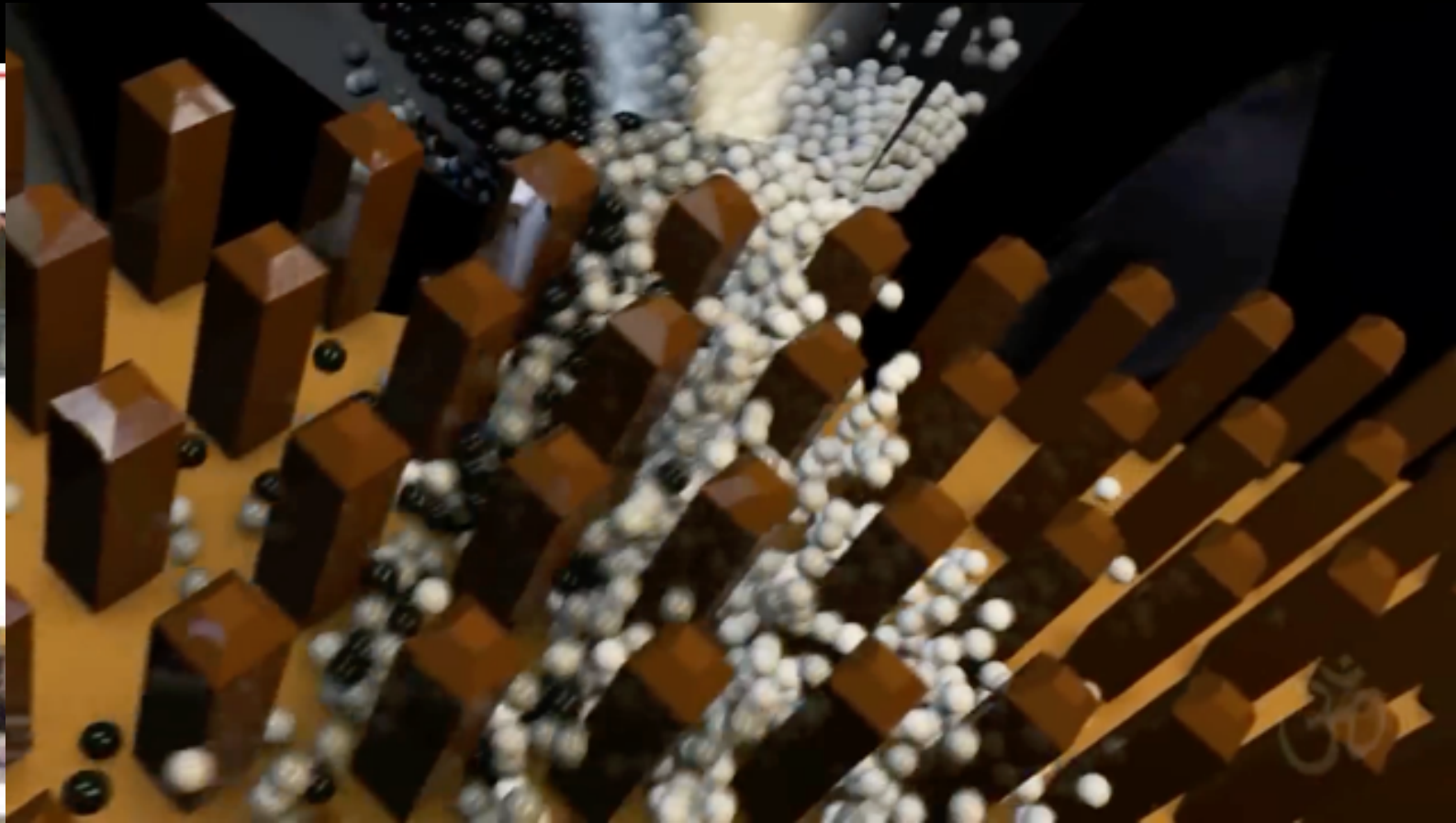# machine learning: next ethical steps

Patrick van der Smagt
Machine Learning Research Lab
Volkswagen Group, Munich

Eövtös Loránd University, Budapest
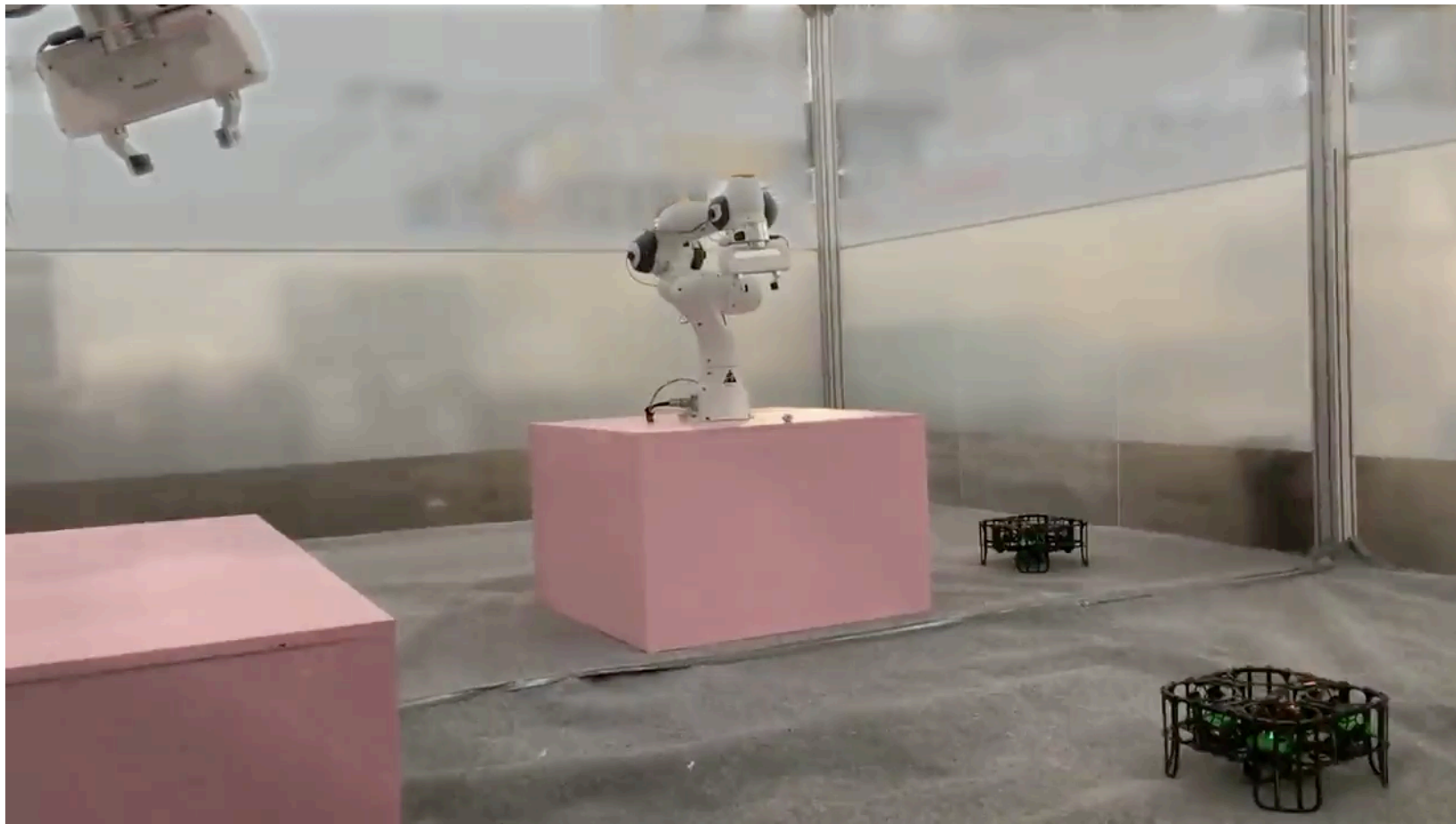Ludwig Maximilian University, Munich

"Artificial Intelligence" in 2020:
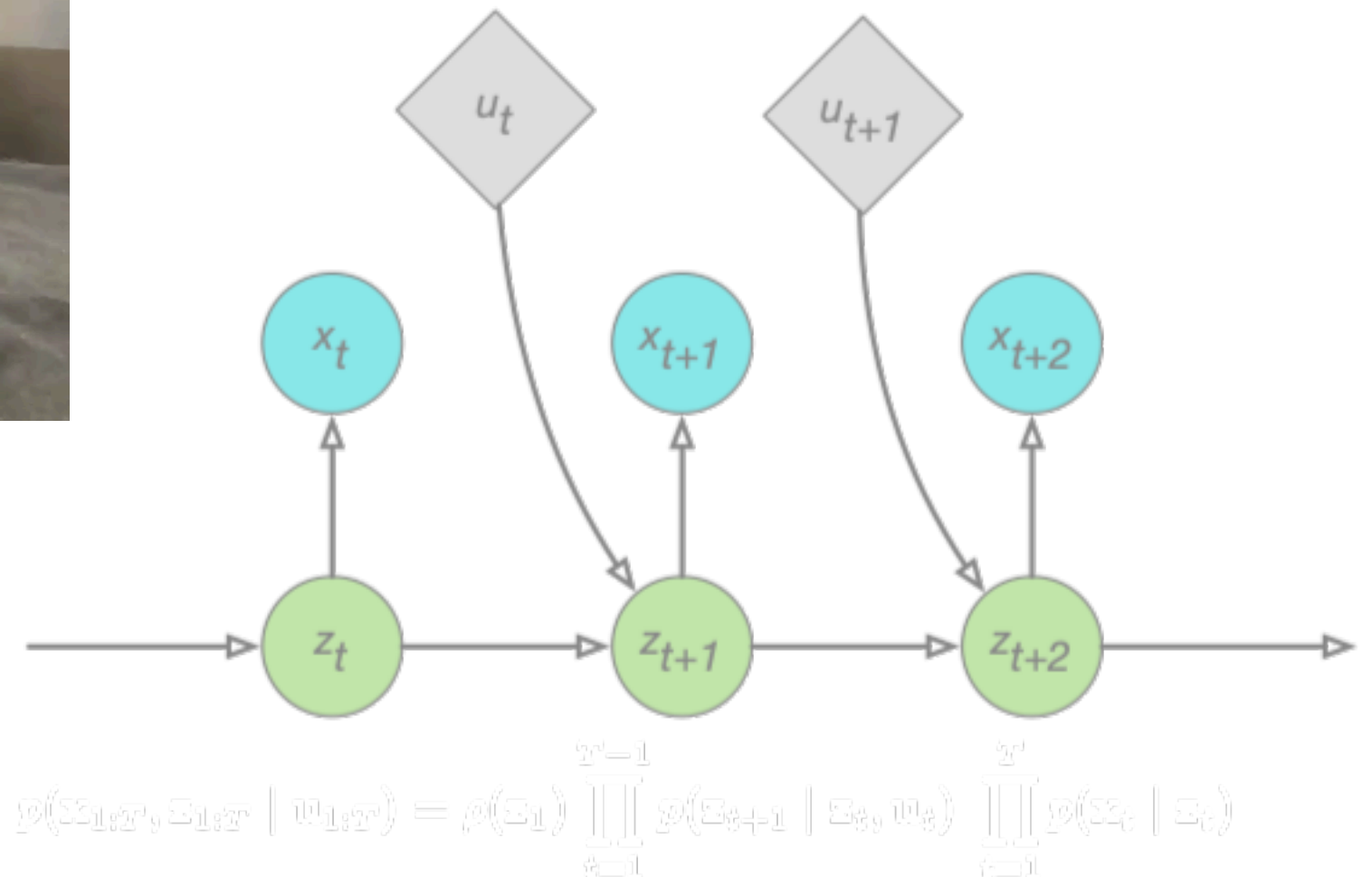**the great success of supervised machine learning**

# control with neural networks

probabilistic
neural networks
learn to
predict the future....

# …to learn a map of the environment…

# 2015, Google Photo scandal

black people were categorised as "gorillas"

one of the first AI public scandals

Google acknowledged, and promised to fix

# 2016, Microsoft twitter chatbot

# gender bias in images

# 2017, Amazon reported on its discriminating HR tool

Amazon HR used 10 years of CVs / hiring decisions to train the algorithm

https://www.theguardian.com/technology/2018/oct/10/amazon-hiring-ai-gender-bias-recruiting-engine

https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G

## 2018, Uber's deadly crash

A jaywalking pedestrian was lethally hit by a self-driving car while pushing her bike

Not only data can go wrong!

Allegedly, Chinese startups and government are using AI to automate racial profiling and persecution of the Uyghurs minority

https://www.nytimes.com/2019/04/14/technology/china-surveillance-artificial-intelligence-racial-profiling.html

# ethical framework for trustworthy AI

advocates that AI should be…

# ethical framework for trustworthy AI

advocates that AI should be…

self-assessment
is no
lasting solution

despite many attempted initiatives...

scandals continue happening

As we reach the limits of self-assessed ethical AI...

the public loses trust in companies and institutions

etami:

*ethical and **t**rustworthy*

*a**rtificial and **m**achine **i**ntelligence*

founded on HLEG report (April 8, 2019)

52 participants

4 of which are in etami

# etami

**Goal:** development and foundation of a European Organisation for the Conformity Assessment of Ethical and Trustworthy Artificial and Machine Intelligence.

**Planned results:**
- Conformity assessment criteria for ethical artificial intelligence based on the "The Ethics Guidelines for Trustworthy Artificial Intelligence" (published in 4/2019 by the EC).
- Certification methods, tools, corresponding guidelines for industrial process organisation
- Create/Support a standard for the conformity assessment of AI.
- Foundation of an organisation for standardisation and certification of Ethical and Trustworthy Artificial and Machine Intelligence.

**Project partners:**
- *academic:* Eötvös Loránd University; KU Leuven; TU Berlin; TU Munich; UnternehmerTUM
- *corporate:* ABB; ATOS; AVL; Continental; Deutsche Bahn; Leonardo; Poste Italiane; Siemens; Volkswagen Group

**Project lead:**
Prof. Dr. Patrick van der Smagt, Volkswagen Group Machine Learning Research Lab
Dr. Djalel Benbouzid, Volkswagen Group Machine Learning Research Lab

**we advocate third-party assessment**

**trustworthiness by design**

third-party assessment

avoid conflicts of interest

benefit from the public trust

**"do not reinvent the wheel" principle**

existing HLEG ethical guidelines

base on already-existing standards

**novel processes and certification**

novel AI audit and certification

# the landscape

**rule enforcement** (vertical axis) / **rule making** (horizontal axis)

|  | academia | companies | standard agency/ institution | policy makers |
|---|---|---|---|---|
| **law** | -/- | -/- | -/- | EU, national governments (Estonia, ...) |
| **standard** | -/- | TÜV, SGS, Bureau Verita, CIO Strategy Council | IEEE SA, CEN-CENELEC ISO/IEC, DIN, VDE  etami | Malta certification, Denmark responsibility seal, Australia certification scheme |
| **audit** | TBD | KPMG, The Institute of Internal Auditors (IIA), Deloitte, … | TÜV, SGS, Bureau Veritas, IIA  etami | -/- |
| **self ass-essment** | Future of Humanity Institute (Oxford), OpenAI | PAI (Partner-ship on AI), Asilomar Conference on Beneficial AI, Open AI | TBD | -/- |

OCEANIS

etami today

# the landscape

**rule enforcement**

| | academia | companies | standard agency/ institution | policy makers |
|---|---|---|---|---|
| **law** | -/- | -/- | -/- | EU, national governments (Estonia, ...) |
| **standard** | -/- | TÜV, SGS, Bureau Verita, CIO Strategy Council | IEEE SA, CEN-CENELEC ISO/IEC, DIN, VDE<br><br>etami | Malta certification, Denmark responsibility seal, Australia certification scheme |
| **audit** | TBD | KPMG, The Institute of Internal Auditors (IIA), Deloitte, … | TÜV, SGS, Bureau Veritas, IIA<br><br>etami | -/- |
| **self ass-essment** | Future of Humanity Institute (Oxford), OpenAI | PAI (Partner-ship on AI), Asilomar Conference on Beneficial AI, Open AI | TBD | -/- |

OCEANIS

**rule making**

etami planned