# Reflecting opportunities and risks – AI ethics and its broad range of issues
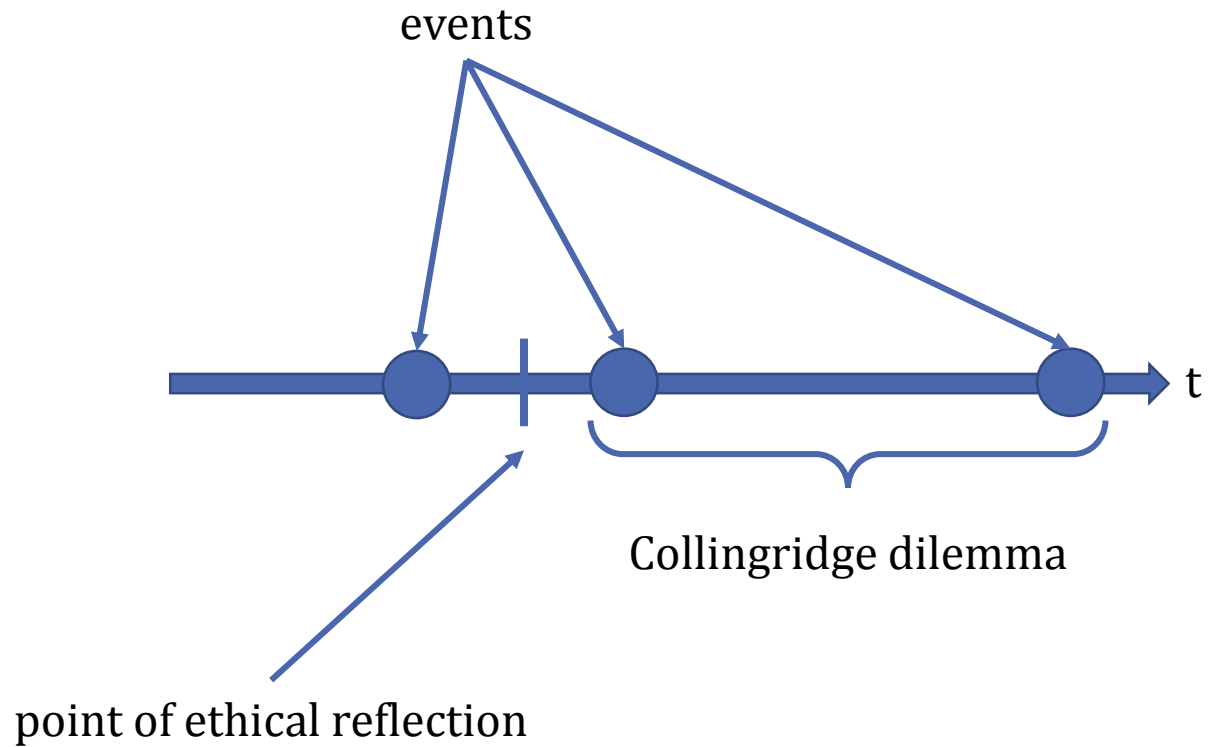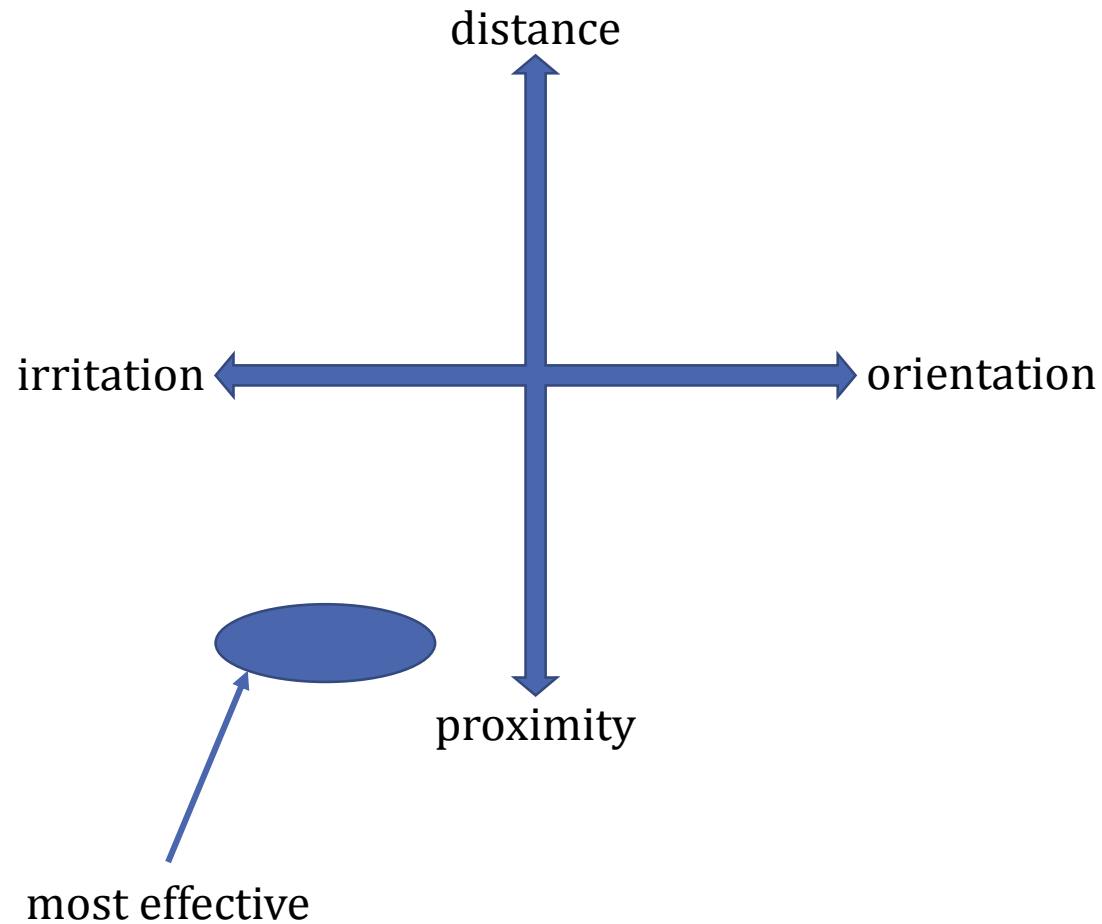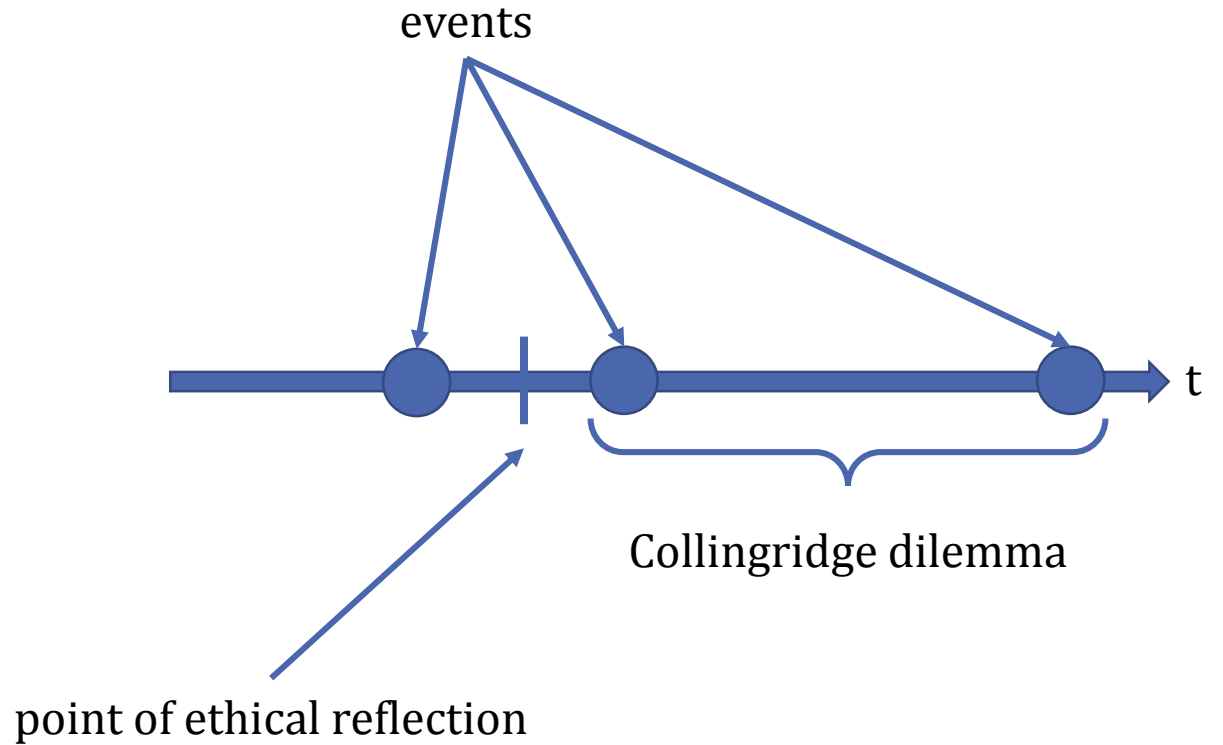
**Dr. Thilo Hagendorff**

*University of Tuebingen*
*Cluster of Excellence Machine Learning*

# Dimensions of AI ethics

distance

irritation

orientation

proximity

most effective

events

t

Collingridge dilemma

point of ethical reflection

# Dimensions of AI ethics

# Dimensions of AI ethics

# Dimensions of AI ethics

distance

irritation ← → orientation

proximity

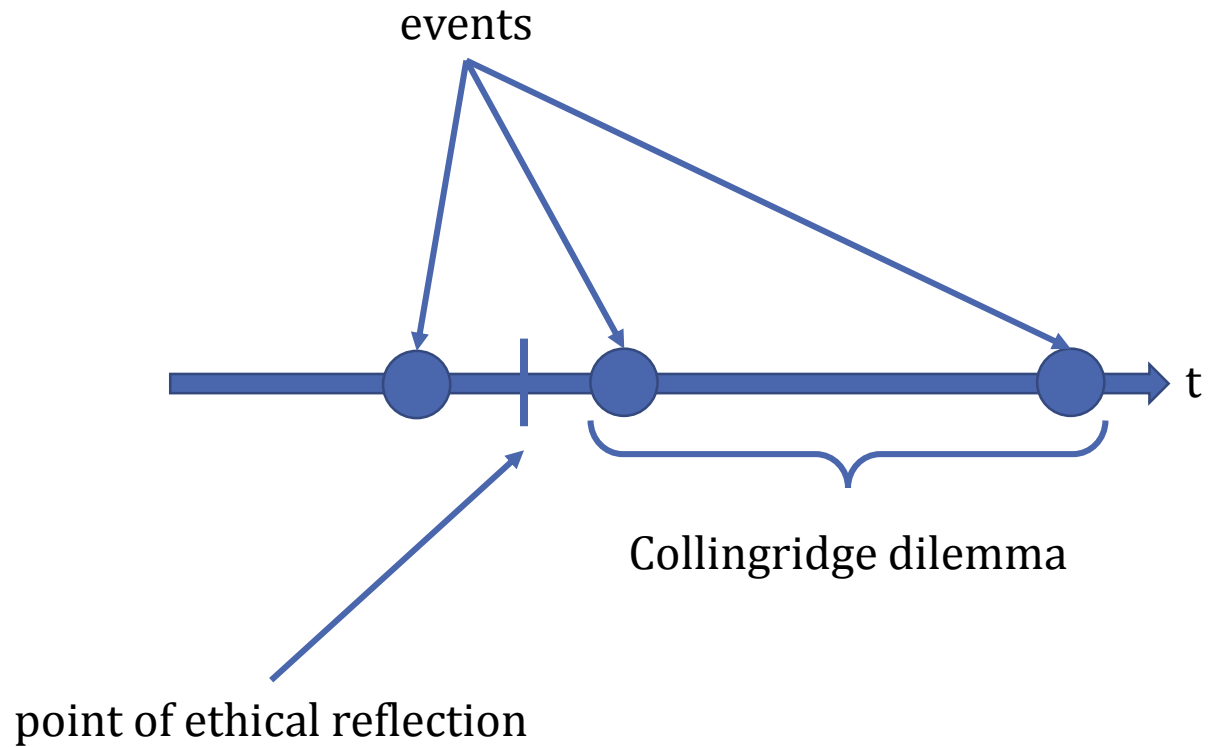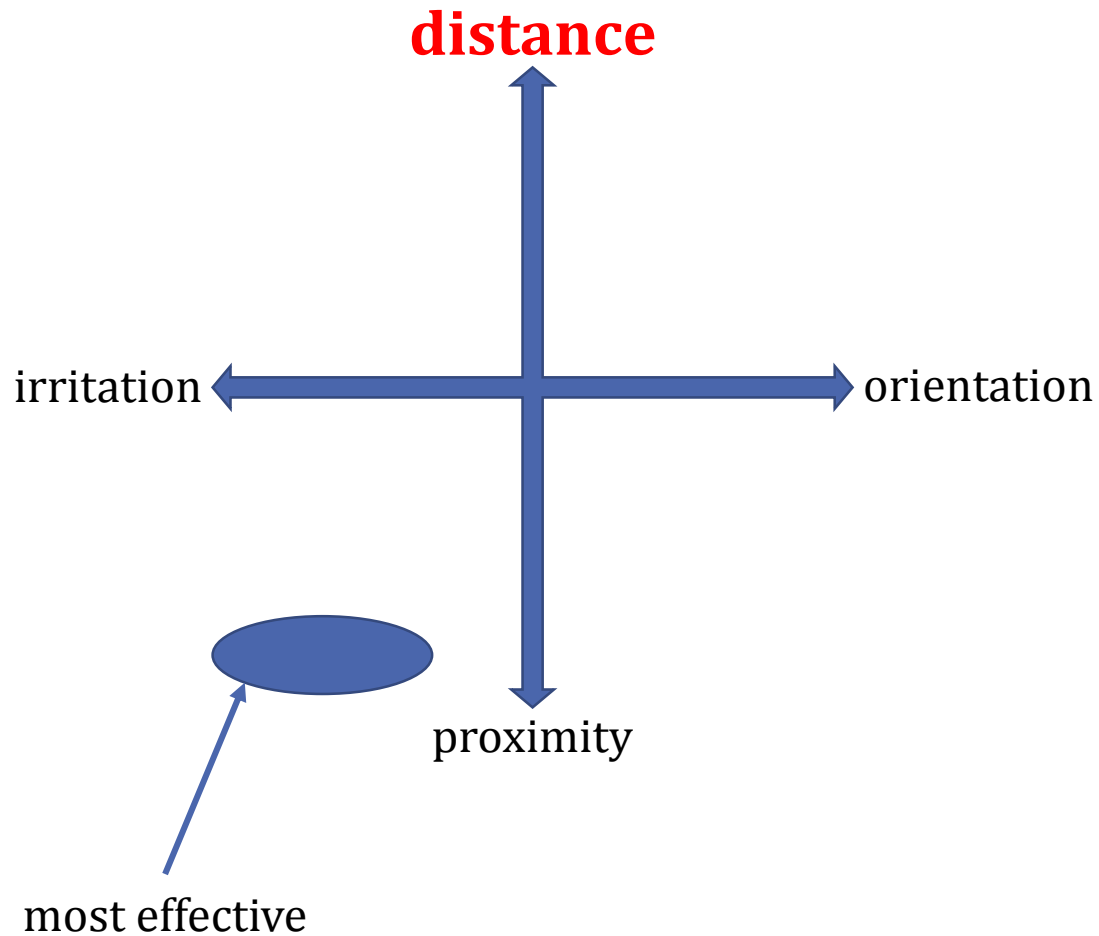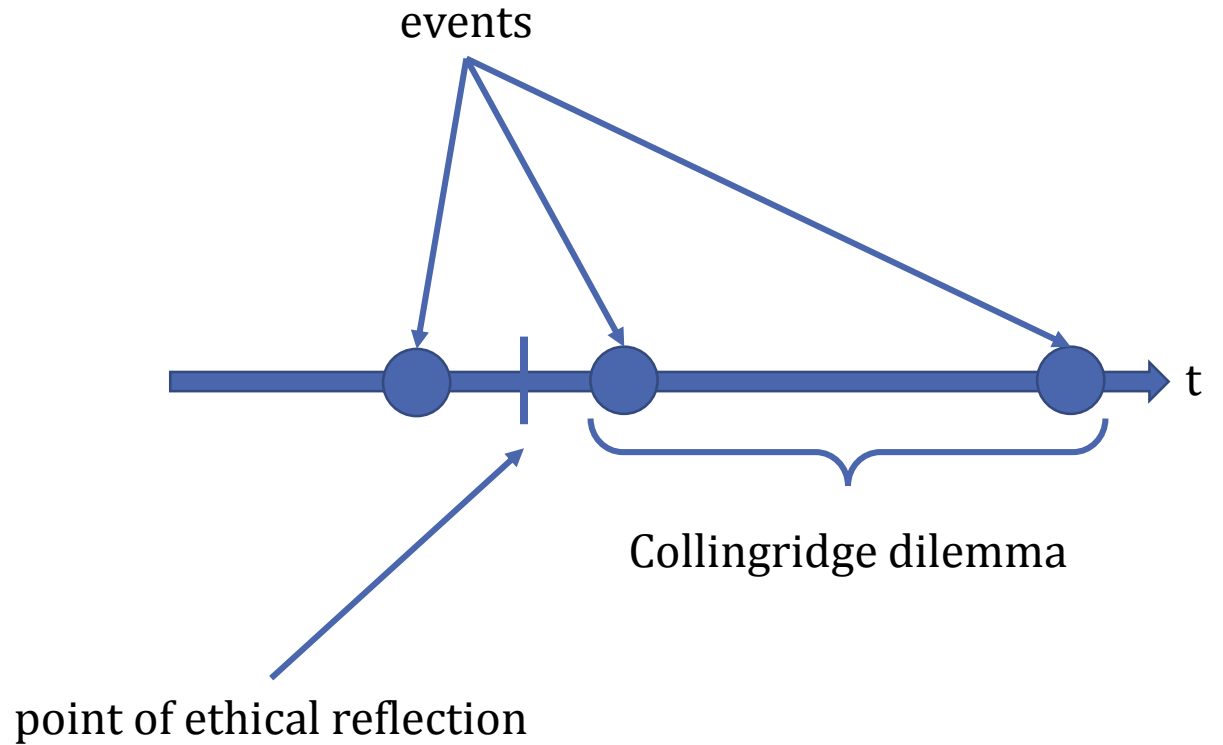most effective

events

t

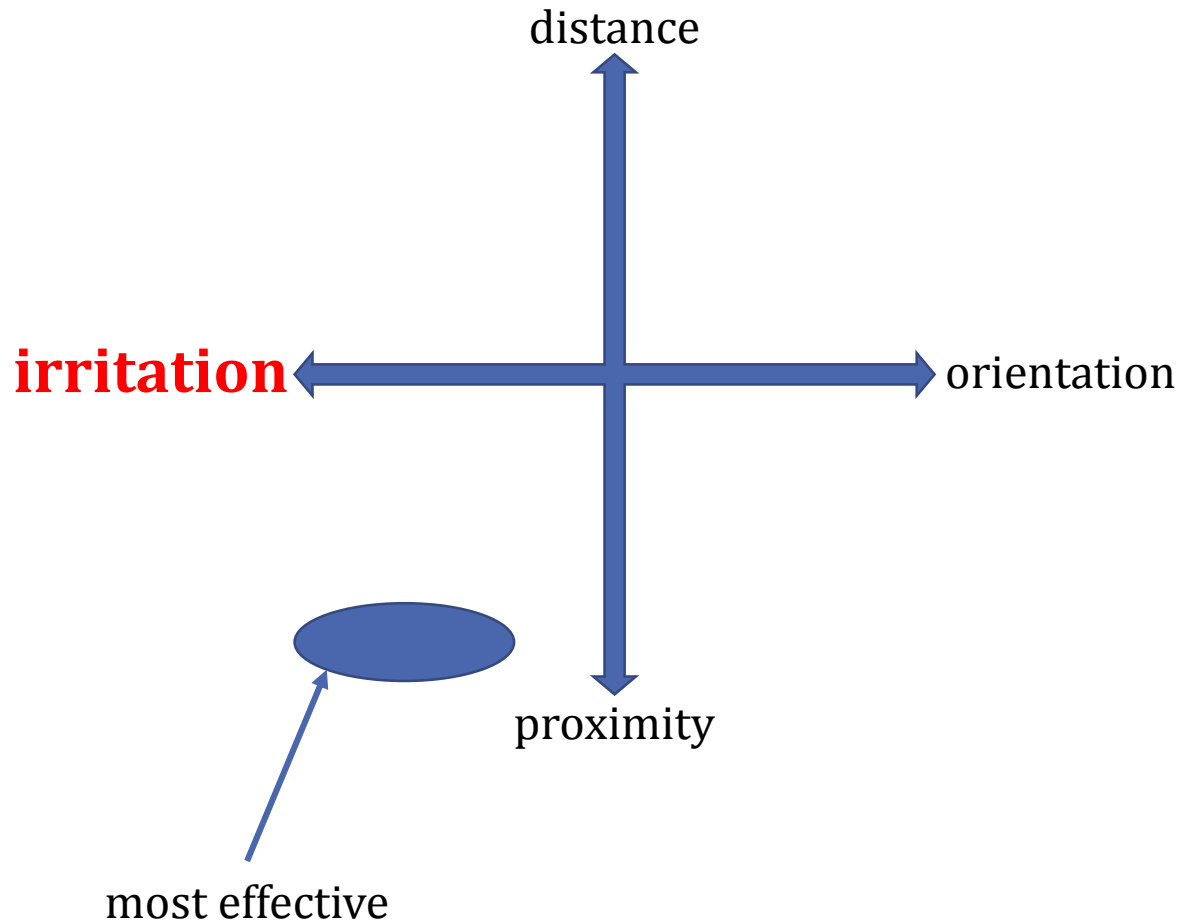point of ethical reflection

Collingridge dilemma

# Dimensions of AI ethics

# Dimensions of AI ethics

# Dimensions of AI ethics

distance

irritation ← → orientation

proximity

most effective

events

point of ethical reflection

Collingridge dilemma

t

# Dimensions of AI ethics

# Demands for AI/ML Ethics



**The future of AI relies on a code of ethics**

Matthew Howard  @mattdhoward  /  10 months ago

NETZPOLITIK.ORG

Technologie

**Keine roten Linien: Industrie entschärft Ethik-Leitlinien für Künstliche Intelligenz**

WIRED  SUBSCRIBE

TOM SIMONITE  BUSINESS  05.16.18  04:32 PM

**TECH FIRMS MOVE TO PUT ETHICAL GUARD RAILS AROUND AI**

Forbes

3,466 views  |  Mar 27, 2019, 01:24pm

**The Growing Marketplace For AI Ethics**

Forbes Insights  **Insights Team** Insights Contributor
FORBES INSIGHTS With Intel AI

Google Trend for "AI ethics"

100 %

50 %

2013   2014   2015   2016   2017   2018

# Guidelines

- Hagendorff, Thilo (2019): The Ethics of AI Ethics. An Evaluation of Guidelines. in: arXiv:1903.03425v1, pp. 1–15.

## The Ethics of AI Ethics
### An Evaluation of Guidelines

Dr. Thilo Hagendorff

University of Tuebingen
International Center for Ethics in the Sciences and Humanities
thilo.hagendorff@uni-tuebingen.de

**Abstract** - Current advances in research, development and application of artificial intelligence (AI) systems have yielded a far-reaching discourse on AI ethics. In consequence, a number of ethics guidelines have been released in recent years. These guidelines comprise normative principles and recommendations aimed to harness the "disruptive" potentials of new AI technologies. Designed as a comprehensive evaluation, this paper analyzes and compares these guidelines highlighting overlaps but also omissions. As a result, I give a detailed overview of the field of AI ethics. Finally, I also examine to what extent the respective ethical principles and values are implemented in the practice of research, development and application of AI systems – and how the effectiveness in the demands of AI ethics can be improved.

**Keywords** - artificial intelligence, machine learning, ethics, guidelines, implementation

## 1 Introduction

The current AI boom is accompanied by constant calls for applied ethics, which are meant to harness the "disruptive" potentials of new AI technologies. As a result, a whole body of ethical guidelines has been developed in recent years collecting principles, which technology developers should adh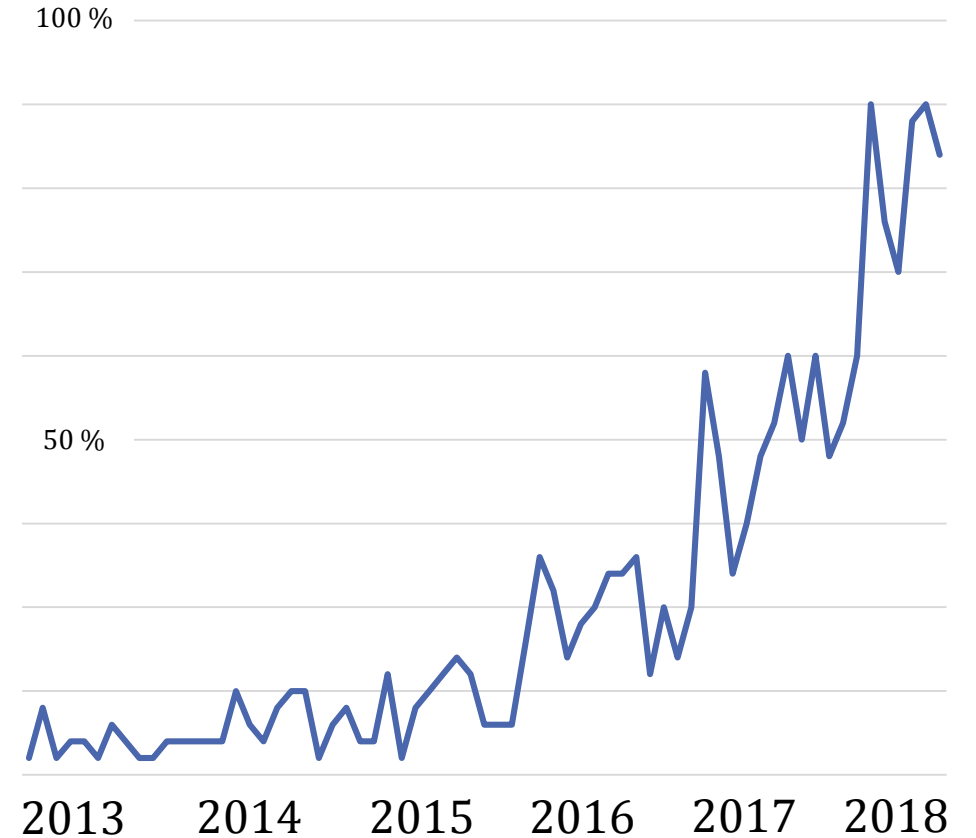ere to as far as possible. However, the critical question arises: Do those ethical guidelines have an actual impact on human decision-making in the field of AI and machine learning? The short answer is: No, most often not. This paper analyzes fifteen of the major AI ethics guidelines and issues recommendations on how to overcome the relative ineffectiveness of these guidelines.

AI ethics – or ethics in general – lacks mechanisms to reinforce its own normative claims. Of course, the enforcement of ethical principles may involve reputational losses in the case of misconduct, or restrictions on memberships in certain professional bodies. Yet altogether, these mechanisms are rather weak and pose no eminent threat. Researchers, politicians, consultants, managers and activists have to deal with this essential weakness of ethics. However, it is also a reason why ethics is so appealing to many AI companies and institutions. When companies or research institutes formulate their own ethical guidelines, regularly incorporate ethical considerations into their public relations work, or adopt ethically motivated "self-commitments", efforts to create a truly binding legal framework are continuously discouraged. Ethics guidelines of the AI industry serve to suggest to legislators that internal self-governance in science and industry is sufficient, and that no specific laws are necessary to mitigate possible technological risks and to eliminate scenarios of abuse (Calo 2017). And even when more concrete laws concerning AI systems are demanded, as recently done by Google (Google 2019), these demands remain relatively vague and superficial.

Science- or industry-led ethics guidelines, as well as other concepts of self-governance, may serve to pretend that accountability can be devolved from state authorities and democratic institutions upon the respective sectors of science or industry. Moreover, ethics can also simply serve the purpose of calming critical voices from the public, while simultaneously the criticized practices are maintained within the organization. The association "Partnership on AI" (2018) which brings together companies such as Amazon, Apple, Baidu, Facebook, Google, IBM and Intel is exemplary in this context. Companies can highlight their membership in such associations whenever the

1

# Guidelines

| ethical aspect | The European Commission's High-Level Expert Group on Artificial Intelligence | The Malicious Use of Artificial Intelligence | AI4People | The Asilomar AI Principles | AI Now 2016 Report | AI Now 2017 Report | AI Now 2018 Report | Principles for Accountable Algorithms and a Social Impact Statement for Algorithms | Montréal Declaration for Responsible Development of Artificial Intelligence | Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems | ITI AI Policy Principles | Microsoft AI principles | Artificial Intelligence at Google | Everyday Ethics for Artificial Intelligence | Partnership on AI | number of mentions |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| privacy protection | ● | ● | ● | ● | ● | ● | ● | | ● | ● | ● | ● | ● | ● | ● | 14 |
| accountability | ● | ● | ● | ● | ● | | ● | ● | ● | ● | ● | ● | ● | | ● | 13 |
| fairness, non-discrimination, justice | ● | | ● | ● | ● | | ● | ● | ● | ● | ● | ● | ● | | ● | 13 |
| transparency, openness | ● | ● | | ● | ● | ● | ● | | ● | ● | ● | | | | ● | 10 |
| safety, cybersecurity | ● | ● | ● | ● | ● | | | | ● | | ● | | ● | ● | ● | 10 |
| common good, sustainability | | | ● | ● | ● | ● | ● | | ● | ● | | | | ● | ● | 9 |
| explainability, interpretabiliy | ● | | ● | ● | ● | ● | ● | | ● | ● | | | | | | 8 |
| human oversight, control, auditing | ● | | ● | | ● | | ● | | ● | ● | ● | | ● | | | 8 |
| dual-use problem, military, AI arms race | | ● | | | | ● | ● | | | ● | ● | ● | | | | 6 |
| solidarity, inclusion, social cohesion | | | ● | ● | ● | ● | | | ● | | | ● | | | ● | 6 |
| science-policy link | | ● | | | | | ● | | ● | | ● | | | ● | | 5 |
| field-specific deliberations (health, military, mobility etc.) | | ● | | | ● | ● | ● | | | ● | | | | | | 5 |
| diversity in the field of AI | | | ● | ● | ● | | ● | ● | ● | | | | | | | 5 |
| public awareness, education about AI and its risks | | | ● | ● | | | | ● | | ● | | | ● | | ● | 5 |
| future of employment | | | ● | ● | | | | ● | | ● | | ● | | | | 4 |
| human autonomy | ● | | | | ● | | | | | ● | | ● | | | | 4 |
| protection of whistleblowers | | | | | | | ● | | | | | | | | | 1 |
| hidden costs (labeling, clickwork, material resources etc.) | | | | | | | ● | | | | | | | | | 1 |
| affiliation (government, industry, science) | government | science | science | science | science | science | science | science | science | industry | industry | industry | industry | industry | industry | |
| number of ethical aspects | 8 | 7 | 11 | 11 | 11 | 9 | 11 | 5 | 11 | 10 | 8 | 6 | 6 | 5 | 8 | |

# Guidelines

| Ethical aspect | Count |
|---|---|
| privacy protection | 14 |
| accountability | 13 |
| fairness, non-discrimination, justice | 13 |
| transparency, openness | 10 |
| safety, cybersecurity | 10 |
| common good, sustainability | 9 |
| explainability, interpretabiliy | 8 |
| human oversight, control, auditing | 8 |
| dual-use problem, military, AI arms race | 6 |
| solidarity, inclusion, social cohesion | 6 |
| science-policy link | 5 |
| field-specific deliberations (health, military, mobility etc.) | 5 |
| diversity in the field of AI | 5 |
| public awareness, education about AI and its risks | 5 |
| future of employment | 4 |
| human autonomy | 4 |
| protection of whistleblowers | 1 |
| hidden costs (labeling, clickwork, material resources etc.) | 1 |

Columns of the source matrix:
- The European Commission's High-Level Expert Group on Artificial Intelligence
- The Malicious Use of Artificial Intelligence
- AI4People
- The Asilomar AI Principles
- AI Now 2016 Report
- AI Now 2017 Report
- AI Now 2018 Report
- Principles for Accountable Algorithms and a Social Impact Statement for Algorithms
- Montréal Declaration for Responsible Development of Artificial Intelligence

affiliation (government, industry, science): government, science, science, science, science, science, science, science, science

number of ethical aspects: 8, 7, 11, 11, 11, 9, 6, 11

# Guidelines

| guidelines | The European Commission's High-Level Expert Group on Artificial Intelligence | The Malicious Use of Artificial Intelligence | AI4People | The Asilomar AI Principles | AI Now 2016 Report | AI Now 2017 Report | AI Now 2018 Report | Principles for Accountable Algorithms and a Social Impact Statement for Algorithms | Montréal Declaration for Responsible Development of Artificial Intelligence | Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems | ITI AI Policy Principles | Microsoft AI principles | Artificial Intelligence at Google | Everyday Ethics for Artificial Intelligence | Partnership on AI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| notes on technical implementations | yes, but very few and superficial | yes, relatively comprehensive | none | none | none | none | none | none | none | yes, but very few and superficial | none | none | none | none | none |
| proportion of women among authors (f/m) | (8/10) | (5/21) | (5/8) | ns | (4/2) | (3/1) | (6/4) | (1/12) | (8/10) | ns | ns | ns | ns | (1/2) | ns |
| length (number of words) | 16,546 | 34,017 | 8,609 | 646 | 11,530 | 18,273 | 25,759 | 13.59 | 4,754 | 40,915 | 2,272 | 75 | 882 | 4,488 | 1,481 |
| affiliation (government, industry, science) | government | science | science | science | science | science | science | science | science | industry | industry | industry | industry | industry | industry |
| number of ethical aspects | 8 | 7 | 11 | 11 | 11 | 9 | 11 | 5 | 11 | 10 | 8 | 6 | 6 | 5 | 8 |

| | | | | | | | | | | | | | | | | total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| accountability | | | | | | | | | | | | | | | | 13 |
| fairness, non-discrimination, justice | | | | | | | | | | | | | | | | 13 |
| transparency, openness | | | | | | | | | | | | | | | | 10 |
| safety, cybersecurity | | | | | | | | | | | | | | | | 10 |
| common good, sustainability | | | | | | | | | | | | | | | | 9 |
| explainability, interpretabiliy | | | | | | | | | | | | | | | | 8 |
| human oversight, control, auditing | | | | | | | | | | | | | | | | 8 |
| dual-use problem, military, AI arms race | | | | | | | | | | | | | | | | 6 |
| solidarity, inclusion, social cohesion | | | | | | | | | | | | | | | | 6 |
| science-policy link | | | | | | | | | | | | | | | | 5 |
| field-specific deliberations (health, military, mobility etc.) | | | | | | | | | | | | | | | | 5 |
| diversity in the field of AI | | | | | | | | | | | | | | | | 5 |
| public awareness, education about AI and its risks | | | | | | | | | | | | | | | | 5 |
| future of employment | | | | | | | | | | | | | | | | |
| human autonomy | | | | | | | | | | | | | | | | |
| protection of whistleblowers | | | | | | | | | | | | | | | | 1 |
| hidden costs (labeling, clickwork, material resources etc.) | | | | | | | | | | | | | | | | 1 |
| affiliation (government, industry, science) | government | science | science | science | science | science | science | science | science | industry | industry | industry | industry | industry | industry | |
| number of ethical aspects | 8 | 7 | 11 | 11 | 11 | 9 | 11 | 5 | 11 | 10 | 8 | 6 | 6 | 5 | 8 | |

# Guidelines

| Ethical aspect | Count |
|---|---|
| privacy protection | 14 |
| accountability | 13 |
| fairness, non-discrimination, justice | 13 |
| transparency, openness | 10 |
| safety, cybersecurity | 10 |
| common good, sustainability | 9 |
| explainability, interpretabiliy | 8 |
| human oversight, control, auditing | 8 |
| dual-use problem, military, AI arms race | 6 |
| solidarity, inclusion, social cohesion | 6 |
| science-policy link | 5 |
| field-specific deliberations (health, military, mobility etc.) | 5 |
| diversity in the field of AI | 5 |
| public awareness, education about AI and its risks | 5 |
| future of employment | 4 |
| human autonomy | 4 |
| protection of whistleblowers | 1 |
| hidden costs (labeling, clickwork, material resources etc.) | 1 |

Columns (source guidelines): The European Commission's High-Level Expert Group on Artificial Intelligence; The Malicious Use of Artificial Intelligence; AI4People; The Asilomar AI Principles; AI Now 2016 Report; AI Now 2017 Report; AI Now 2018 Report; Principles for Accountable Algorithms and a Social Impact Statement for Algorithms; Montréal Declaration for Responsible Development of Artificial Intelligence

| affiliation | government | science | science | science | science | science | science | science | science |
| number of ethical aspects | 8 | 7 | 11 | 11 | 11 | 9 | | 5 | 11 |

# Privacy (14/15)

- personality analysis, image recognition, disease prediction etc.

Private traits and attributes are predictable from digital records of human behavior

Michal Kosinski, David Stillwell, and Thore Graepel

Computer-based personality judgments are more accurate than those made by humans

Wu Youyou, Michal Kosinski, and David Stillwell

Deep neural networks are more accurate than humans at detecting sexual orientation from facial images

Facebook language predicts depression in medical records

Johannes C. Eichstaedt, Robert J. Smith, Raina M. Merchant, Lyle H. Ungar, Patrick Crutchley, Daniel Preoţiuc-Pietro, David A. Asch, and H. Andrew Schwartz

Linguistic Features Identify Alzheimer's Disease in Narrative Speech

Kathleen C. Fraser, Jed A. Meltzer and Frank Rudzicz

# Accountability

- who can be held legally responsible?
- AI systems as "e-persons"

# Fairness (13/15)

- algorithmic discrimination
- bias in training data
- solutions provided by FAT ML community

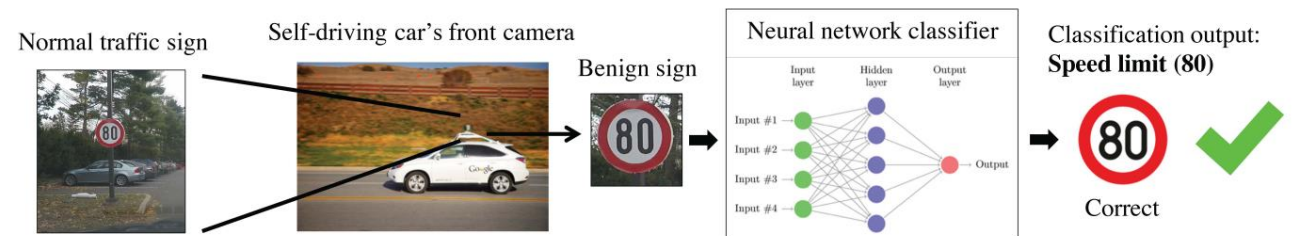# Transparency

- problem of non-transparent organizations dealing with AI
- information asymmetries

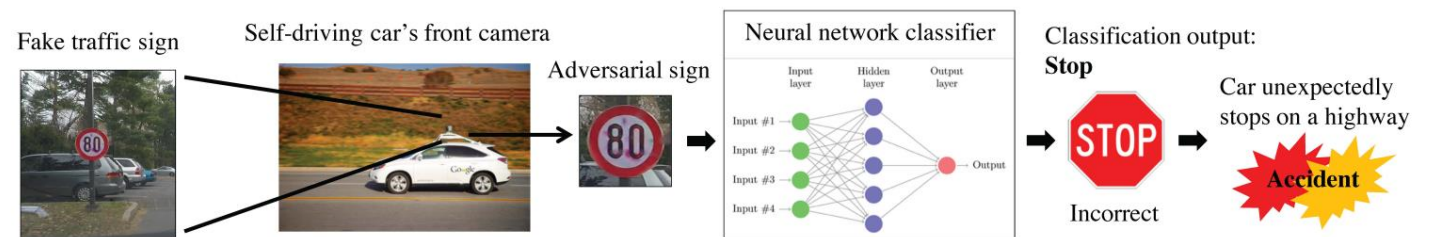# Safety

- dealing with security vulnerabilities
- data poisoning attacks, adversarial examples etc.



(a) Operation of the computer vision subsystem of an AV under *benign conditions*

(b) Operation of the computer vision subsystem of an AV under *adversarial conditions*

# Common good (9/15)

- idea of AI fostering sustainability goals
- AI4Good, Beneficial AI etc.

# Explainability (8/15)

- black box problems
- XAI



input → black box machine learning model → output

Perturbed Instances | P(tree frog)
0.85
0.00001
0.52

Original Image P(tree frog) = 0.54

Locally weighted regression

Query

Explanation

# Human oversight

- developing auditing mechanisms
- human in the loop

# Dual-use problem

- machine learning as "general purpose technology"
- opposing the military use of AI

# Solidarity, social cohesion (6/15)

- AI and social media

- speaking against filter bubbles, micro targeting, radicalization etc.

# Science-policy link (5/15)



- multistakeholder approach
- connecting science, industry and politics

# Field-specific deliberations

- AI in specific social systems or fields
- medicine, military, mobility etc.

# Diversity in the field of AI (5/15)

- diversity crisis in the AI sector

- statistics show blatant inequalities



The Gender Imbalance in AI Research Across 23 Countries

| Country | % MEN | | % WOMEN |
|---|---|---|---|
| Taiwan | 73.91% | | 26.09% |
| Netherlands | 79.17% | | 20.83% |
| France | 85.12% | | 14.88% |
| Denmark | 85.29% | | 14.71% |
| Austria | 85.71% | | 14.29% |
| Japan | 85.71% | | 14.29% |
| China | 85.93% | | 14.07% |
| USA | 86.57% | | 13.43% |
| Singapore | 87.88% | | 12.12% |
| South Korea | 87.93% | | 12.07% |
| Russia | 89.47% | | 10.53% |
| Canada | 89.61% | | 10.39% |
| Italy | 90.00% | | 10.00% |
| Switzerland | 90.28% | | 9.72% |
| Spain | 91.30% | | 8.70% |
| Israel | 91.80% | | 8.20% |
| United Kingdom | 91.81% | | 8.19% |
| Germany | 92.42% | | 7.58% |
| Australia | 92.50% | | 7.50% |
| Belgium | 92.86% | | 7.14% |
| India | 94.44% | | 5.56% |
| Finland | 95.65% | | 4.35% |
| Sweden | 100.00% | | 0.00% |

TOTAL AVERAGE 88% MEN

ELEMENT AI

*Among 4000 researchers who have been published at the leading conferences NIPS, ICML or ICLR in 2017

# Public awareness, education about AI

- creation of educational curricula
  and public awareness activities

# Future of employment

- ideas about robot taxes, universal basic income etc.

# Human autonomy (4/15)

- not using AI for behavior manipulation
- nudging, micro targeting, personalized online advertising, captology, etc.

# Protection of whistleblowers

- need for better protection

# Hidden costs (1/15)

- labeling factories (clickwork), content moderation, energy, material resources etc.

# Guidelines

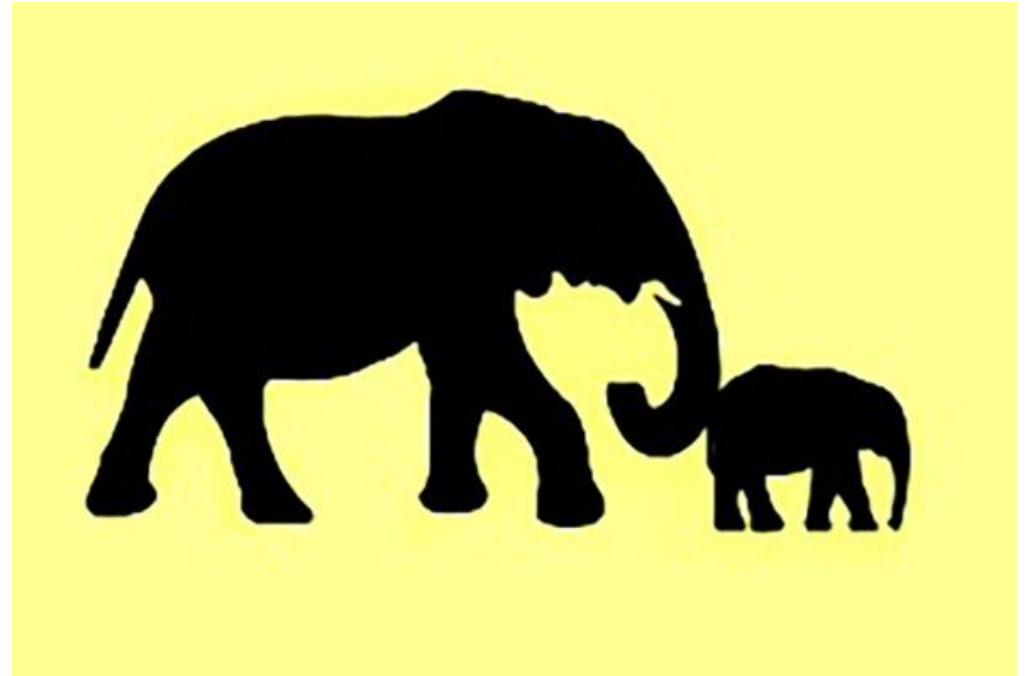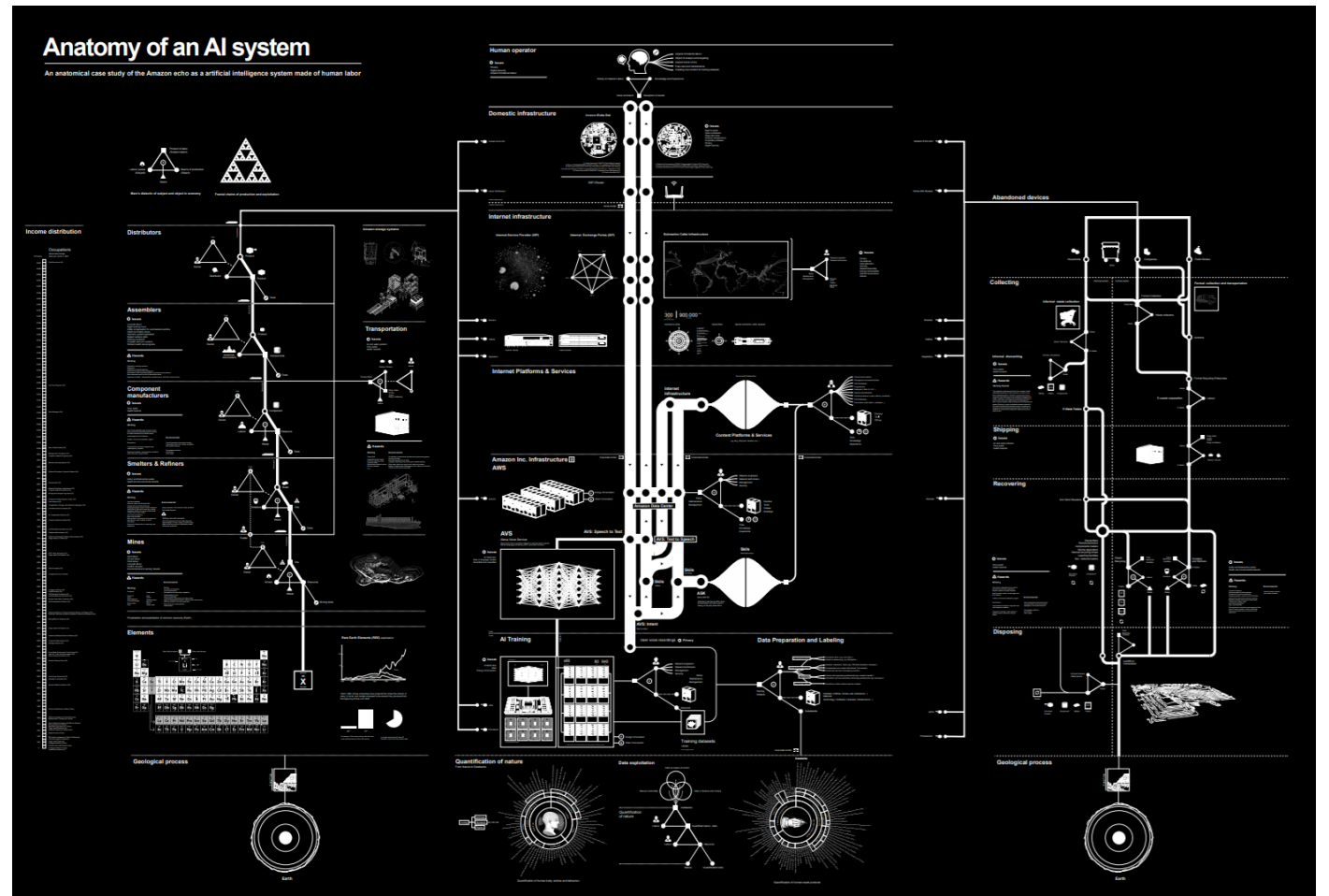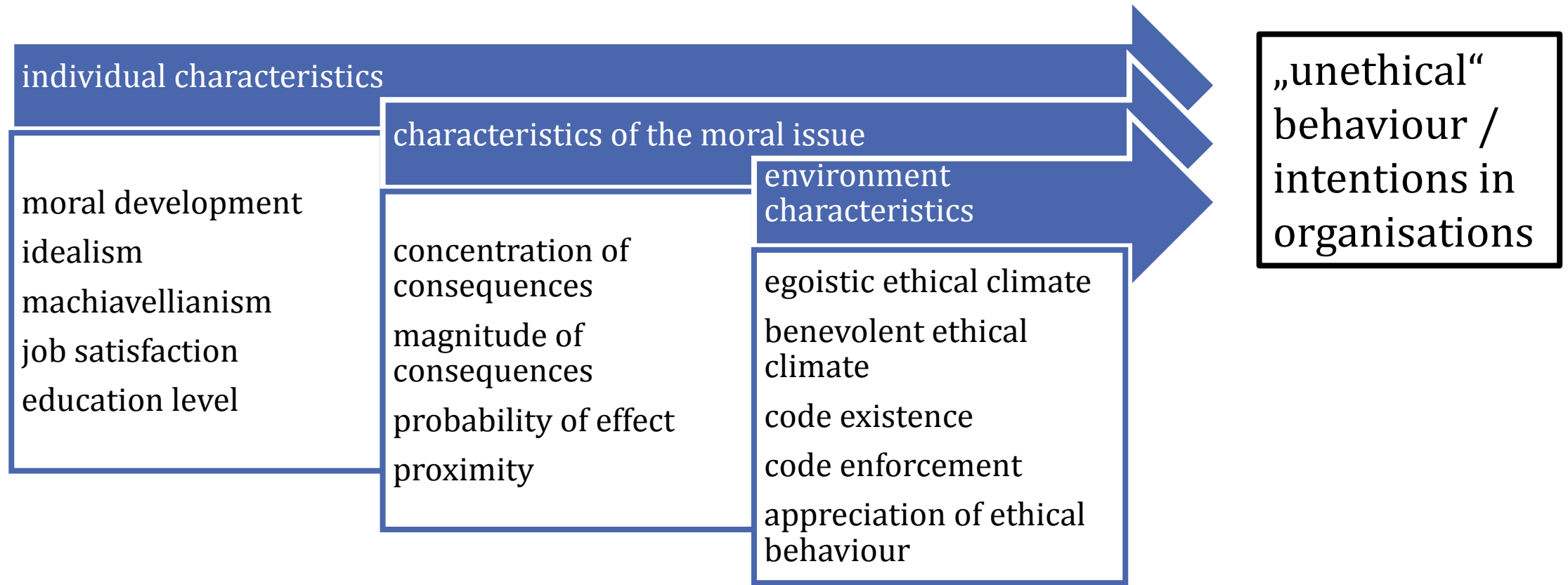| Ethical aspect | The European Commission's High-Level Expert Group on Artificial Intelligence | The Malicious Use of Artificial Intelligence | AI4People | The Asilomar AI Principles | AI Now 2016 Report | AI Now 2017 Report | AI Now 2018 Report | Principles for Accountable Algorithms and a Social Impact Statement for Algorithms | Montréal Declaration for Responsible Development of Artificial Intelligence | Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems | ITI AI Policy Principles | Microsoft AI principles | Artificial Intelligence at Google | Everyday Ethics for Artificial Intelligence | Partnership on AI | number of mentions |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| privacy protection | x | x | x | x | x | x | x |  | x | x | x | x | x | x | x | 14 |
| accountability | x | x | x | x | x |  | x | x | x | x | x | x |  | x | x | 13 |
| fairness, non-discrimination, justice | x |  | x | x | x | x | x | x | x | x | x | x | x | x |  | 13 |
| transparency, openness | x | x |  | x | x | x | x | x | x | x |  |  | x | x |  | 10 |
| safety, cybersecurity | x | x | x | x | x |  |  |  | x | x | x | x | x |  | x | 10 |
| common good, sustainability |  | x | x | x |  |  | x |  | x | x | x |  | x | x | x | 9 |
| explainability, interpretabiliy | x |  | x | x |  | x |  | x | x | x | x |  |  |  |  | 8 |
| human oversight, control, auditing | x |  | x |  | x |  | x | x | x | x |  |  | x |  |  | 8 |
| dual-use problem, military, AI arms race |  | x |  | x |  |  | x |  |  | x |  |  | x |  |  | 6 |
| solidarity, inclusion, social cohesion |  |  | x |  | x |  | x |  | x |  |  | x |  |  | x | 6 |
| science-policy link | x | x |  |  |  | x |  |  |  | x |  |  |  |  |  | 5 |
| field-specific deliberations (health, military, mobility etc.) | x |  |  |  |  |  |  |  |  | x |  |  |  |  |  | 5 |
| diversity in the field of AI |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  | 5 |
| public awareness, education about AI and its risks | x |  | x |  |  |  |  |  |  | x |  |  |  | x |  | 5 |
| future of employment |  |  |  |  |  |  | x |  |  | x |  |  |  | x | x | 4 |
| human autonomy | x |  |  |  |  |  |  |  | x | x |  |  |  |  |  | 4 |
| protection of whistleblowers |  |  |  |  |  |  | x |  |  |  |  |  |  |  |  | 1 |
| hidden costs (labeling, clickwork, material resources etc.) |  |  |  |  |  |  |  | x |  |  |  |  |  |  |  | 1 |
| affiliation (government, industry, science) | government | science | science | science | science | science | science | science | science | industry | industry | industry | industry | industry | industry | |
| number of ethical aspects | 8 | 7 | 11 | 11 | 11 | 9 | 11 | 5 | 11 | 10 | 8 | 6 | 6 | 5 | 8 | |

# „Unethical" behaviour (Kish-Gephart et al. 2010)

individual characteristics

moral development
idealism
machiavellianism
job satisfaction
education level

characteristics of the moral issue

concentration of consequences
magnitude of consequences
probability of effect
proximity

environment characteristics

egoistic ethical climate
benevolent ethical climate
code existence
code enforcement
appreciation of ethical behaviour

„unethical" behaviour / intentions in organisations

# Dr. Thilo Hagendorff

University of Tuebingen
Cluster of Excellence Machine Learning

thilo.hagendorff@uni-tuebingen.de
www.thilo-hagendorff.info