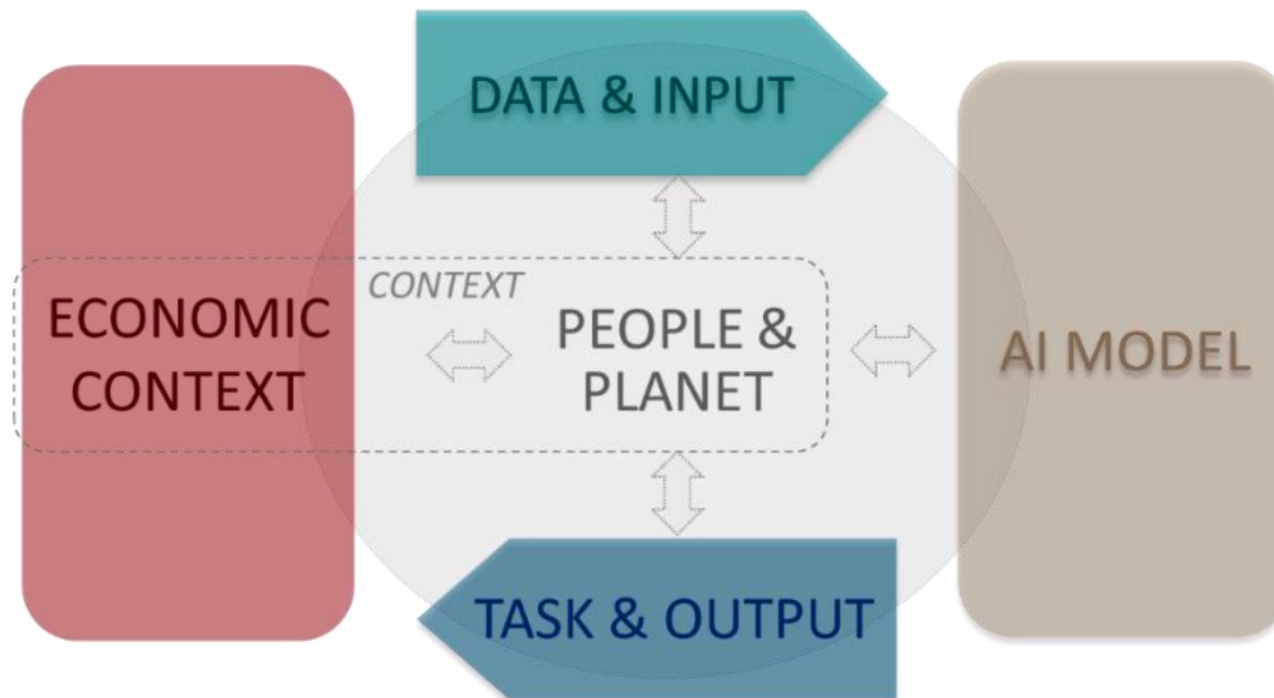# Responsible Data Science

Session 6: 25.05.2023, 15.15 – 19.30 h
MA Seminar, SoSe 2023, Hasso-Plattner Institut

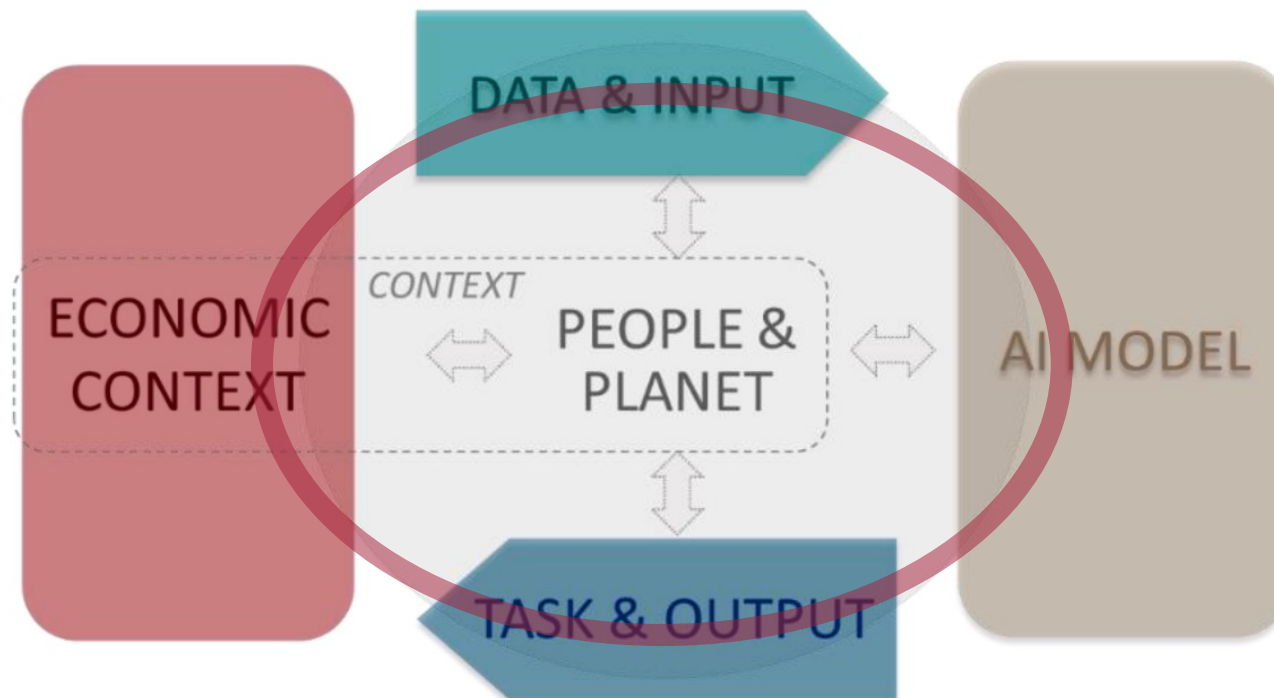Dr. Simon David Hirsbrunner

# Today

| topic | time |
|---|---|
| Introduction | 15h15 |
| Student presentation and discussion: Anatomy of AI (Pia Rissom) | 15h30 |
| Discussion: Sustainable AI | 16h30 |
| —— Break —— | 17h00 |
| Input: risk assessments / reflexivity | 17h30 |
| Exercise: simulation | 17h45 |
| —— Break —— | 18h30 |
| Seminar papers: brainstorming session in small groups | 18h45 |
| End | 19h30 |

# People and planet



Source: key high-level dimensions of the OECD Framework for the classification of AI Systems 2022

# People and planet



Source: key high-level dimensions of the OECD Framework for the classification of AI Systems 2022

# Presentation and discussion

Crawford, K., & Joler, V. (2019). Anatomy of an AI System. *Virtual Creativity*, *9*(1), 117–120.

(presentation by Pia Rissom)

HPI Hasso Plattner Institut
Digital Engineering · Universität Potsdam

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# AI as wicked problem?

As described by Rittel and Webber (1974), wicked problems have 10 important characteristics:

1) They do not have a definitive formulation.

2) They do not have a "stopping rule." In other words, these problems lack an inherent logic that signals when they are solved.

3) Their solutions are not true or false, only good or bad.

4) There is no way to test the solution to a wicked problem.

5) They cannot be studied through trial and error. Their solutions are irreversible so, as Rittel and Webber put it, "every trial counts."

6) There is no end to the number of solutions or approaches to a wicked problem.

7) All wicked problems are essentially unique.

8) Wicked problems can always be described as the symptom of other problems.

9) The way a wicked problem is described determines its possible solutions.

10) Planners, that is those who present solutions to these problems, have no right to be wrong. Unlike mathematicians, "planners are liable for the consequences of the solutions they generate; the effects can matter a great deal to the people who are touched by those actions."

# Presentation and discussion

van Wynsberghe, A. (2021). Sustainable AI: AI for sustainability and the sustainability of AI. *AI and Ethics, 1*(3), 213–218. https://doi.org/10.1007/s43681-021-00043-6

HPI Hasso Plattner Institut
Digital Engineering · Universität Potsdam

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
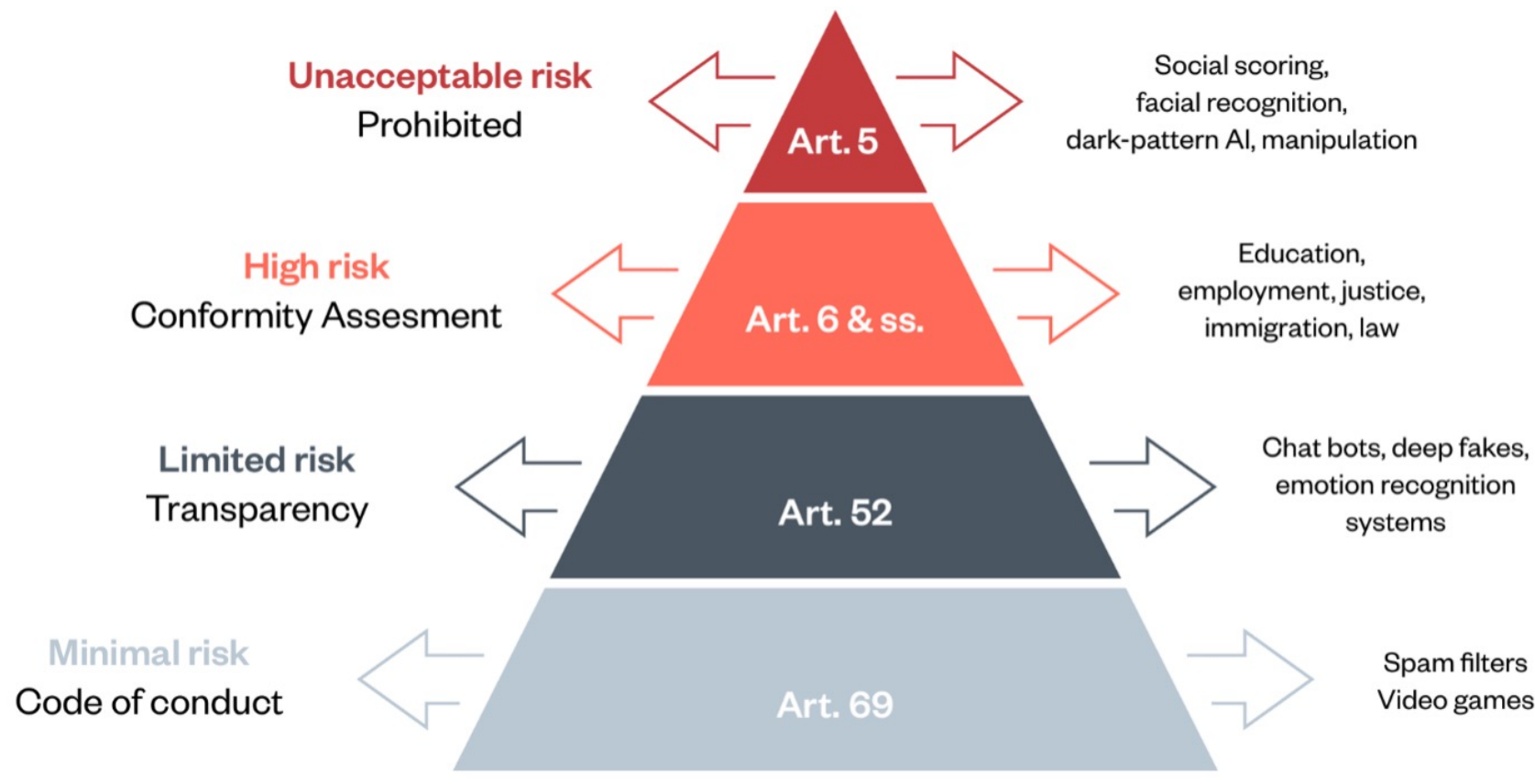TÜBINGEN

# AI risks for stakeholders

**Definition of risk**
„general probability of negative consequences to actions"
(see Cambridge Dictionary, 2022 in Lütge et al. 2022).

**Algorithm-related risks for stakeholders**
❏ Access to goods, benefits, or services ❏ Financial ❏
Property/material resources ❏ Reputation ❏ Emotional ❏
Life/security ❏ Privacy ❏ Liberty ❏ Rights/intellectual property
(see City and County of San Francisco's Ethics and Algorithms
Toolkit).

HPI Hasso Plattner Institut
Digital Engineering · Universität Potsdam

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# Risk-based approach of the (proposed) EU AI regulation



**Unacceptable risk** — Prohibited — Art. 5 — Social scoring, facial recognition, dark-pattern AI, manipulation

**High risk** — Conformity Assesment — Art. 6 & ss. — Education, employment, justice, immigration, law

**Limited risk** — Transparency — Art. 52 — Chat bots, deep fakes, emotion recognition systems

**Minimal risk** — Code of conduct — Art. 69 — Spam filters, Video games

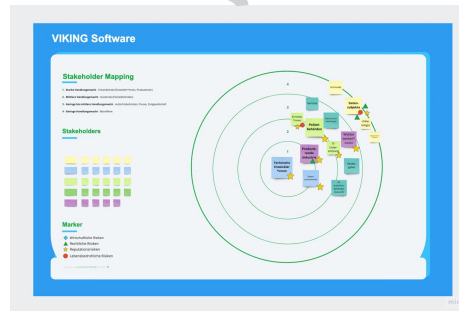Source: https://www.adalovelaceinstitute.org/resource/eu-ai-act-explainer/

## Tool
# Stakeholder and risk mapping



### Grad der Handlungsmacht

1 - **Starke Handlungsmacht** - Entwickelnde (Entwickler*innen, Produzenten)

2 - **Mittlere Handlungsmacht** - Nutzende (Polizeibehörden)

3 - **Geringe bis mittlere Handlungsmacht** - Aufsichtsbehörden, Presse, Zivilgesellschaft
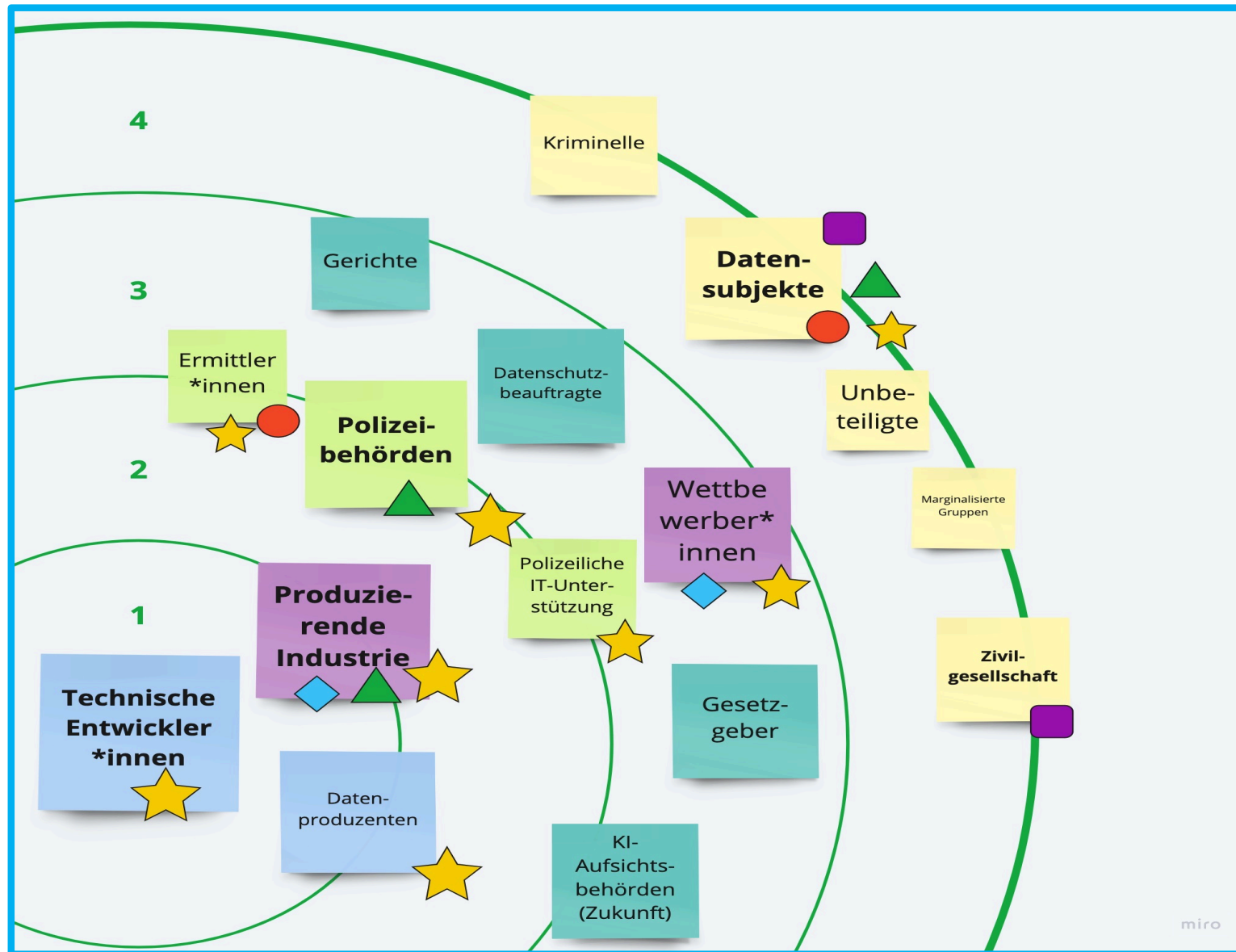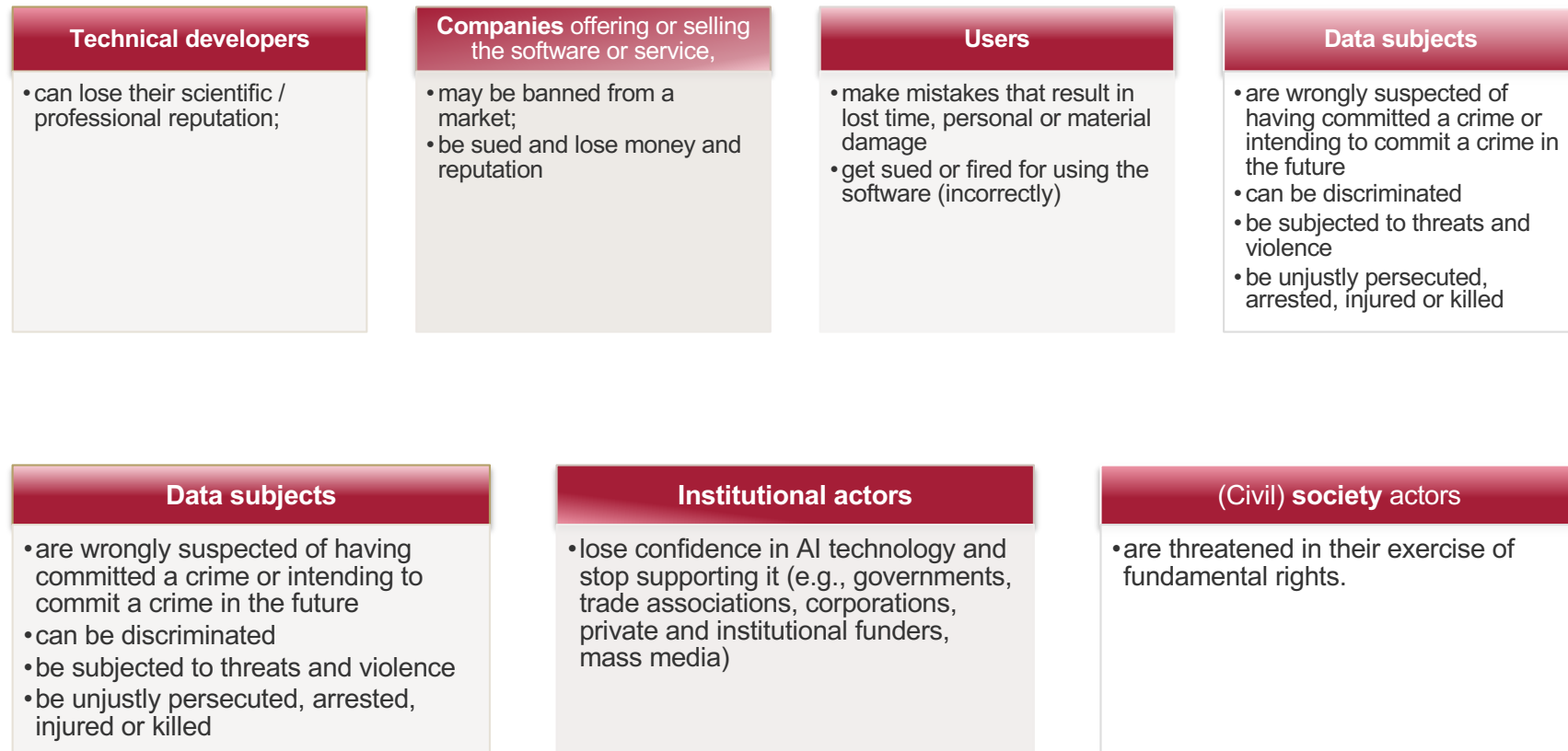
4 - **Geringe Handlungsmacht** - Betroffene

### Risiken

◆ Wirtschaftliche Risiken
▲ Juristische Risiken
★ Reputationsrisiken
● Lebensbedrohliche Risiken
■ Grundrechtliche/ethische Risiken

# Risks for different stakeholder groups

| Technical developers | Companies offering or selling the software or service, | Users | Data subjects |
|---|---|---|---|
| • can lose their scientific / professional reputation; | • may be banned from a market;<br>• be sued and lose money and reputation | • make mistakes that result in lost time, personal or material damage<br>• get sued or fired for using the software (incorrectly) | • are wrongly suspected of having committed a crime or intending to commit a crime in the future<br>• can be discriminated<br>• be subjected to threats and violence<br>• be unjustly persecuted, arrested, injured or killed |

| Data subjects | Institutional actors | (Civil) society actors |
|---|---|---|
| • are wrongly suspected of having committed a crime or intending to commit a crime in the future<br>• can be discriminated<br>• be subjected to threats and violence<br>• be unjustly persecuted, arrested, injured or killed | • lose confidence in AI technology and stop supporting it (e.g., governments, trade associations, corporations, private and institutional funders, mass media) | • are threatened in their exercise of fundamental rights. |

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# **1** **Exercice**

Come back to your stakeholder mapping addressing Palantir's AIP (Artificial Intelligence Platform) for defense.

Identify concrete risks and attach them to different stakeholders on the map (post-its).

# Reflexivity in data science

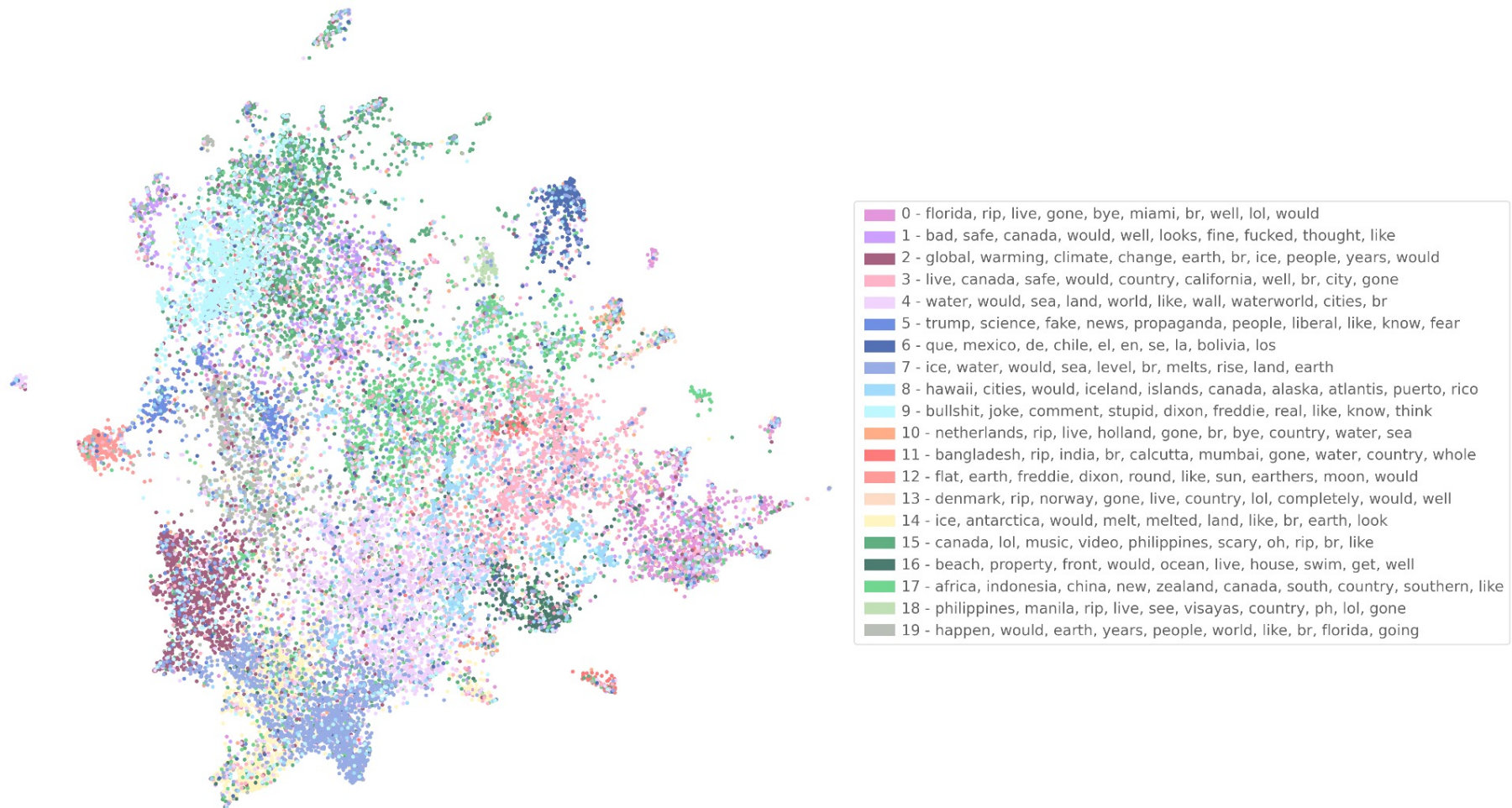Reflexivity understood as critical self-awareness should be a cornerstone of a *responsible data science.*

"Reflexive research cultivates a critical self-awareness, including itself among its objects of study and developing useful concepts for reflecting on the research as it is happening." (Agre 1997a: 27)

# An experiment in reflexivity (I)



0 - florida, rip, live, gone, bye, miami, br, well, lol, would
1 - bad, safe, canada, would, well, looks, fine, fucked, thought, like
2 - global, warming, climate, change, earth, br, ice, people, years, would
3 - live, canada, safe, would, country, california, well, br, city, gone
4 - water, would, sea, land, world, like, wall, waterworld, cities, br
5 - trump, science, fake, news, propaganda, people, liberal, like, know, fear
6 - que, mexico, de, chile, el, en, se, la, bolivia, los
7 - ice, water, would, sea, level, br, melts, rise, land, earth
8 - hawaii, cities, would, iceland, islands, canada, alaska, atlantis, puerto, rico
9 - bullshit, joke, comment, stupid, dixon, freddie, real, like, know, think
10 - netherlands, rip, live, holland, gone, br, bye, country, water, sea
11 - bangladesh, rip, india, br, calcutta, mumbai, gone, water, country, whole
12 - flat, earth, freddie, dixon, round, like, sun, earthers, moon, would
13 - denmark, rip, norway, gone, live, country, lol, completely, would, well
14 - ice, antarctica, would, melt, melted, land, like, br, earth, look
15 - canada, lol, music, video, philippines, scary, oh, rip, br, like
16 - beach, property, front, would, ocean, live, house, swim, get, well
17 - africa, indonesia, china, new, zealand, canada, south, country, southern, like
18 - philippines, manila, rip, live, see, visayas, country, ph, lol, gone
19 - happen, would, earth, years, people, world, like, br, florida, going

(Hirsbrunner et al. 2022)

HPI Hasso Plattner Institut
Digital Engineering · Universität Potsdam

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# An experiment in reflexivity (II)



(Hirsbrunner et al. 2022)

# An experiment in reflexivity (III)



(Hirsbrunner et al. 2022)

# An experiment in reflexivity (IV)

```python
# Two lists of Strings (pieces of text) are created
text_list_1 = ['Global warming',
               'Trump is bad.'
               ]


text_list_2 = ['Climate change',
               'Trump is a bad ass.'
               ]
```
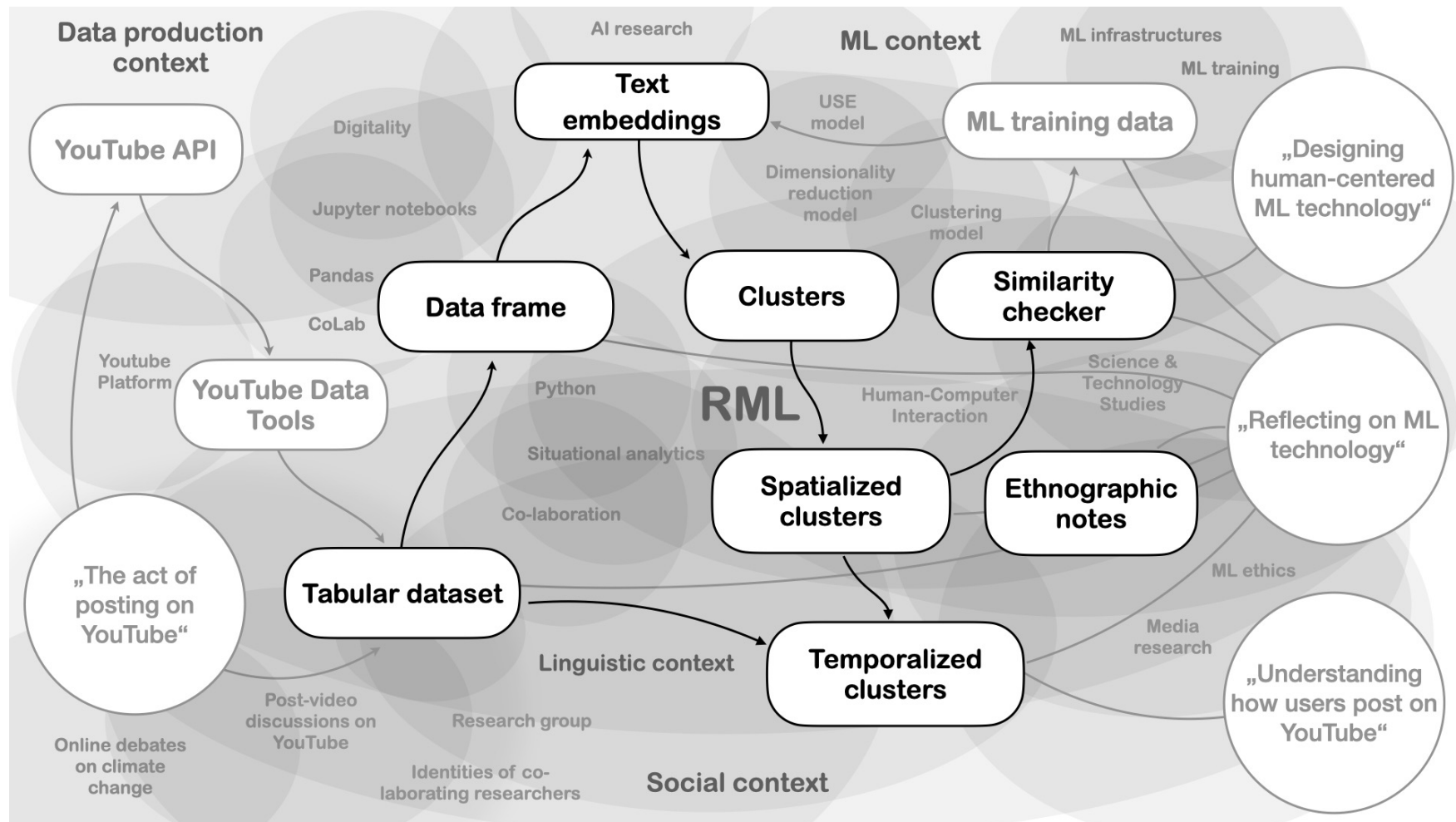
|   | Texts 1 | Texts 2 | Scores |
|---|---------|---------|--------|
| 0 | Global warming | Climate change | 0.646544 |
| 1 | Trump is bad. | Trump is a bad ass. | 0.908402 |

(Hirsbrunner et al. 2022)

# An experiment in reflexivity (III)



(Hirsbrunner et al. 2022)

INTERNATIONALES ZENTRUM
FÜR ETHIK IN DEN
WISSENSCHAFTEN (IZEW)

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# 2 Exercice

Take up the role of one of the stakeholders. Make some notes about the interests of this stakeholder.

Defend the interests of the stakeholder within the staged design process.

## Sources

See entire list of course references on Github:
 https://github.com/simonsimson/responsible-data-science/tree/main/slides