# ETHICAL ASPECTS IN RESEARCH TOPICS

**Mihály Héder**
*Department of Philosophy and History of Science*
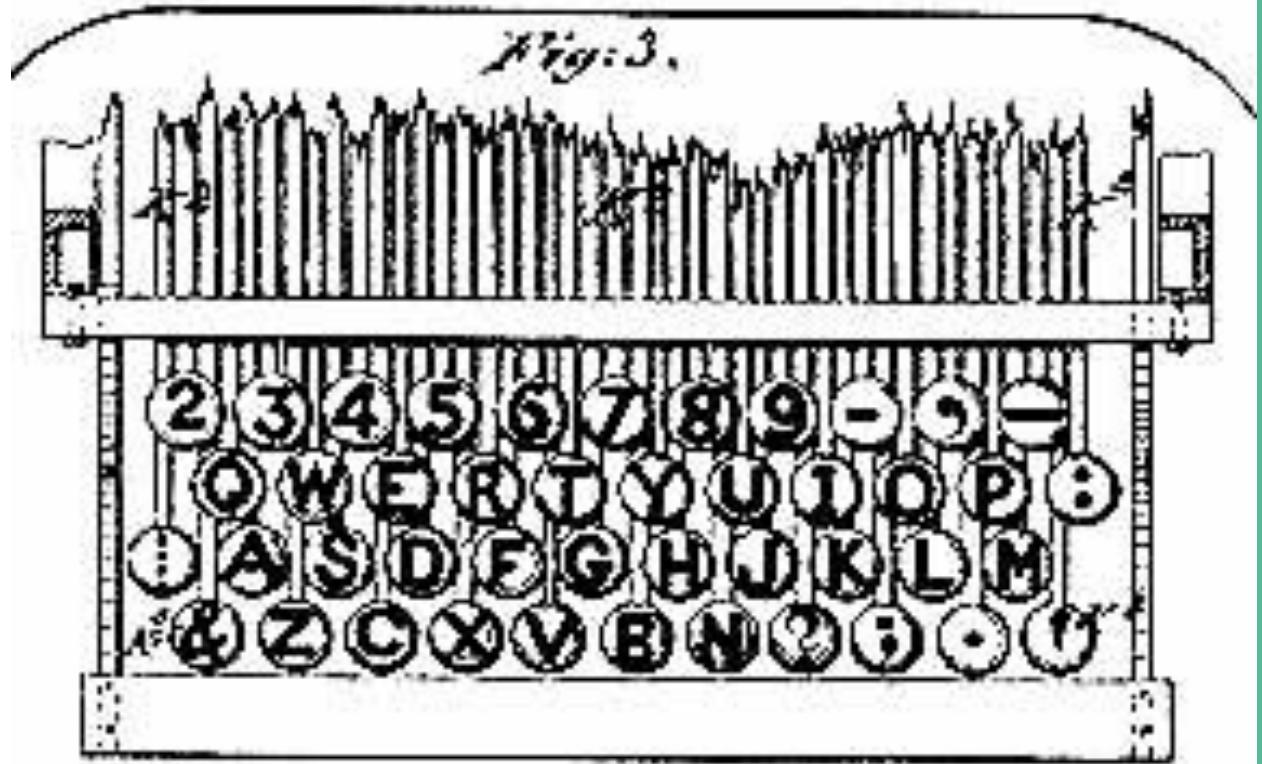*www.filozofia.bme.hu*
*Budapest University of Technology and Economics*

**HU Utrecht, 27 January 2025.**

humancentered-ai.eu

# WHY are we here?
# (on a human-centered AI event)?

# The history of the QWERTY keyboard layout



Fig: 3.

# The first typewriter

- The first typewriter was created by Christopher Latham Sholes in the 1860s
- There were a lot of technological problems:
  - The typewriter easily jammed, especially if the user was fast enough
- Sholes started with the alphabetic order first
  - He re-arranged the keyboard to a random layout (QWERTY) in order to **slow down the user**

# QWERTY's diffusion

- In the **1880s** there were concurrent keyboard layouts, QWERTY was one of them
- By 1890 there was no technological reason to use QWERTY
  - But it was bought by Remington and spread everywhere called "universal"

# Locked in to QWERTY

- By the 20th century society was locked in to QWERTY
- There are three reasons of this
  - Economics of scale (mass manufacturing)
  - Interdependence of technology
  - Irredeemable investments

# Collingridge-dilemma

- There lies a deep tension in the logic of technology development:
  - When a **technology is new** it is easy to modify
    - In theory, big problems could be avoided
    - BUT there is not enough information to discover the exact problems
  - In case of an **established, ubiquitous technology** it is easy to see the problems
    - But it is hard to make changes



**example: Dichlorodiphenyltrichloroethane (DDT)**

# AI vs. technological lock-in

# Why do we need to deal with AI specifically?

- **Potential technology lock-in:**
  - The impact is greater: fewer and fewer **technology owners** can influence a larger and larger slice of the world
  - AI is **software**
    - consequently: the marginal cost of its "multiplication" is negligible.
    - e.g. if there will be an effective AI diagnostic system / psychologist / accountant it will be more worthwhile to copy it than to rewrite it

# Why do we need to deal with AI specifically?

- **Potential technology lock-in:**
  - Moreover, AI is Software-as-a-Service, so it can be provided from anywhere in the world in most cases
    - except locally controlled robots
  - AI developers start in a certain direction, which will be difficult to change later
    - similar situation: e.g. gmail accounts - raise your hand if you don't have one

# Solution attempt #1
# Human-Centered Design of machines?

hcaim
human centred
artificial intelligence
masters

# What is not human-centered design?



Source: Elias Beck. *'Child Labor in the Industrial Revolution'*. History Crunch. December 30, 2021. https://www.historycrunch.com/child-labor-in-the-industrial-revolution.html#/

# What is not human-centered design?

Gilbreth chronocyclograph of motions necessary
to move and file sixteen boxes full of glass, n.d.
From: Mike Mandel, *Making Good Time: Scientific
Management, the Gilbreths, Photography and
Motion, Futurism* (Santa Cruz, CA: California
Museum of Photography, University of California,
Riverside, 1989), 26.

# What is not human-centered design?



TIME

☰     SUBSCRIBE

BUSINESS • TECHNOLOGY

Exclusive: OpenAI Used Kenyan Workers on Less Than $2 Per Hour to Make ChatGPT Less Toxic

15 MINUTE READ

This image was generated by OpenAI's image-generation software, Dall-E 2. The prompt was: "A seemingly endless view of African workers at desks in front of computer screens in a printmaking style." TIME does not typically use AI-generated art to illustrate its stories, but chose to in this instance in order to draw attention to the power of OpenAI's technology and shed light on the labor that makes it possible.   Image generated by Dall-E 2/OpenAI

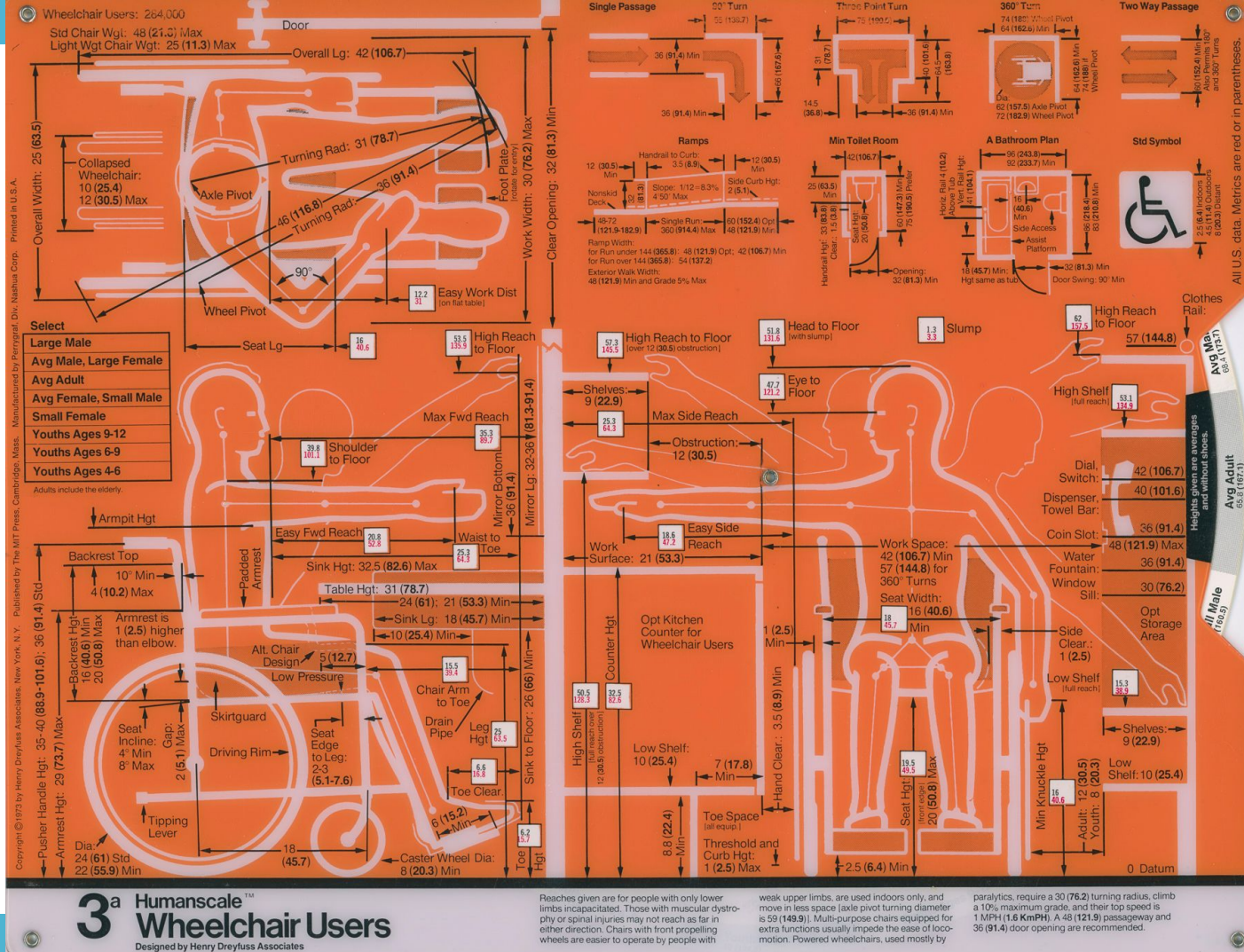BY **BILLY PERRIGO**

JANUARY 18, 2023 7:00 AM EST

hcaim — human centred artificial intelligence masters

# What is the machine maker's role on all of this?

- Self-absolving strategy 1
  - "The imperative of technology"
    - This is what **efficiency dictates**
- Self-absolving strategy 2
  - This is **what the customer wants**
- Self-absolving strategy 3
  - It **was** **legal** at the time

- Human-Centered Design: **rejection of all the above**
  - **humanities toolkit**
    - **argumentation**
    - **critical thinking**

# Knowledge about Humans

Image: Henry Dreyfuss Associates, *Humanscale* selector 3a "Wheelchair Users," 1974. Plastic, paper, and metal. Milwaukee Art Museum Research Center.

Source: Hanna Pivo, 20th-Century Tools for Measuring Time and Bodies April 19, 2019
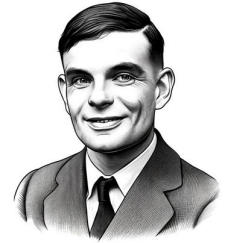https://blog.mam.org/2019/04/19/20th-century-tools-for-measuring-time-and-bodies/
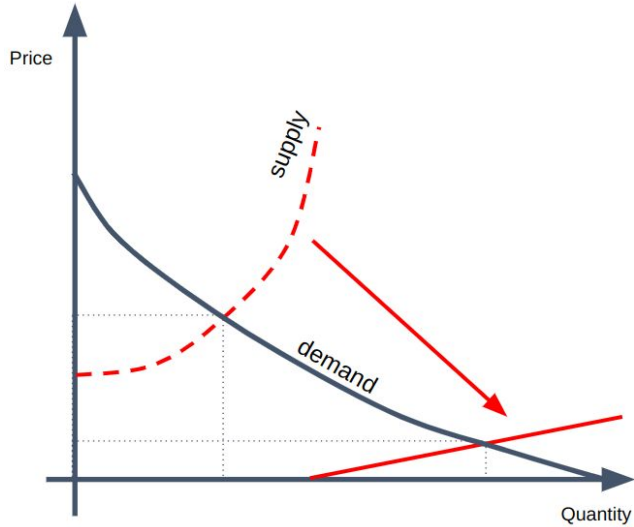
# Ethical topics in AI

# AI Ethics Problem Set:

## THE FUTURE OF WORK



Image credit: Freepik premium

# Why is AI special?
# Future of Work

- AI threatens jobs
- This could result in another economic transformation similar to the **industrial revolutions**
  - from the point of view of today's people, the industrial revolutions represent the historical antecedent of a comfortable, technological life,
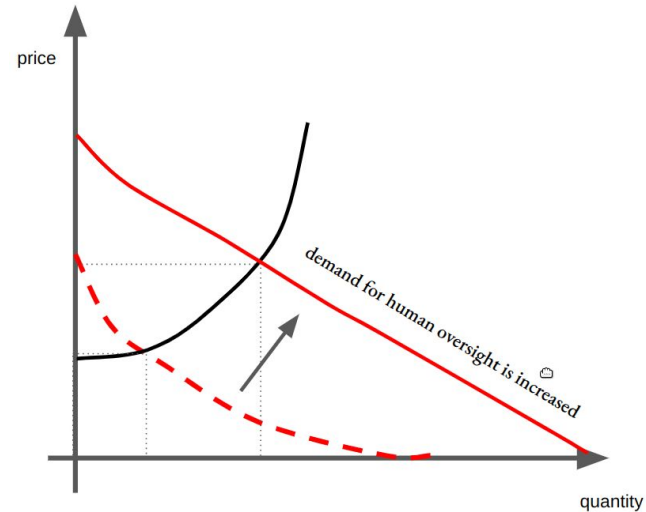  - but those who experienced it could also be accompanied by serious disruption and impoverishment

# Problem Set: The Artificial Worker

# Machine Ethics

# AI Ethics Problem Set:
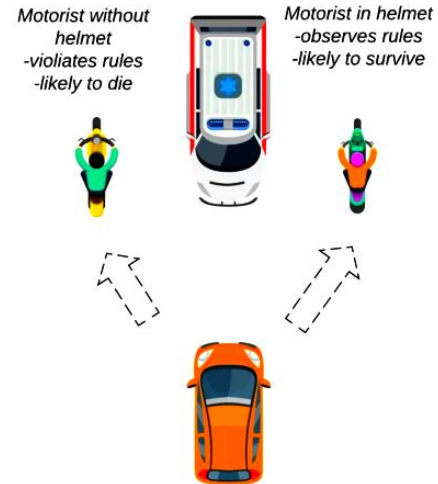
## MACHINE ETHICS



Motorist without helmet
-violiates rules
-likely to die

Motorist in helmet
-observes rules
-likely to survive

P. Lin helmet problem

# Why is AI special? Machine Ethics

- Designing the behavior of an artificial person (machine ethics)
  - we give the machine autonomy, since we precisely want it to think and make decisions for us
    - analogy: *raising children*
    - we want to avoid: biases, opacity
  - We place a machine in a "trusted" position in an unprecedented way
  - It arises that that artificial agent develops its own goals and does not only serve our goals

# Topics of machine ethics

- Transparency
  - explainability
- Fairness
  - bias
  - lack of discrimination
- Alignment
  - wider social **values** like wellbeing
    - brings in politics

# Fairness metrics

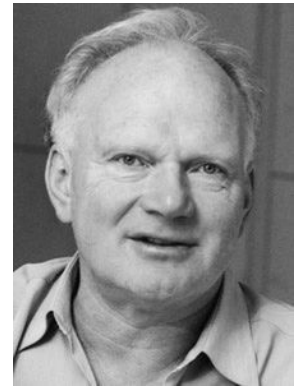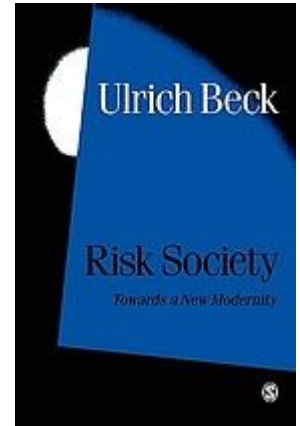| | | Prediction | | |
|---|---|---|---|---|
| | Total Population | Predicted Positive | Predicted Negative | $Prevalence = \dfrac{\sum TP}{\sum Total\ Population}$ |
| Ground Truth | Ground Truth Positive (GTP) | True Positive (TP) | False Negative (FN) | True Positive Rate, Sensitivity, Recall $TPR = \dfrac{\sum True\ Positive}{\sum GTP}$ |
| | Ground Truth Negative (GTN) | False Positive (FP) | True Negative (TN) | False Positive Rate, Fallout $FPR = \dfrac{\sum False\ Positive}{\sum GTN}$ |

- Equal Opportunity
  - $TPR_0 \sim= TPR_1$
- Equalized Odds
  - $TPR_0 \sim= TPR_1$
  - $FPR_0 \sim= FPR_1$

Demand from society
The social control of technology?

humancentered-ai.eu

# Ulrich Beck's analysis (1980's)

- A main feature of modern society is that it is **preoccupied with the future**, and
  - especially the negative scenarios, that is 'risks'
- Catastrophes were formerly attributed to bad luck or divine acts
  - but not in Humanity's control
- Now that our control seems greater (modern science) the **responsibility is ours**
  - this in turn undermines the institutions of modern society, e.g. trust in science

1944-2015

# Modernity 2 (Beck's society)

- The victories of the first modernity (taking risks) have a boomerang effect
- Taking risks not serve us anymore
- So we enter **reflexive modernity**
  - Pesticide
  - Ozone
  - Nuclear
  - Toxins
  - CFC
  - Plastics, etc.
- Plus, we anticipate even more negative consequences
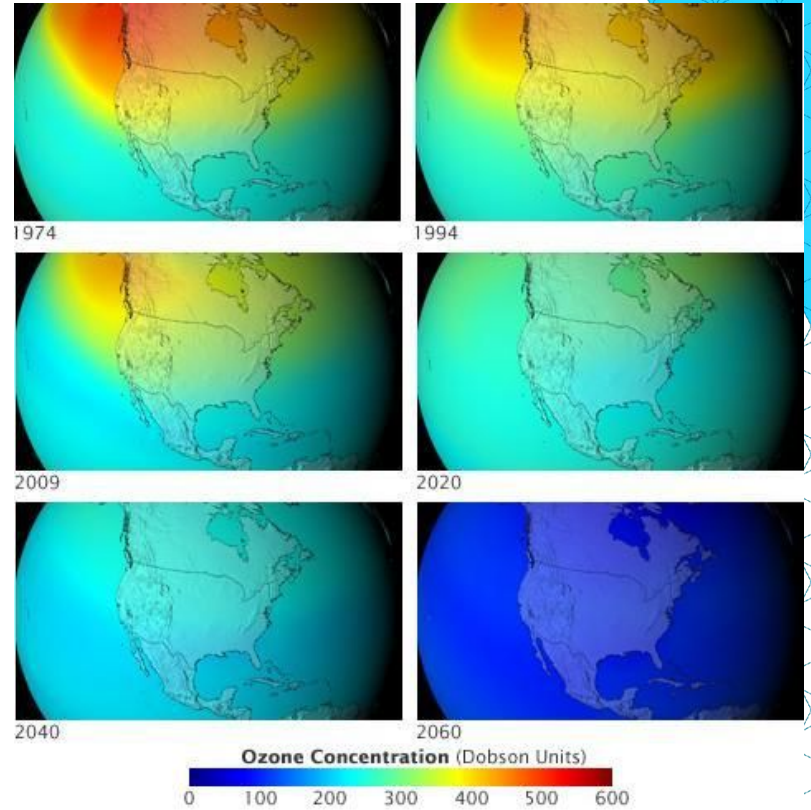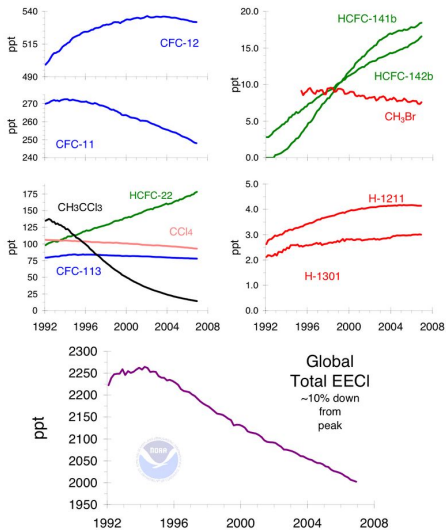  - AI, GMO

# Modernity 2

- No clear culprits
  - many of us are implicated in these negative effects
  - we need modern technology to even identify and tackle risks
- The new risks are far more evenly distributed
  - Modernity 1
    - Living next to a factory was risky, if rich you could move away
  - Modernity 2
    - Ozone, global warming
    - arguably you can only temporarily can avoid these risks with your personal wealth
- **You will operate in a more regulated environment than any generation before**
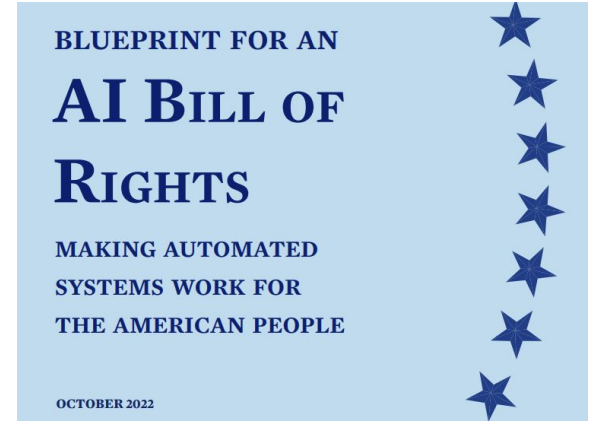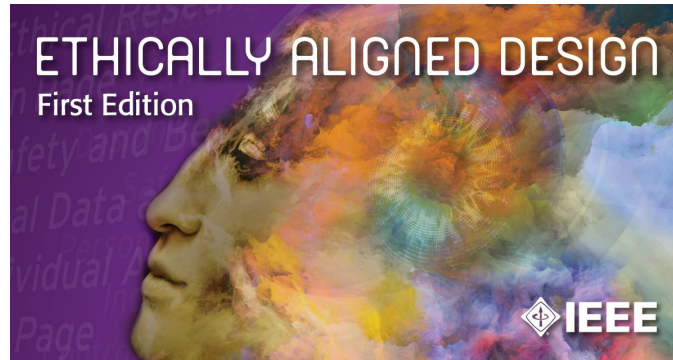
Solution attempt #2
Regulation

humancentered-ai.eu

# Reasons to regulate

- Avoid harm
- Pushback
  - "regulation stifles innovation"
    - answer1: there are worse things than slowed innovation)
    - answer2: Porter's hypothesis





Ozone Concentration (Dobson Units)

**Chlorofluorocarbon (CFC)**
**Wikimedia Commons: Ozone Layer depletion**
**if the CFC was not banned**

# Soft Law / Recommendations



The European Commission's
**High-Level Expert Group on Artificial Intelligence**

AI

DRAFT
**Ethics Guidelines for Trustworthy AI**



**ETHICALLY ALIGNED DESIGN**
First Edition

IEEE



BLUEPRINT FOR AN
**AI BILL OF RIGHTS**

MAKING AUTOMATED
SYSTEMS WORK FOR
THE AMERICAN PEOPLE

OCTOBER 2022

Recommendations, Guidelines and Law

# AI Act

Official Journal
of the European Union

EN

L series

2024/1689

12.7.2024

**REGULATION (EU) 2024/1689 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL**

of 13 June 2024

laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act)

https://artificialintelligenceact.eu/assessment/eu-ai-act-compliance-checker/

# Why are we here - revisited

hcaim

human centred
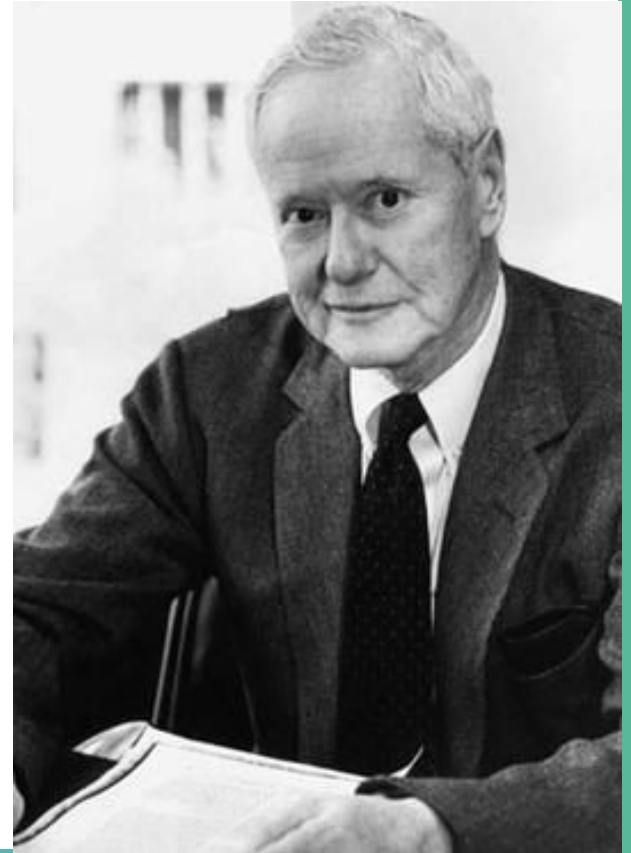artificial intelligence
masters

# Why are we here: to learn

- To understand the stakes
  - **technological lock-in**
- To learn about the methods
  - **machine ethics**, guidelines, technical standards
- To understand the demands of society
  - **rules** & regulations

# Why are we here: to think & challenge

- Personal Integrity
  - Challenging the
    - **customer**,
    - the **technology** and
    - the **rules**
    - **yourself**
  - Ability to walk away
- It takes knowledge of
  - history
  - argumentation
  - ethics
- CUDOS - the ethos of science
  - **C**ommon discoveries
  - **U**niversal knowledge
  - **D**isinterestedness
  - **O**rganized **S**cepticism

Robert K. Merton
source:
WikiMedia
Commons

# Why are we here: to meet

- everyone needs acknowledgement
- in human-centered AI your are most likely to get this from your peers

retroreflective vests

bicycle retroreflectors