# Ethical AI Frameworks and their Limits

Erich Prem

**eu|te|ma**
TECHNOLOGY MANAGEMENT

UNIVERSITY OF VIENNA

# Overview

A short introduction to ethics

Modelling people

Ethical issues of AI

Frameworks and principles

From principles to practice

Final remarks

# Exercises

a. ChatGPT (hallucination)
b. Bias (creditworthiness)
c. Decision-making (trolley problem)
d. Illegal texts (generative AI)
e. Illegal images (image classification)
f. Responsibility (model cards)

# A short introduction to ethics

*Compassion is the basis of morality.*

A. Schopenhauer

# Philosophy of morality

*Morality is an informal public system applying to all rational persons, governing behaviour that affects others, and includes what are commonly known as the moral rules, ideals and virtues and has the lessening of evil and harm as its goal.*
(Bernard Gert)

εθος – custom (behaviour)

ηθος – character (attitude towards behaviours)

descriptive, normative, applied, metaethics

**Some common virtues**
  truthfulness
  courage
  honesty
  impartiality
  reliability
  …
**Ideals:** e.g., justice
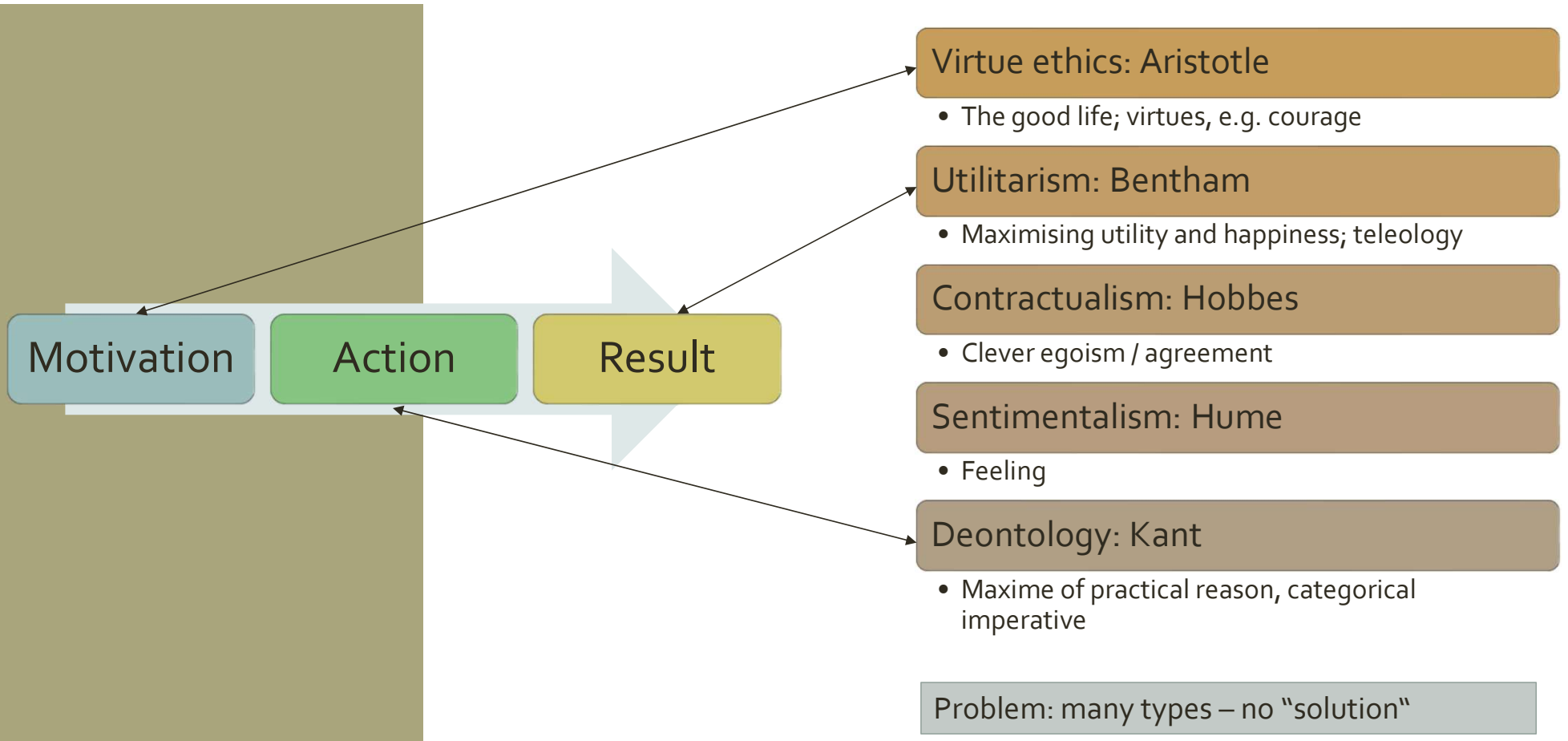
**Some common harms**
  death
  pain
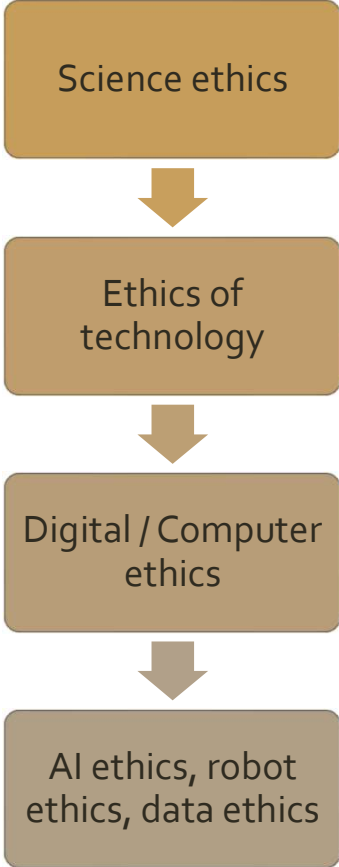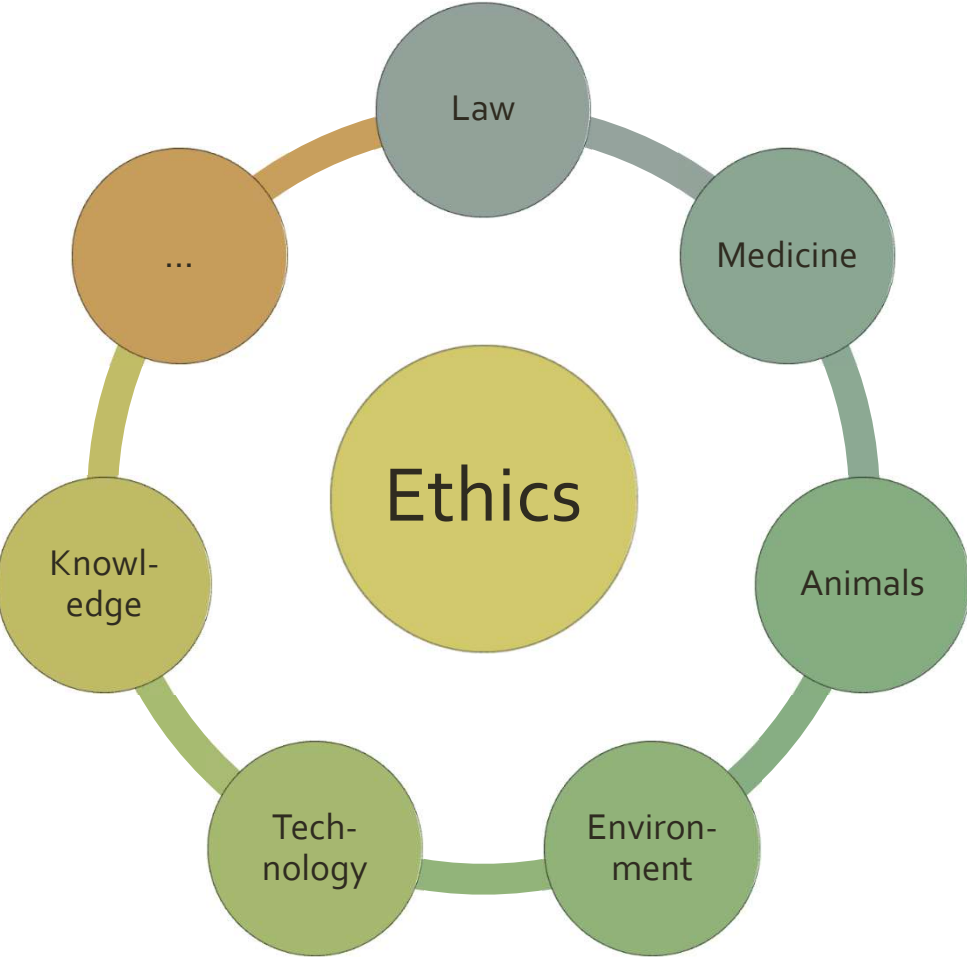  disability
  loss of freedom
  loss of pleasure
  loss of rights
  …

# Types of ethics

Motivation → Action → Result

**Virtue ethics: Aristotle**
- The good life; virtues, e.g. courage

**Utilitarism: Bentham**
- Maximising utility and happiness; teleology

**Contractualism: Hobbes**
- Clever egoism / agreement

**Sentimentalism: Hume**
- Feeling

**Deontology: Kant**
- Maxime of practical reason, categorical imperative

Problem: many types – no "solution"

# Domain ethics and digital ethics

# Computer ethics

A separate field of ethics?

- Ubiquity of computer technology
- *Open, malleable* technology (J.H. Moor) vs. *old* ethical problems in new clothes (Deborah Johnson)?
- New aspects (e.g. internet, robotics, AI, data science)

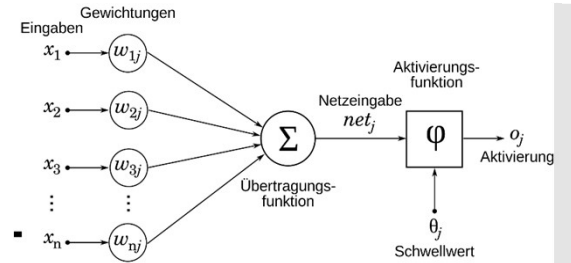| | | | |
|---|---|---|---|
| Guidelines for programmers and IT-specialists | Safety | Cybercrime | IPR and ownership of software |
| **Privacy and anonymity** | Responsibility | Networks, virtual societies, globalisation | Technical dependability |
| Distribution fairness | Power, democracy, participation | Computers and education | Automation, labour, and work |
| Accessibility | Robot ethics | **AI, algorithmic decision making** | **Autonomous systems** |

# Modelling people

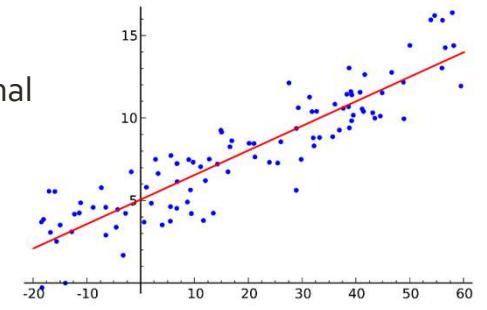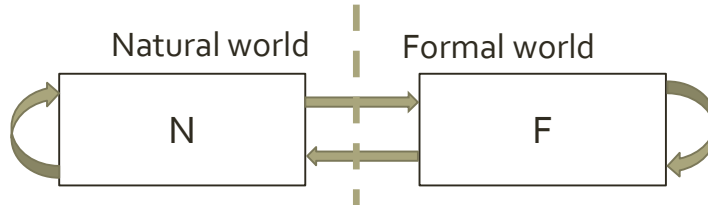*Ethics is in origin the art of recommending to others the sacrifices required for cooperation with oneself.*

B. Russel

# The modelling problem



Model
- Correct
- Relevant
- Simple

Natural world

Formal world

$N$

$F$

Natural law

Formal rules

Eingaben, Gewichtungen, Netzeingabe $net_j$, Aktivierungsfunktion, Übertragungsfunktion, Schwellwert $\theta_j$, Aktivierung $o_j$

BACKWARD CHAINING

GOAL: Make $20.00

RULE: If the lawn is shaggy and the car is dirty and you mow the lawn and wash the car, then Dad will give you $20.00

Does the lawn need mowing?   Does the car need washing?

Do you have a mower?   hose?   bucket?   rags?

gas?   electric?   push?

*** The inference engine will test each rule or ask the user for additional information.

# What changes if N=human, modelling people?

Complexity limits our models in what we can know, predict, or control – and in some cases what we *should* do.

Natural world | Formal world

Natural law

N

F

Formal rules

Should we **know** a person's
- Gender, income, religion, sexuality
- Online searches
- Pharmaceutical shopping?

Should we **predict** a person's
- Talent
- Time of death
- Likelihood of getting STDs
- Unemployment?

Should we **control** a person's
- Exercise routine?
- Learning capacity?
- Eating habits?
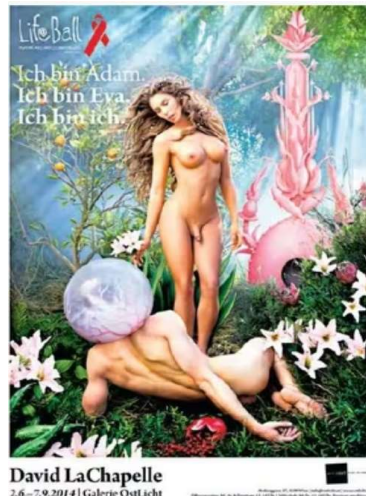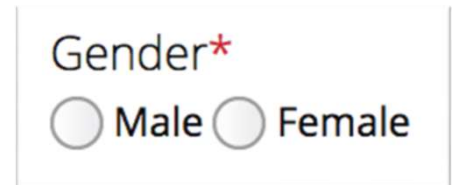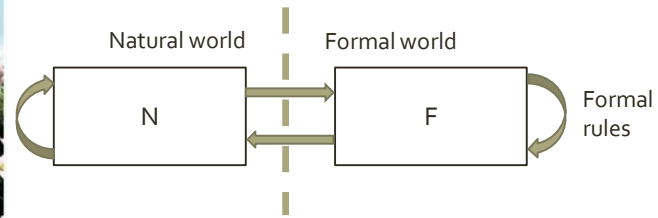
# Observables of complex systems are a choice.



Photo: Life Ball, David LaChapelle

Natural world    Formal world

N    F    Formal rules

Gender*

○ Male ○ Female

Our choices have epistemic and ethical consequences:
What gets counted counts.

# How data analytics and AI may impact on people's lives

## Decide on people

- Denying loan
- Losing a job
- Out secrets to family
- Increase insurance premium
- Objectify individuals as a mere category

## Perpetuate trends

- Continue the past by taking decisions based on the past (e.g., admission)

## Influence people

- Trigger behaviours such as voting, buying, spending, …

AI may help to make existing ethical issues explicit. Addressing bias goes beyond just a 'correction' – it can be a response of society to change the future.

# Classification of human behaviour



Support tools for toilet for people with disabilities or dementia

Use of **depth sensors** instead of camera

TU Wien Institute of Visual Computing Computer Vision Lab

## How much should we know?

Violations of privacy may cause

- Degradation (dignity)

- Potential to exploit

Not merely a legal issue, also an ethical concern about autonomy.

TU Wien Institute of Visual Computing
Computer Vision Lab

TU WIEN
TECHNISCHE
UNIVERSITÄT
WIEN
Vienna University of Technology

# Should companies...

**How?**

**Benefit?**

**Business?**

**Value**

- *Build models* of employees based on their medical records and digital traces to predict their level of absence from the firm or to offer gym classes?
- Should we *monitor* what people watch on television to improve program planning and advertising?
- Should we *predict* a teenagers pregnancy to catch the moment she starts buying new products and is a promising target for special offers?
- Should we *identify* homosexual couples to offer them special offers they might like for vacation?
- Should we equip a car with an electronic black box and *tracker* to offer reduced insurance premiums or *disable* cars to drive Saturday night?

# Ethical issues of AI

*The voice of reason is quiet.*

S. Freud

# Ethical Issues of AI Systems



Safety



Fairness



Transparency



Privacy



Responsibility



Work

How "sensitive" and problematic data is, depends on the context.
Contexts changes over time while data may be persistent
even when it becomes out of date or recognised as wrong.

Data can become very dangerous….

…when the context changes.

**Health**
- treatment from your doctor about the onset of Alzheimer
- data flow to employer

- unemploy-ment

**Dating**
- Grindr or Twitter traces
- a visit to Kuwait or Egypt

- incarceration

**Communication**
- joking online, political critique
- change of politics

- persecution

**Religion**
- minority group
- change in government

- death

# Which discrimination…is fair?

**Insurance**

- People with big houses pay more for their insurance. Fair?

**Add data analysis**

- In some cases, young drivers pay more when renting a car. Fair?

**Add smart controls**

- Young men cause more accidents. Should they pay less for not driving at night or weekends? Use a monitoring device?

**Add AI**

- Propose a safer routing to driver or charge fee if driver decides for dangerous route?

Moral, ethics, or politics?

# The ethics of anticipatory models or the right to a future

## Extending the past into the future



| Past | | Model | | Action | | Future | |
|------|--|-------|--|--------|--|--------|--|
| • Mostly male engineers | | • Engineers are mostly men | | • Choose men as engineers | | • Mostly male engineers | |

Moral, ethics, or politics?

# Trolley Problem and Autonomous Driving



Ethical problems of intervention in human decision-making (already for driver assistance systems)

- Limitation of autonomy (action)
- Creating machine autonomy?
- Business case?

**Selected ethical issues of language models (ChatGPT)**

| Privacy issues and data leaks | Authorship, plagiarism | Work conditions, alienation |
|---|---|---|
| Misinformation | Manipulation, deceit | Censorship |
| Fairness and bias | Security | Power, democracy |

# Exercises A (hallucination), D (illegal texts)

**Content can be illegal or restricted for publication**

- Publicly **denying the holocaust** or distributing Nazi symbols.

- **Inciting terrorist** acts, instructions for illegal actions, recruiting members for terrorist associations (online providers in the EU)

- Publishing a guide to the manufacture of **drugs** may be punishable as aiding the manufacture.

- Participation in a **suicide**

- **Child pornography** (note: pornography is not illegal): depicting sexual acts with children or their genitals (regardless of how they are generated!).

- **Intellectual property** (content, logos, software…)

- **Personal rights** (images)

# Ethical Dilemma: Freedom of expression versus online hate speech

Are **private** platforms **public spaces**?

- Should/must they protect freedom of speech, art, science?

- Regulate content through rules, standards, AI-based filters

**State regulation** as a measure **against** illegal content.

- Question of "disturbing" content such as false reports, strong opinion (so-called "harmful content")?

- Regulation often triggers the use of algorithmic methods for content moderation, usually not prescribed.

- Question of the objective of discourse
  - Counter examples: local media, professional forums (e.g. LinkedIn)
  - Wiener Aktionismus / Viennese Actionism

## AI-based discourse and content moderation

- AI-algorithms for the identification of *problematic* content
  - Many mistakes, simple approaches: difficult technical problem
  - Who has the right to define what should be deleted?
  - What rights should people have whose contributions are deleted?
  - Discourse power: platform collaborate with undemocratic states

- Significant erroneous deletion
  - little information about practice of deletion
  - few options for appeal
  - few pro-freedom regulations (i.e. "rights to publish").

- Alternatives
  - Education
  - De-anonymisation
  - Ombudsperson
  - Effective recourse mechanisms

# Omnia vincit amor

## The externalisation of intention

- Art, pornography or medicine?

- Reducibility of pornography to nudity?

- Question of images and intentions (not depicted).

- cf. debate about chat control in the EU: automatic scanning of communication for child pornography.

Michelangelo Merisi da Caravaggio 1602

https://de.wikipedia.org/wiki/Datei:Caravaggio_-_Cupid_as_Victor_-_Google_Art_Project.jpg

# Frameworks and principles

*Medicine rests upon four pillars – philosophy, astronomy, alchemy, and ethics.*

Paracelsus

# Ethical principles: principlism

**Belmont report** (April 18, 1979)
https://www.hhs.gov/ohrp/regulations-and-policy/belmont-report/read-the-belmont-report/index.html

| Principle | Example application |
|---|---|
| Respect for persons | Informed consent |
| Beneficience | Weighing risks and benefits |
| Justice | Selection of test subjects |

| | |
|---|---|
| autonomy | non-maleficience |
| beneficence | justice |

Tom Beauchamp, James Childress
Orientation at four principles

# Research ethics: origin in medical ethics

Hippocratic oath of ethics, traditionally by physicians – still relevant practice (today Declaration of Geneva)

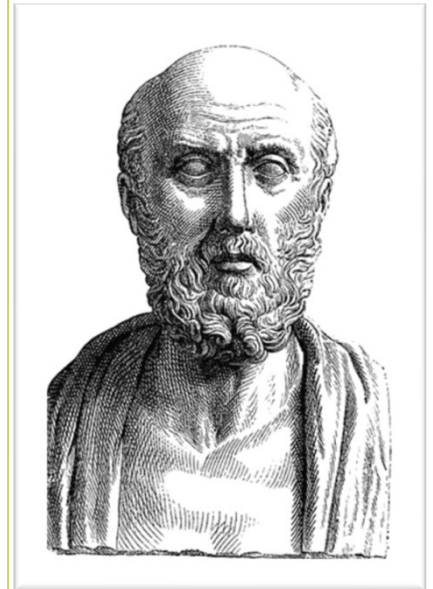Promise to uphold ethical standards specific to medical practice, e.g.,

- Confidentiality

- Non-maleficience ("First do no harm")

AS A MEMBER OF THE MEDICAL PROFESSION:
I SOLEMNLY PLEDGE to dedicate my life to the service of humanity;
THE HEALTH AND WELL-BEING OF MY PATIENT will be my first consideration;
I WILL RESPECT the autonomy and dignity of my patient;
I WILL MAINTAIN the utmost respect for human life;
I WILL NOT PERMIT considerations of age, disease or disability, creed, ethnic origin, gender, nationality, political affiliation, race, sexual orientation, social standing or any other factor to intervene between my duty and my patient;
I WILL RESPECT the secrets that are confided in me, even after the patient has died;
I WILL PRACTISE my profession with conscience and dignity and in accordance with good medical practice;
I WILL FOSTER the honour and noble traditions of the medical profession;
I WILL GIVE to my teachers, colleagues, and students the respect and gratitude that is their due;
I WILL SHARE my medical knowledge for the benefit of the patient and the advancement of healthcare;
I WILL ATTEND TO my own health, well-being, and abilities in order to provide care of the highest standard;
I WILL NOT USE my medical knowledge to violate human rights and civil liberties, even under threat;
I MAKE THESE PROMISES solemnly, freely, and upon my honour.

Unidentified engraver - 1881 Young Persons' Cyclopedia of Persons and Places. Wikipedia.org

# Research ethics: origin in horrific research

- **1892 Albert Neisser injects girls and women with serum from syphilis patients without their consent**
  - **Note: Neisser believed in a bacterial infection**

- **Public discussion:  argument between scientists and opponents including public, law;**

- **Regulation by the Prussian cultural ministry that included a requirement of consent and a ban on experiments with children.**



Doctor drawing blood from a patient as part of the Tuskegee Syphilis Study.

National Archives Atlanta, GA (U.S. government)

**Nazi experiments on jews and other concentration camp inmates**

- **1946/47 Nuremberg Code after WWII**
  - **10 conditions considered essential: voluntary consent, good for society, …**

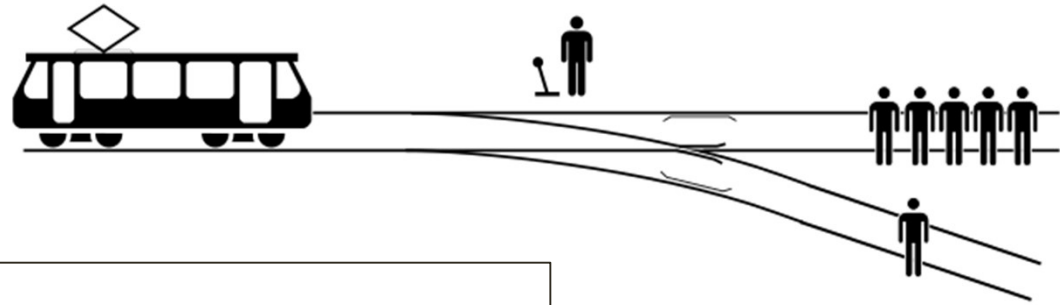- **1964 Declaration of Helsinki**
  - **For all involved actors "to protect the life, health, dignity, integrity, right to self-determination, privacy, and confidentiality of personal information of research subjects (WMA 2013 §11)**

- **1932-1972 Tuskegee Syphilis Study**
  - **US Public Health Service and CDC on 400 African Americans**
  - **participants were not informed of their infection**

# The Trolley Problem: an ethical dilemma
(Engisch 1930)

**Variants**

- Fat man (Thomson 1976)

- Transplantation (Thomson 1985)
  - Healthy donor or patients

- Autonomous vehicles (Lin 2013)
  - Driver or pedestrians

cf. experiments with opinions, e.g. "Moral Machine" online quiz (MIT) with 9 dilemmas or TV shows with viewer participation.

Clarification of different ethical positions: utilitarian versus deontological ethics; positive versus negativ duties (virtue ethics).

Also cultural variation (e.g. saving younger over older).

Not a "solution" of moral problems, e.g. for driving.

Awad, E., Dsouza, S., Kim, R. *et al.* The Moral Machine experiment. *Nature* **563,** 59–64 (2018).
https://doi.org/10.1038/s41586-018-0637-6

# Exercise C

Oversimplified questions can/should be rejected.

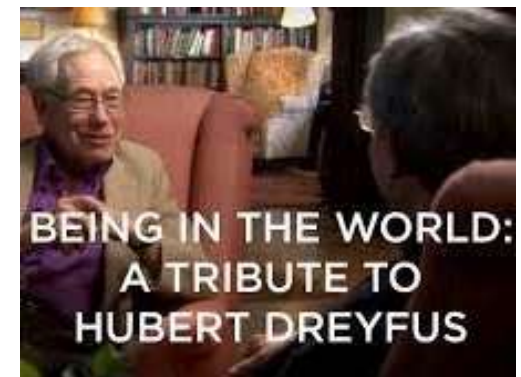- Trolley problem assumes a technical solution. Abstraction from actual issues and strategies of car makers.
  - Very different strategies in practice (reduce energy)
  - Protection of passengers

- No counting up of human lives: irreconcilable with human dignity.

- Kantian categorical imperative: Human as an end-in-itself "Autonomy is hence the reason for the dignity of human and any reasoning nature"

# Role of context



- Humans take the whole context in account
- Surpasses capabilities of today's AI by far
- Question of formalisability of human action

- Hubert L. Dreyfus

- AI critique: human action cannot be described in formal rules

- Requires human purpose: goals, experience, valuation, understanding

- Question of (human) embodiment

- -> anthropocentric ethics



BEING IN THE WORLD: A TRIBUTE TO HUBERT DREYFUS

# Ethical framework principles

- **Transparency** (including explainability, understandability, disclosure etc.)

- **Justice** and fairness (including consistency, inclusion, equality, bias, diversity, remedy, redress etc.)

- **Non-maleficence** (security, safety, precaution, prevention, integrity etc.)

- **Responsibility** (accountability, liability)

- **Privacy**

- **Beneficence** (well-being, peace, social good, common good)

- **Freedom** & autonomy (consent, choice, self-determination, liberty, empowerment)

- **Trust**

- **Sustainability** (environment, energy)

- **Dignity**

- **Solidarity** (social security, cohesion)

# Large number of "ethics frameworks"…

Table 2 Comparison of ethical principles in recent publications demonstrating the emerging consensus of 'what' ethical AI should aspire to be

| AI4People (published November 2018) (Floridi et al. 2018) | Five principles key to any ethical framework for AI (L Floridi and Clement-Jones 2019) | Ethics Guidelines for Trustworthy AI (Published April 2019) (European Commission 2019) | Recommendation of the Council of Artificial Intelligence (Published May 2019) (OECD 2019b) | Beijing AI Principles for R&D (Published May 2019) ('Beijing AI Principles' 2019) |
|---|---|---|---|---|
| Beneficence | AI must be beneficial to humanity | Respect for human autonomy | Inclusive growth, sustainable development and well-being | **Do good:** (covers the need for AI to promote human society and the environment) |
| Non-Maleficence | AI must not infringe on privacy or undermine security | Prevention of harm | Robustness, security and safety | **Be responsible:** (covers the need for researchers to be aware of negative impacts and take steps to mitigate them) **Control risks:** (covers the need for developers to improve the robustness and reliability of systems to ensure data security and AI safety) |
| | | | **Human-centred values** and fairness | **For humanity:** (covers the need for AI to serve humanity by conforming to human values including freedom and autonomy) |
| | | Fairness | Human-centred values **and fairness** | **Be diverse and inclusive:** (covers the need for AI to benefit as many people as possible) **Be ethical:** (covers the need to make the system as fair as possible, minimising discrimination and bias) |
| Explicability | AI systems must be understandable and explainable | Explicability | Transparency and explainability Accountability | **Be ethical:** (covers the need for AI to be transparent, explainable and predictable) |

For a more detailed comparison see Floridi and Cowls (2019) and Hagendorff (2019)

| Concepts | Basic notions relevant for debating ethical aspects |
|---|---|
| Principles | Ethical principles (e.g. values) |
| Concerns | Ways in which principles are threatened through AI systems use and development |
| Rules | **Strategies and guidelines for addressing the challenges** |

J. Morley et al. (2019) From what to how. https://ssrn.com/abstract=3830348

# From principles to practice

*I am a part of that power that always wants evil and always creates good. I am the spirit that always denies!*

Mephisto in *Faust*, J.W. Goethe

## What to do about AI to make it "ethical" (in practice)

| | |
|---|---|
| Rules, regulation | Checklists |
| Standards (e.g. IEEE) | Technologies |
| Councils, Boards | Consulting |
| Seals and labels | Good practice |
| Virtues | … |

- Current research topic in the AI academic literature

- Sub-fields of AI/ML, e.g. XAI

- Algorithms mostly for
  - Explainability
  - De-biasing
  - Privacy preservation

- Tools include
  - Data sets
  - Communities
  - Metrics
  - Process models

## Labels provide information about AI models

Inspiration from labels for food, clothing for consumers

**Exercise F**

Shift of responsibility to user

Fiction of consent: experience from
- Terms of Use
- Dark Patterns / GDPR agreement
- Etc.

# Standards

Existing standards for AI/autonomous systems
- Model process for addressing ethical concerns during systems design (IEEE 7000-2021)
- Transparency of autonomous systems (IEEE 7001-2021)
- Data privacy process (IEEE 7002-2022)
- Algorithmic bias considerations (IEEE P7003)
- Standards on child and student data governance (IEEE P7004)
- ...



IEEE P7000 https://ethicsinaction.ieee.org/p7000/

# From what to how: proposals

| Summaries | Notions | Procedures | Code | Infrastructure | Education | Ex-post assessment and agreement |
|---|---|---|---|---|---|---|
| Overviews and introductions | Frameworks and concepts | Process models | Algorithmic methods | Data sets | Training and tutorial | Audit |
| Case studies and examples | Criteria and checklists | Guidelines and codes of practice | Design patterns | Online communities | | License model |
| | Declarations | Standards | Software libraries | | | |
| | Metrics | | Software assistants | | | |
| Good practice | Regulation | Consulting | | Ethics councils and boards | Coaching | Labels, warnings, consent management |

Erich Prem (2023) From Ethical AI Frameworks to Tools: A review of approaches. In: AI and Ethics.

# Fairness

Dozens of notions of fairness: many have mathematical interpretations.
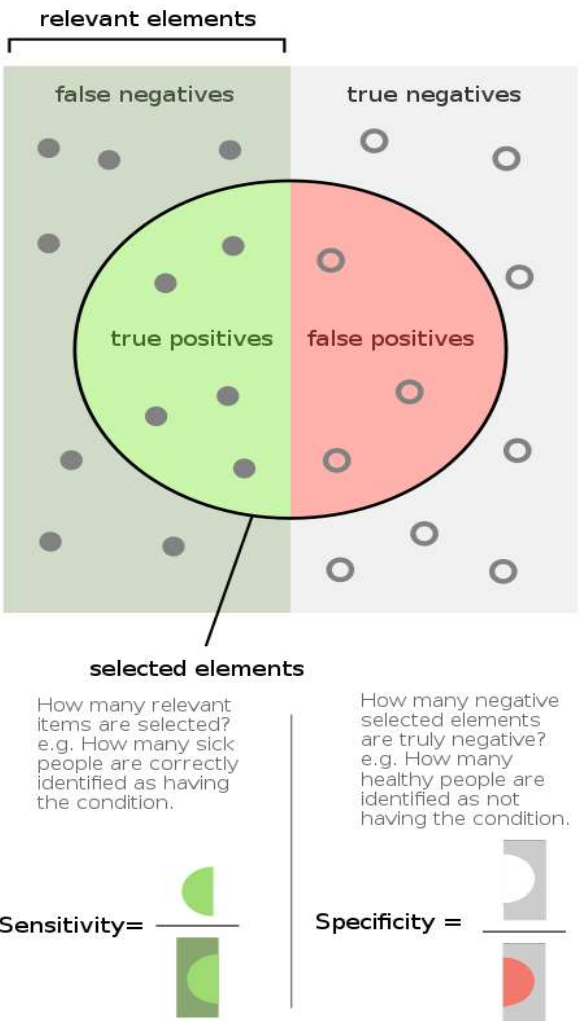
- Justice: adherence to the standards agreed in a society
- *Fairness: related evaluative judgement whether a decision (action) is morally right*
    - subjective
    - underlying idea of "all humans are equal"

But: is fairness "just" a mathematical notion?

In computer models the question is often **unavoidable,** i.e. in selecting a model, shaping the error function etc.

M. Seng Ah Lee, L. Floridi, J. Singh (2021) Formalising trade-offs beyond algorithmic fairness: lessons from ethical philosophy and welfare economics. https://ssrn.com/abstract=3679975

# For example, biases: what is *really* fair?

**Exercise B**



relevant elements

false negatives | true negatives

true positives | false positives

selected elements

How many relevant items are selected? e.g. How many sick people are correctly identified as having the condition.

How many negative selected elements are truly negative? e.g. How many healthy people are identified as not having the condition.

Sensitivity=

Specificity =

Assume: modelling default risk of a lender on a loan. Scenario: supervised learning, some "inappropriate" attribute present, e.g. race, gender, social status

- False positives (FP): lost opportunity (predicted default, but would have repaid)
- False negative (FN): lost revenue (predicted repayment, but defaulted)

Various error rates:
- True positive rate, sensitivity, probability that an actual positive will test positive. : (TPR)=TP/(TP+FN)
- True negative rate, specificity: (TNR)=TN/(FP+TN)
- False positive rate, fall-out: (FPR)=FP/(FP+TN)=1-TNR
- False negative rate (FNR)=FN/(FN+TP)=1-TPR
- Positive predictive value, precision: (PPV)=TP/(TP+FP)

M. Seng Ah Lee, L. Floridi, J. Singh (2021) Formalising trade-offs beyond algorithmic fairness: lessons from ethical philosophy and welfare economics. https://ssrn.com/abstract=3679975

# Which inequality is fair? A selection of ideas…

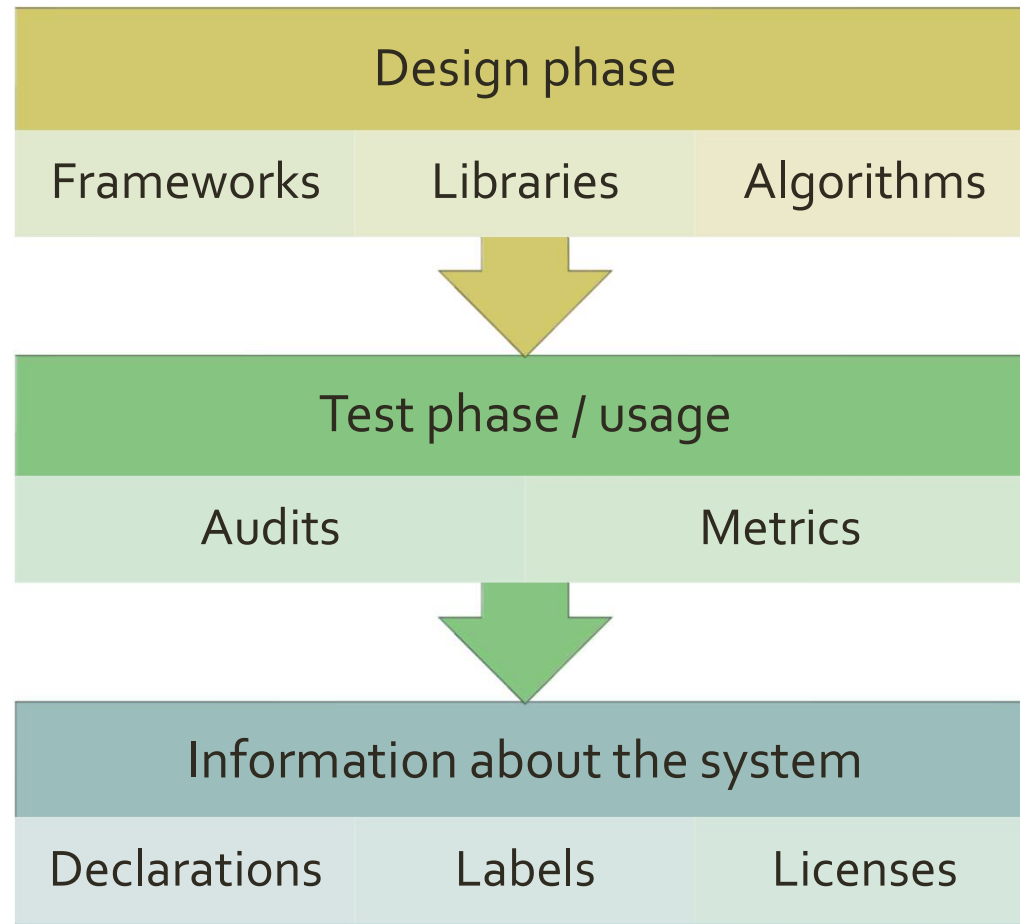| Fairness metric (literature) | Equalising | Intuition/example |
| --- | --- | --- |
| Maximise total accuracy | N/A | Most accurate model gives people the loan and interest they 'deserve' by minimising errors |
| Demographic parity, group fairness | Outcome | Black and white applicants have same loan approval rates |
| Equal opportunity | FNR | Among creditworthy applications, black and white applicants have similar approval rates |
| Predictive equality | FPR | Among defaulting applicants, black and white have similar rates of denied loans |
| Equal odds | TPR, TNR, PPV | Both of the above: Among creditworthy applicants, probability of predicting repayment is the same regardless of race |
| Counterfactual fairness | Prediction in counterfactual scenario | For each individual, if they were a different race, the prediction would be the same |
| Individual fairness | Outcome for 'similar' individuals | Each individual has the same outcome as another 'similar' individual of a different race |

**Not all inequalities can be removed.**

M. Seng Ah Lee, L. Floridi, J. Singh (2021) Formalising trade-offs beyond algorithmic fairness: lessons from ethical philosophy and welfare economics. https://ssrn.com/abstract=3679975

# Types of discrimination

| Inequality type | Example |
| --- | --- |
| Natural | Disability at birth |
| Socioeconomic | Parents' assets |
| Talent | Skills |
| Preference | Saving behaviour |
| Treatment | Job market discrimination |

Question of discrimination

- Certain characteristics should not result in disadvantages (often they have in the past)
  - ethnicity, gender, religion, age, disability, sexual orientation
- Often targets a change in society (policies)
- Distinction of in/acceptable inequalities, (non-)explainable discrimination, ir/relevant features
  - Income: relevant feature
  - Gender: irrelevant
- In practice very difficult!
- Modern proposal: include only attributes that an individual can directly influence. (No one should be treated worse just out of bad luck.

M. Seng Ah Lee, L. Floridi, J. Singh (2021) Formalising trade-offs beyond algorithmic fairness: lessons from ethical philosophy and welfare economics. https://ssrn.com/abstract=3679975

**Tools and methods for various design phases**

Design phase

| Frameworks | Libraries | Algorithms |

Test phase / usage

| Audits | Metrics |

Information about the system

| Declarations | Labels | Licenses |

# Example LLM (e.g. ChatGPT)

**Creation**
- Data sources (quality, legality, ethicality, filtering…)
- Design issues (anthropomorphising)

↓

**Use**
- Usage, influence, effects, dangers
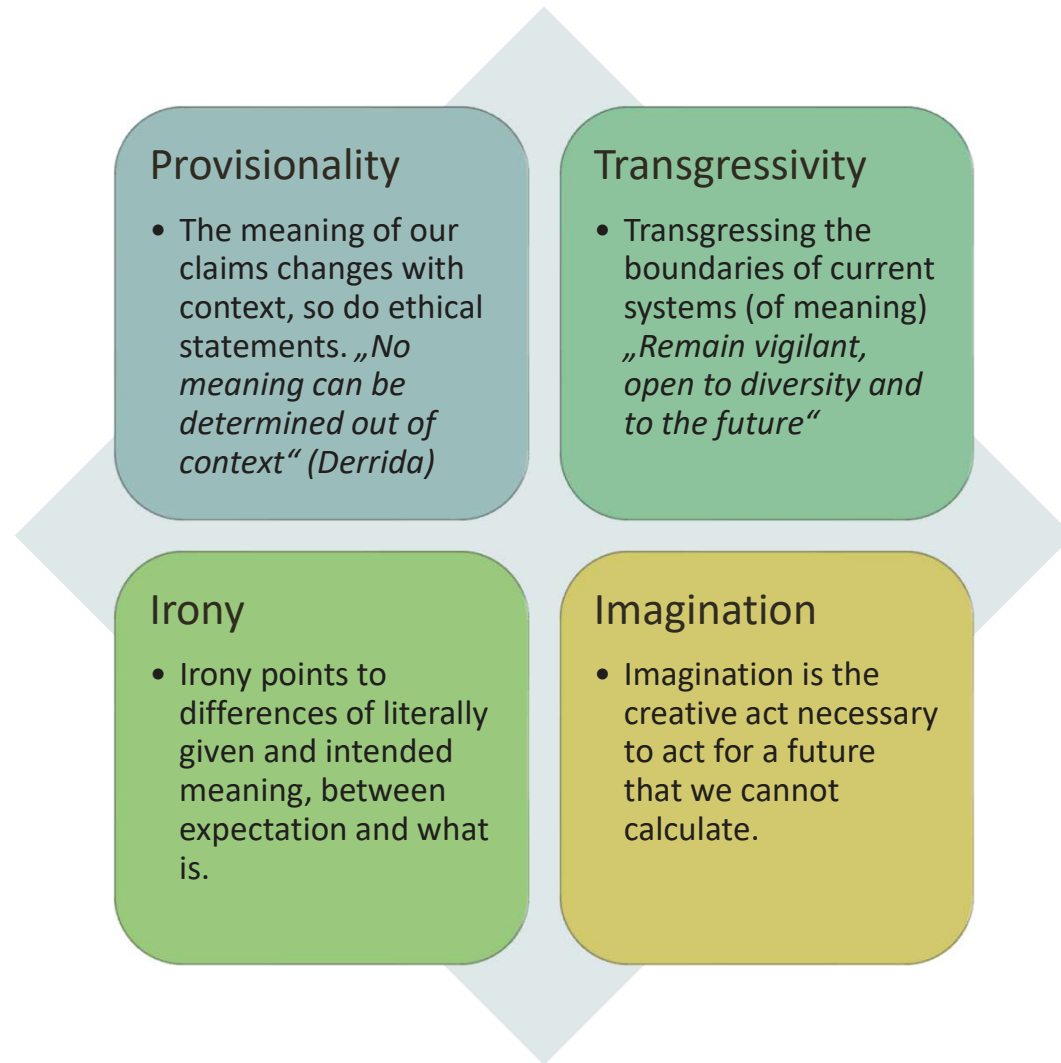
↓

**Power**
- Implications, politics, geopolitics



PKBnews.in
**Is Man Killed By AI? Belgian Man Commits Suicide After T… Chatbot**
A Belgian man has reportedly died by suicide after chatting with an AI-powered chatbot for six weeks. According to statements by his wife to…
vor 1 Tag

Euronews
**Man ends his life after an AI chatbot 'encouraged' him to s… himself to stop climate change**
A Belgian man reportedly ended his life following a six-week-long conversation about the climate crisis with an artificial intelligence (AI)…
vor 2 Wochen

VICE
**'He Would Still Be Here': Man Dies by Suicide After Talking… Chatbot, Widow Says**
A Belgian man recently died by suicide after chatting with an AI chatbot on an app called Chai, Belgian outlet La Libre reported.
vor 2 Wochen

Interesting Engineering
**Belgian woman blames ChatGPT-like chatbot ELIZA for he…**

# Final remarks

*Whereof one cannot speak, thereof one must be silent.*

L. Wittgenstein

# Four principles of an ethics for complex systems

## Provisionality
- The meaning of our claims changes with context, so do ethical statements. „*No meaning can be determined out of context*" (Derrida)

## Transgressivity
- Transgressing the boundaries of current systems (of meaning) „*Remain vigilant, open to diversity and to the future*"

## Irony
- Irony points to differences of literally given and intended meaning, between expectation and what is.

## Imagination
- Imagination is the creative act necessary to act for a future that we cannot calculate.

# What is digital humanism?

**DIGITAL** HUMANISM

Digital humanism is an initiative to actively shape digitization so that people and society are the focus.

Digital humanism is a call to use digital technologies to protect human rights and develop democracy.

Digital humanism acknowledges the key role of digital technologies for progress and innovation and seeks to expand it to sustain and expand our social achievements.

https://dighum.ec.tuwien.ac.at/dighum-manifesto/

# Contact me



Dr.phil. Dr.tech. Erich Prem (MBA)
*Managerial economist*


www.erichprem.at
prem at eutema.com
@ErichPrem


eutema GmbH
www.eutema.com


Association for digital humanism
www.digitalhumanism.at



https://dighum.ec.tuwien.ac.at/per
spectives-on-digital-humanism/

eu|te|ma
TECHNOLOGY MANAGEMENT

UNIVERSITY OF VIENNA

TU WIEN
TECHNISCHE UNIVERSITÄT WIEN
Vienna University of Technology