# (HOW) SHOULD AI BE REGULATED?

**Prof. Dr. Virginia Dignum**

**Chair of Social and Ethical AI - Department of Computer Science**

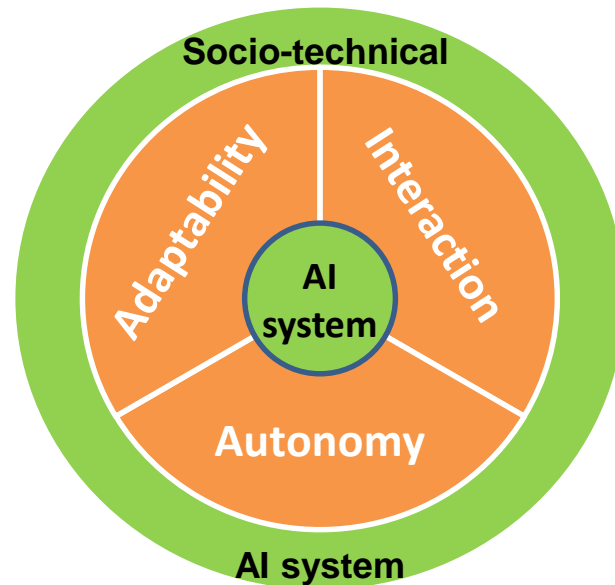**Email: virginia@cs.umu.se - Twitter: @vdignum**

UMEÅ UNIVERSITY

# WHAT IS AI?

- A technology
  - Currently mostly pattern matching (stochastic, non-deterministic)

- A field of science
  - Model intelligence as means to understand intelligence

- An entity
  - Magic, all-knowing, all-powerful

UMEÅ UNIVERSITY

# WHAT IS AI?

- Not just algorithm

- Not just machine learning

- But

- AI applications are not alone
  - Socio-technical AI systems

# AI IS NOT INTELLIGENCE!

- What AI systems cannot do (yet)
  - Common sense reasoning
    - Understand context
    - Understand meaning
  - Learning from few examples
  - Learning general concepts
  - Combine learning and reasoning

- What AI systems can do (well)
  - Identify patterns in data
    - Images
    - Text
    - Video
  - Extrapolate those patterns to new data
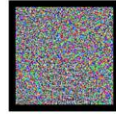  - Take actions based on those patterns

# AI IS NOT INTELLIGENCE!



"panda"  +  Adversarial Noise  =  "gibbon"

# AI IS NOT INTELLIGENCE!

# WHAT ARE THE RISKS?

bias and prejudice

discrimination

## misuse

loss of self-determination

removing human responsibility

devaluation of human skills

lack of control

UMEÅ UNIVERSITY

# WHAT ARE THE GAINS?

- Medicine

- Climate

- Education

- Work

- Communication

BRIGHT FUTURE AHEAD

UMEÅ UNIVERSITY

# WHAT IS RESPONSIBLE AI?

Responsible AI is

- Ethical

- Lawful

- Reliable

- Beneficial

Responsible AI recognises that

- AI systems are artefacts

- We set the purpose

UMEÅ UNIVERSITY

# RESPONSIBLE AI

- AI can potentially do a lot. <span style="color:red">Should it?</span>

- Who should decide?

- Which values should be considered? Whose values?

- How do we deal with dilemmas?

- How should values be prioritized?

- .....

> Essential question: shoud we use AI here?

UMEÅ UNIVERSITY

# AI AND ETHICS - SOME CASES

- Not just trolley problems!!

- Automated manufacturing
  - How can technical advances combined with education programs (human resource development) help workers practice new sophisticated skills so as not to lose their jobs?

- Chatbots
  - Mistaken identity (is it a person or a bot?)
  - Manipulation of emotions / nudging / behaviour change support

- Automated decision making
  - Accuracy versus explainability
  - E.g. 95% accuracy but no explanation or 80% accuracy but always explains

UMEÅ UNIVERSITY

# PRINCIPLES AND GUIDELINES

## Responsible / Ethical / Trustworthy....



https://ec.europa.eu/digital-single-market/en/high-level-expert-group-artificial-intelligence



https://ethicsinaction.ieee.org



https://www.oecd.org/going-digital/ai/principles/
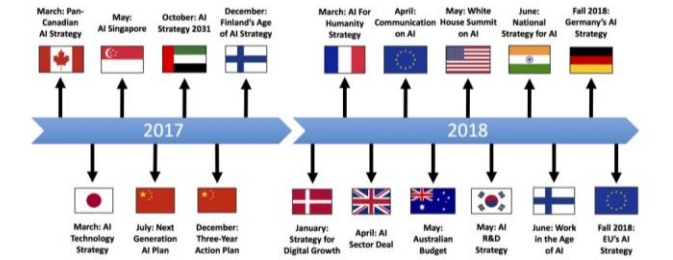
UMEÅ UNIVERSITY

# MANY INITIATIVES (AND COUNTING...)

- Strategies / positions
    - IEEE Ethically Aligned Design
    - European Union
    - OECD
    - WEF
    - Council of Europe
    - National strategies:
        - Tim Dutton, https://medium.com/politics-ai/an-overview-of-national-ai-strategies-2a70ec6edfd
    - ...
- Declarations
    - Asilomar
    - Montreal
    - ...



Check Alan Winfield blog:
http://alanwinfield.blogspot.com/2019/04/an-updated-round-up-of-ethical.html

| EU HLEG | OECD | IEEE EAD |
|---|---|---|
| • Human agency and oversight<br>• **Technical robustness and safety**<br>• Privacy and data governance<br>• **Transparency**<br>• **Diversity**, non-discrimination and fairness<br>• **Societal and environmental well-being**<br>• **Accountability** | • benefit people and the planet<br>• respects the rule of law, **human rights**, democratic values and **diversity**,<br>• include appropriate safeguards (e.g. human intervention) to ensure a **fair and just society**.<br>• **transparency** and responsible disclosure<br>• **robust, secure and safe**<br>• Hold organisations and individuals **accountable** for proper functioning of AI | • How can we ensure that A/IS do not infringe **human rights**?<br>• effect of A/IS technologies on **human well-being**.<br>• How can we assure that designers, manufacturers, owners and operators of A/IS are responsible and **accountable**?<br>• How can we ensure that A/IS are **transparent**?<br>• How can we extend the benefits and minimize the risks of AI/AS technology being misused? |

# BUT ENDORSEMENT IS NOT (YET) COMPLIANCE

UMEÅ UNIVERSITY

| EU HLEG | OECD | IEEE EAD |
|---|---|---|
| • Human agency and oversight<br>• **Technical robustness and safety**<br>• Privacy and data governance<br>• **Transparency**<br>• **Diversity**, non-discrimination and fairness<br>• **Societal and environmental well-being**<br>• **Accountability** | • benefit people and the planet<br>• respects the rule of law, **human rights**, democratic values and **diversity**,<br>• include appropriate safeguards (e.g. human intervention) to ensure a **fair and just society**.<br>• **transparency** and responsible disclosure<br>• **robust, secure and safe**<br>• Hold organisations and individuals **accountable** for proper functioning of AI | • How can we ensure that A/IS do not infringe **human rights**?<br>• effect of A/IS technologies on **human well-being**.<br>• How can we assure that designers, manufacturers, owners and operators of A/IS are responsible and **accountable**?<br>• How can we ensure that A/IS are **transparent**?<br>• How can we extend the benefits and minimize the risks of AI/AS technology be |
| **regulation** | **observatory** | **standards** |

The promise of AI:
Better decisions

# HOW DO WE MAKE DECISIONS?

# HOW DO WE MAKE DECISIONS TOGETHER?



UMEÅ UNIVERSITY

# DESIGN IMPACTS DECISIONS IMPACTS SOCIETY

- Choices
- Formulation
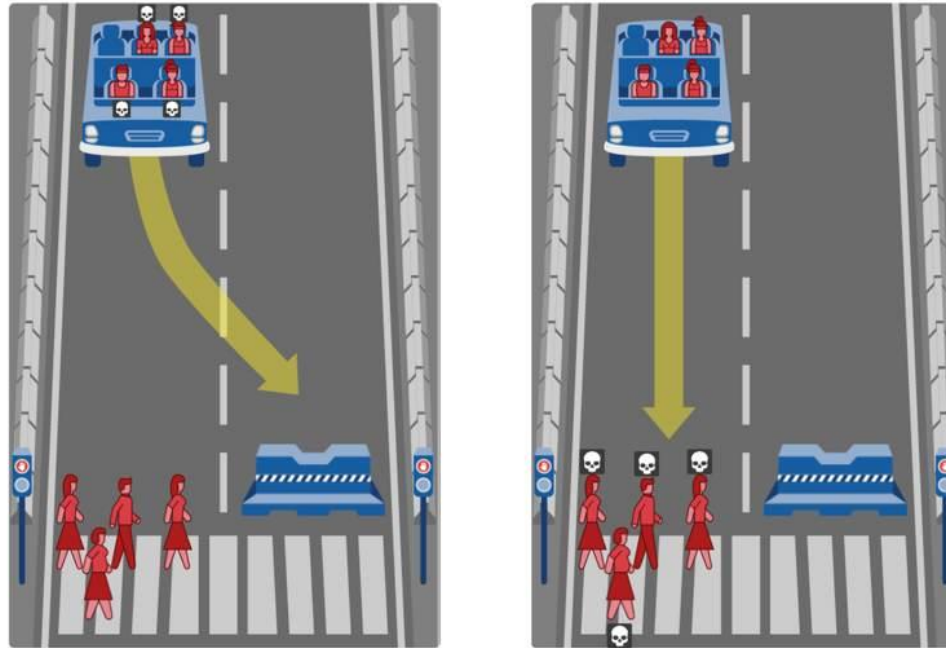- Involvement
- Legitimacy
- Aggregation

UMEÅ UNIVERSITY
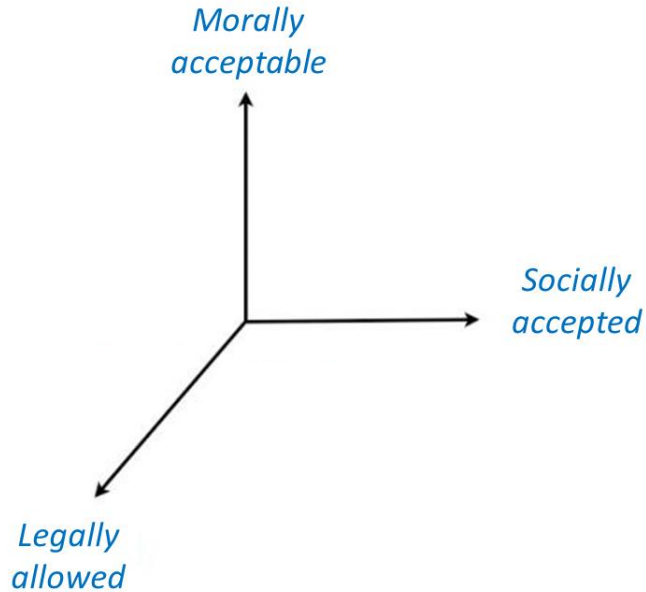
# WHICH DECISIONS SHOULD AI MAKE?

# WHICH DECISIONS SHOULD AI MAKE?



What should the self-driving car do?

# HOW SHOULD AI MAKE DECISIONS?

Morally
acceptable

Socially
accepted

Legally
allowed



UMEÅ UNIVERSITY

# TAKING RESPONSIBILITY

- **<u>in</u>** Design
  - o Ensuring that development <u>processes</u> take into account ethical and societal implications of AI and its role in socio-technical environments

- **<u>by</u>** Design
  - o Integration of ethical reasoning abilities as part of the <u>behaviour</u> of artificial autonomous systems

- **<u>for</u>** Design(ers)
  - o Research integrity of <u>stakeholders</u> (researchers, developers, manufacturers,...) and of institutions to ensure regulation and certification mechanisms

UMEÅ UNIVERSITY

# *IN* DESIGN: PROCESS
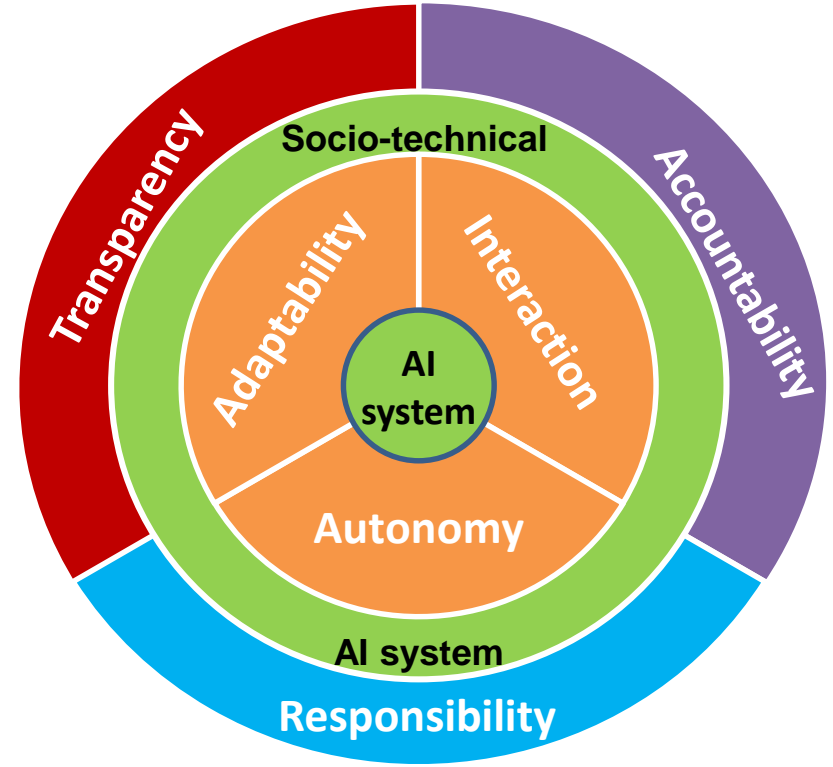
"Do things right, and do the right things."

PETER DRUCKER

UMEÅ UNIVERSITY

# TAKING RESPONSIBILITY: ART

- AI needs ART
  - o **A**ccountability
  - o **R**esponsibility
  - o **T**ransparency



UMEÅ UNIVERSITY

# ETHICS _IN_ DESIGN– DOING IT RIGHT

- Principles for Responsible AI = ART
  - o **A**ccountability
    - ▪ Explanation and justification
    - ▪ Design for values
  - o **R**esponsibility
    - ▪ Autonomy
    - ▪ Chain of responsible actors
    - ▪ Human-like AI
  - o **T**ransparency
    - ▪ Data and processes
    - ▪ Not just about algorithms

- AI systems (will) take decisions that have ethical grounds and consequences
- Many options, not one 'right' choice
- Need for design methods that ensure

UMEÅ UNIVERSITY

# ETHICS *IN* DESIGN: AI – DOING IT RIGHT

- Principles for Responsible AI = ART
  - **A**ccountability
    - Explanation and justification
    - Design for values
  - **R**esponsibility

  - **T**ransparency



- Optimal AI is explainable AI
- Many options, not one 'right' choice

# ETHICS _IN_ DESIGN: AI – DOING IT RIGHT

- Principles for Responsible AI = ART
  - **A**ccountability
    - Explanation and justification
    - Design for values
  - **R**esponsibility
    - Autonomy
    - Chain of responsible actors
    - Human-like AI
  - **T**ransparency

What should the self-driving car

The machine is not responsible!

UMEÅ UNIVERSITY

# ETHICS _IN_ DESIGN: AI – DOING IT RIGHT

- Principles for Responsible AI = ART
  - **A**ccountability
    - Explanation and justification
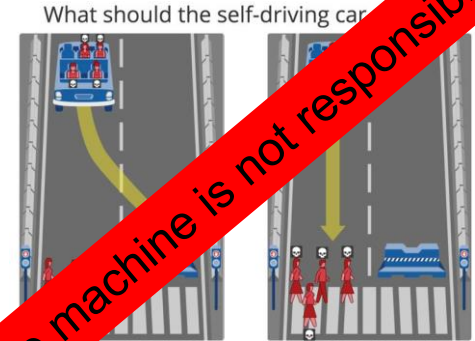    - Design for values
  - **R**esponsibility
    - Autonomy
    - Chain of responsible actors
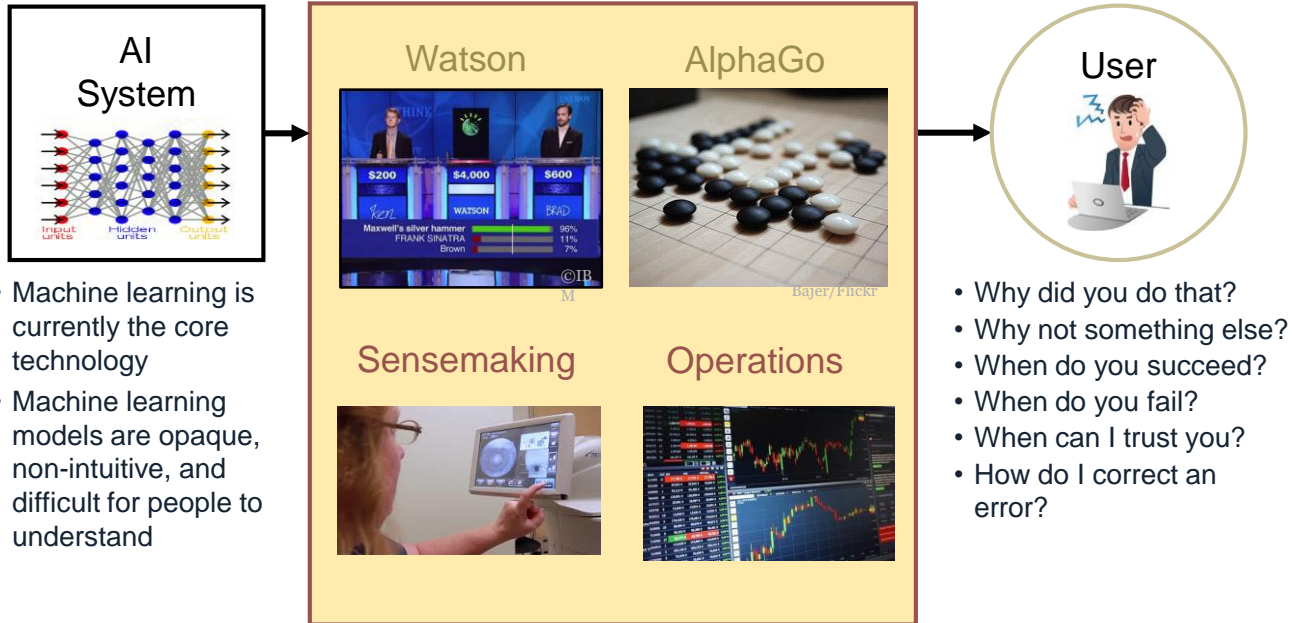    - Human-like AI
  - **T**ransparency
    - Data and processes
    - Algorithms
    - Choices and decisions

UMEÅ UNIVERSITY

# CONCERN: EXPLAINABLE AI

**AI System**



- Machine learning is currently the core technology
- Machine learning models are opaque, non-intuitive, and difficult for people to understand

Watson


©IBM

AlphaGo


Bajer/Flickr

Sensemaking



Operations



**User**



- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

UMEÅ UNIVERSITY

# WHAT IS AN EXPLANATION?

Correct
Compreensible
Timely
Complete
Parsimonous

UMEÅ UNIVERSITY

# NO AI WITHOUT EXPLANATION

- XAI is for the user:
  - Who depends on decisions, recommendations, or actions of the system
  - Just in time, clear, concise, understandable
- XAI is about:
  - provide an explanation of individual decisions
  - enable understanding of overall strengths & weaknesses
  - convey an understanding of how the system will behave in the future
  - convey how to correct the system's mistakes

UMEÅ UNIVERSITY

# *BY* DESIGN: ARTIFICIAL AGENTS

- Can we teach ethics to AI?

- <u>Should</u> we teach ethics to AI?
  - What does that mean?
  - What is needed?

- Decisions matter
  - Our decisions matter

- Good results matter
  - Accurate/efficient/fair/sustainable?

- The dilemmas
  - 95% accurate but not explainable or 80% accurate but explainable?
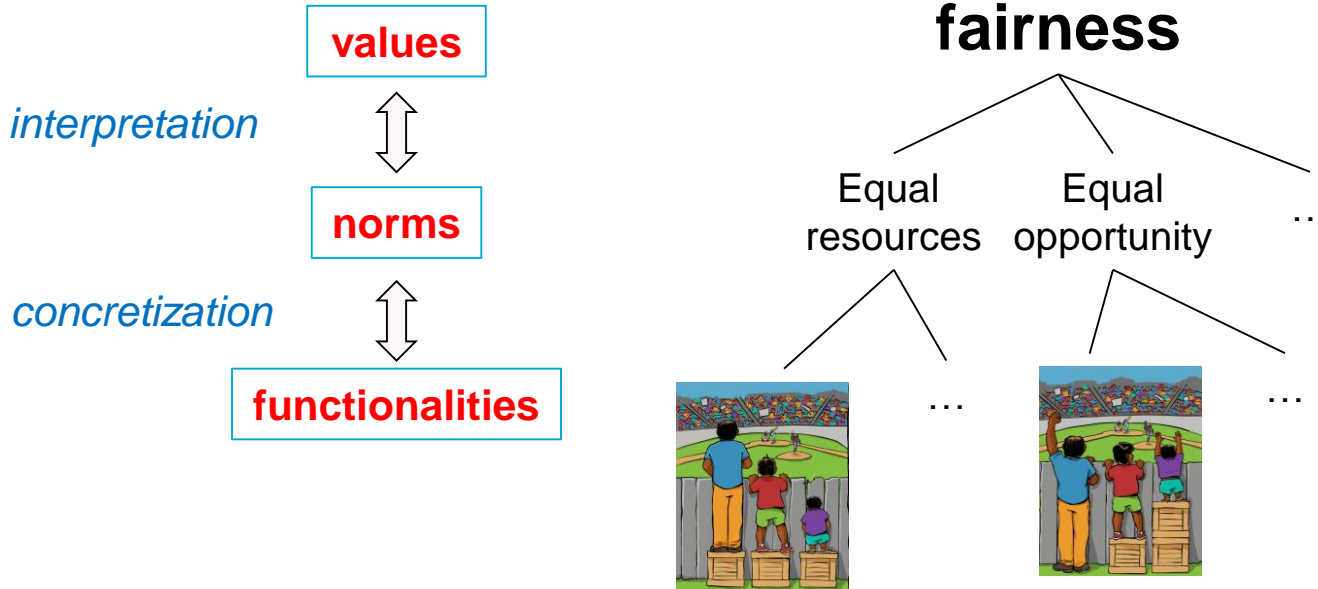  - Energy costs: more AI per joule!

UMEÅ UNIVERSITY

# WHICH VALUES – WHOSE VALUES

- Sources

  - Society (Designer, Users, Owner, Manufacturer)

  - Law: legislation, standards

  - Ethics

- But

  - Who decides who has a say?

  - How to make choices and tradeoffs between conflicting values?

  - How to verify whether the designed system embodies the intended values?
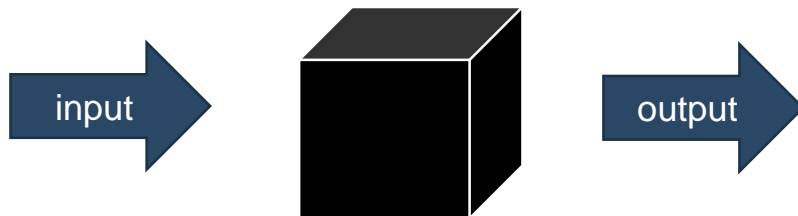
UMEÅ UNIVERSITY

# DECISIONS MATTER!

**values**

*interpretation*

**norms**

*concretization*

**functionalities**

fairness

Equal resources      Equal opportunity      …

      …            …

UMEÅ UNIVERSITY

# ONE PROBLEM

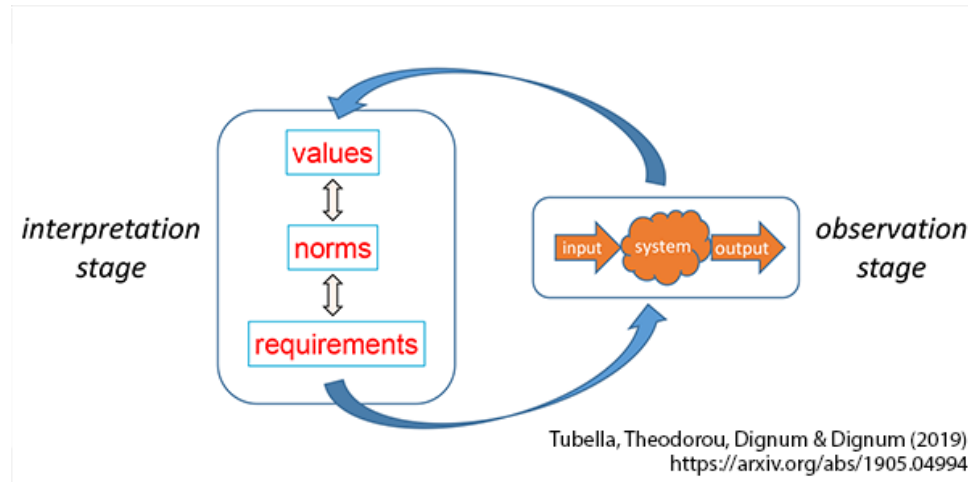- **black boxes** cannot always be avoided



input → output

- Still, we need to **trust** systems. **Check** their **compliance** against our **values**.
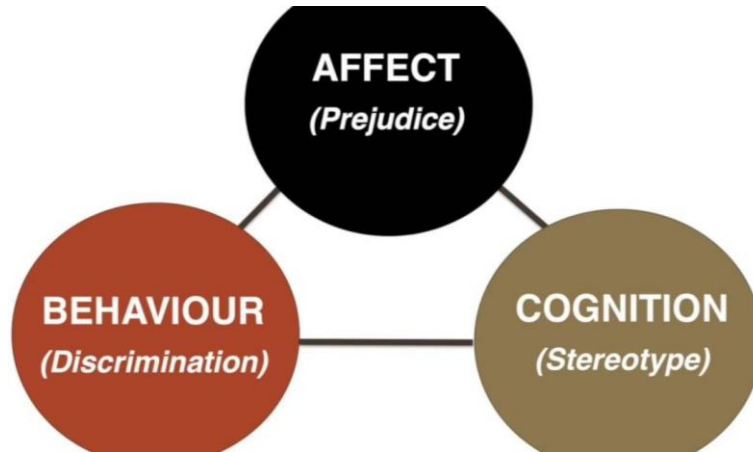
UMEÅ UNIVERSITY

# GLASS BOX APPROACH

- Doing the right thing
    - Elicit, define, agree, describe, report

- Doing it right
    - Explicit values, principles, interpretations, decisions
    - Evaluate input/output against principles



Tubella, Theodorou, Dignum & Dignum (2019)
https://arxiv.org/abs/1905.04994

# CONCERN: BIAS AND DISCRIMINATION

- Bias is inherent on human data – we need bias to make sense of world

- But we dont want AI to be prejudiced!

- Unbiasing: Are we creating new bias ?

# BIAS IS MORE THAN BIASED DATA

- Who is collecting the data?

- Whose data is being collected?

- Which data is collected?
  - Why don't we keep information about the colour of the socks of the data collector?

- How and by who is the data labelled?
  - Data farms, exploitation

- What is the training data?
  - The cheapest and easiest to attain?

UMEÅ UNIVERSITY

# GUIDELINES TO DEVELOP ALGORITHMS RESPONSIBLY

- Who will be affected?

- What are the decision criteria we are optimising for?

- How are these criteria justified?

- Are these justifications acceptable in the context we are designing for?

- How are we training our algorithms?

UMEÅ UNIVERSITY

# *FOR* DESIGN(ERS): PEOPLE

- **Regulation**
- **Certification**
- **Standards**
- **Conduct**

**AI principles are principles for us**

# REGULATION AND CERTIFICATION

- Taking an ethical perspective
    - Ethics is the new green
    - Business differentiation
    - Certification to ensure public acceptance



- Principles and regulation are drive for transformation
    - Better solutions
    - Return on Investment

UMEÅ UNIVERSITY

# Recommendations for trustworthy AI – Main issues

1. Empower and protect humans and society

2. Take up a tailored approach to the AI market

3. Secure a Single European Market for Trustworthy AI

4. Enable AI ecosystems thorough sectoral multi stakeholder alliances

5. Foster the European data economy

6. Exploit the multi-faceted role of the public sector

7. Strengthen and unite Europe's research capabilities

8. Nurture education

9. Adopt a risk-based governance approach to AI and ensure an appropriate regulatory framework

10. Stimulate an open and lucrative investment environment

11. Embrace a holistic way of working

INDEPENDENT
HIGH-LEVEL EXPERT GROUP ON
ARTIFICIAL INTELLIGENCE
SET UP BY THE EUROPEAN COMMISSION

AI

POLICY AND INVESTMENT RECOMMENDATIONS
FOR
TRUSTWORTHY AI

European
Commission

# WHAT IS BEING REGULATED?

- A computational technology for decision making?

- A field of scientific research that studies theories and methods for adaptability, interaction and autonomy of machines?

- An intelligent entity that acts autonomously in (our) environment?

UMEÅ UNIVERSITY

# ALTERNATIVES / COMPLEMENTS TO REGULATION

- Standards (IEEE, ISO)
    - soft governance; non mandatory to follow
    - demonstrate due diligence and limit liability
    - user-friendly integration between products

- Advisory panels and ethics officers
    - Set and monitor ethical guidelines
    - able to veto any projects or deliverables that do not adhere to guidelines

- Assessment lists for trustworthy, ethical, AI (EU)
    - responsible AI is more than ticking boxes
    - Means to assess maturity are needed

- Education and training

- Appeal to civic duty / voluntary implementation (Australia)

# TRUSTWORTHY AI



## CERTIFICATION FOR AI?

# CONCERN: WHO IS DEVELOPING AI?

- 18% researchers at conferences are women

- 80% professors are men

- Workforce
  - Google: 2,5% black, 3,6% Latino, 10% women
  - Facebook: 3,8% black, 5% Latino, 15% women

UMEÅ UNIVERSITY

- Design impacts decisions impacts society impacts design

- Regulation requires understanding what and why regulate

- AI systems are tools, artefacts made by people:
  We set the purpose

- AI can give answers, but we ask the questions