# RESPONSIBLE ARTIFICIAL INTELLIGENCE

Prof. Dr. Virginia Dignum

Chair of Social and Ethical AI - Department of Computer Science
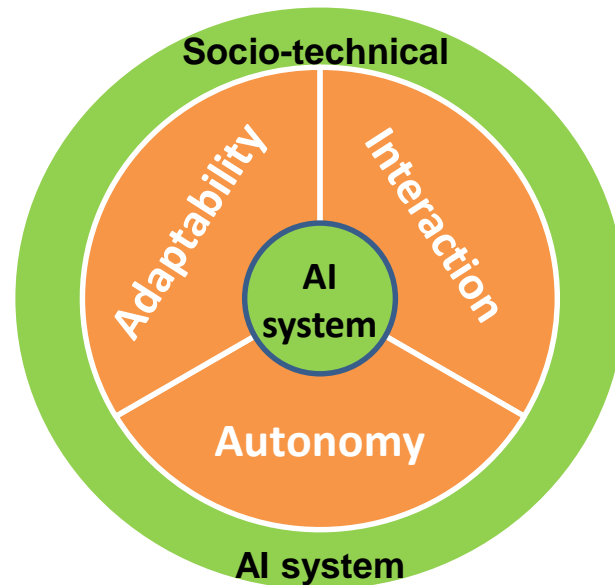
Email: virginia@cs.umu.se - Twitter: @vdignum

UMEÅ UNIVERSITY

# WHAT IS AI?

- Not just algorithm

- Not just machine learning


- But

- AI applications are not alone
    o 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
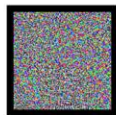  - Take actions based on those patterns

UMEÅ UNIVERSITY

# AI IS NOT INTELLIGENCE!



Adversarial Noise

"panda" + = "gibbon"

# AI IS NOT INTELLIGENCE!

# WHAT IS RESPONSIBLE AI?

Responsible AI is

- Ethical

- Lawful

- Reliable

- Beneficial

Responsible AI recognises that

- AI systems are artefacts

- We set the purpose

- We are responsible!

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?

- .....

# 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/
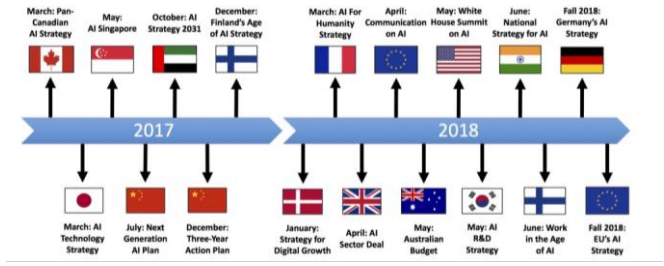
UMEÅ UNIVERSITY

# MANY INITIATIVES (AND COUNTING...)

- Strategies / positions
  - IEEE
  - European Union
  - OECD
  - WEF
  - Council of Europe
  - Many national strategies
  - ...

- Declarations
  - Asilomar
  - Montreal
  - ...



https://arxiv.org/ftp/arxiv/papers/1906/1906.11668.pdf lists 84!

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 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?

# DESIGN IMPACTS DECISIONS IMPACTS SOCIETY







- Choices
- Formulation
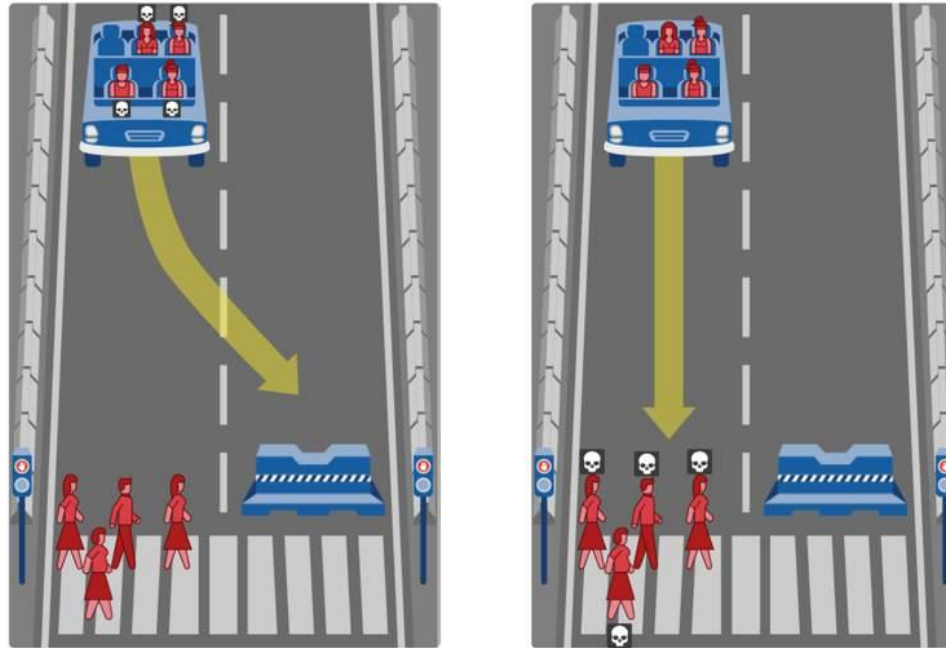- Involvement
- Legitimacy
- Voting system

UMEÅ UNIVERSITY
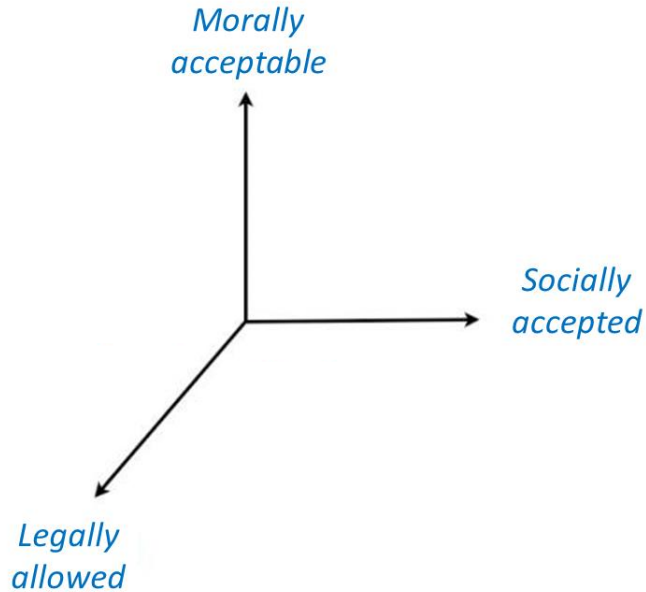
# WHICH DECISIONS SHOULD AI MAKE?

# WHICH DECISIONS SHOULD AI MAKE?



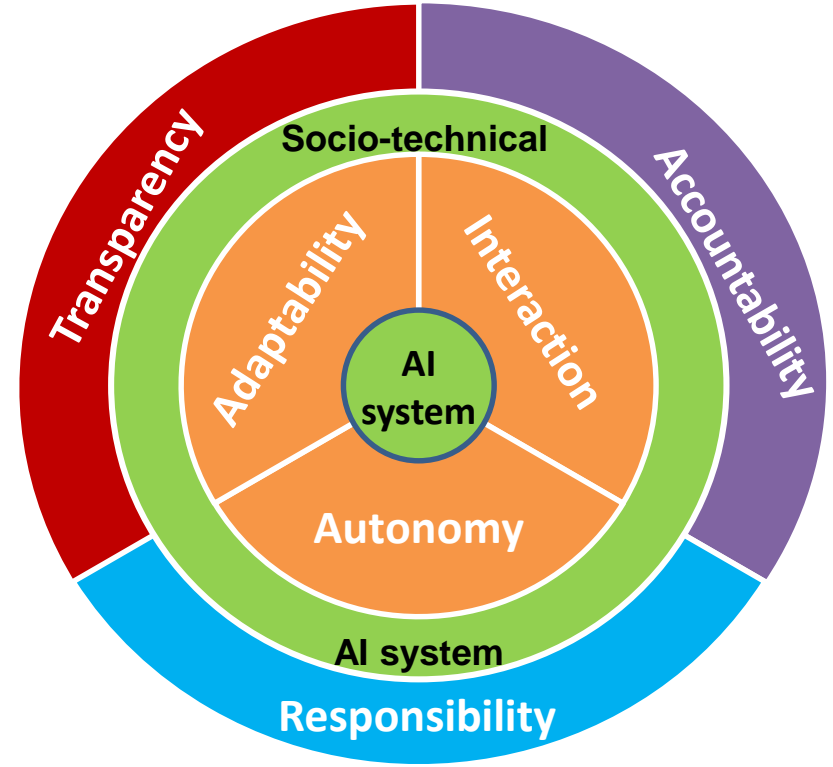What should the self-driving car do?

# 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: ART

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



UMEÅ UNIVERSITY

# ACCOUNTABILITY

- 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
- Explanation is for the user: context matters

UMEÅ UNIVERSITY

# CHALLENGE: NO AI WITHOUT EXPLANATION



- Explanation is for the user:
  - Different needs, different expertises and interests
  - Just in time, clear, concise, understandable, correct

- Explanation is about:
  - individual decisions and the 'big picture'
  - enable understanding of overall strengths & weaknesses
  - convey an understanding of how the system will behave in the future
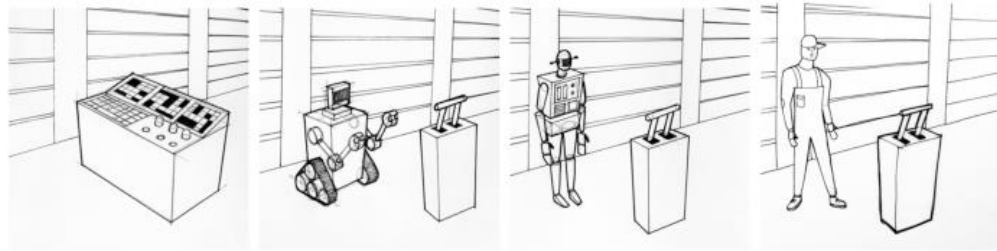  - convey how to correct the system's mistakes

# RESPONSIBILITY

- 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



The machine is not responsible!

# RESPONSIBILITY CHALLENGES

- Chain of responsibility
  - researchers, developerers, manufacturers, users, owners, governments, …
  - Liability and conflict settling mechanisms

- Human-like systems
  - Robots, chatbots, voice…
  - Expectations
  - Vulnerable users
  - Mistaken identity



https://ieeexplore.ieee.org/document/7451743

- Responsibility for choices
  - 95% accurate but no explanation or 80% accurate with explanation?
  - Fairness or sustainability?

UMEÅ UNIVERSITY
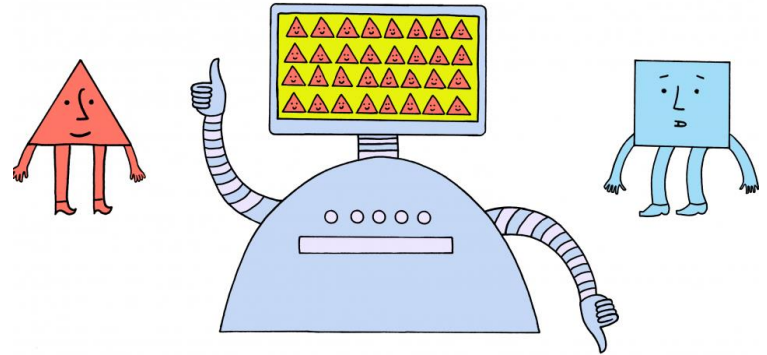
# TRANSPARENCY

- 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

# CHALLENGE: BIAS AND DISCRIMINATION

Remember: AI systems extrapolate patterns from data to take action

- Bias is inherent on human data
  - we need bias to make sense of world
- Bias leads to stereotyping and prejudice
- Bias is more than biased data

UMEÅ UNIVERSITY

# BY DESIGN: ARTIFICIAL AGENTS

- Can AI systems be ethical?
  - What does that mean?
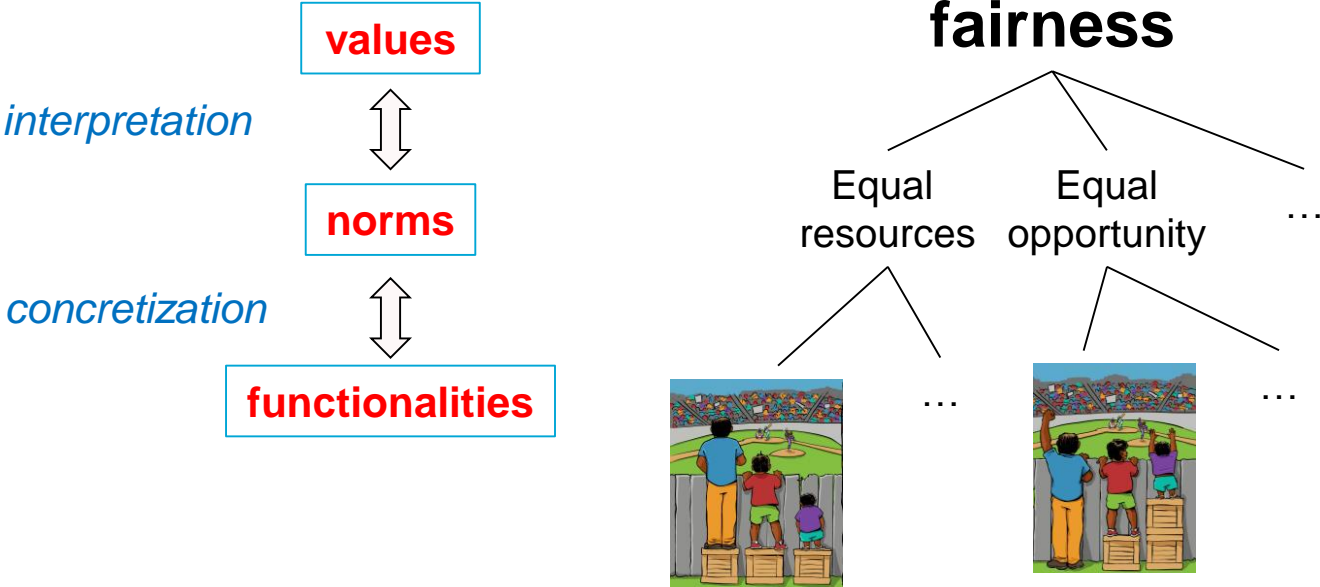  - What is needed?

- Design for values

# ETHICAL BEHAVIOR

- **<u>Should</u>** we teach ethics to AI?
- Understanding ethics
  - Which values? Whose values?
  - Who gets a say?
- Using ethics
  - What is proper action given a values?
  - Are ethical theories of use?
  - How to prioritise values?
  - Is knowing ethics enough?
- Ethical reasoning
  - Many different theories
    - Utilitarian, Kantian, Virtues…)
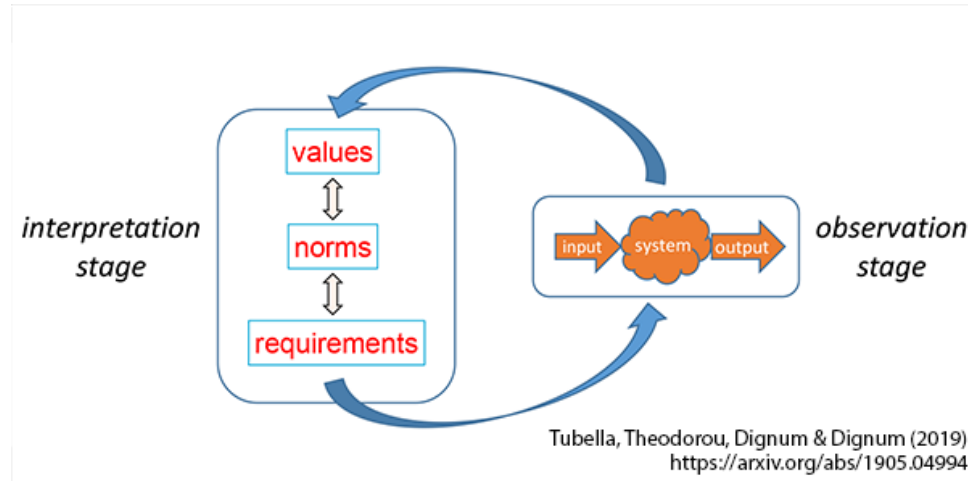  - Highly abstract
  - Do not provide ways to resolve conflicts



UMEÅ UNIVERSITY

# DESIGN FOR VALUES

values

*interpretation*

norms

*concretization*

functionalities

fairness

Equal resources     Equal opportunity     …

…     …

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

# *FOR* DESIGN(ERS): PEOPLE

- Regulation

- Certification

- Standards

- Conduct

- AI principles are principles for us

UMEÅ UNIVERSITY

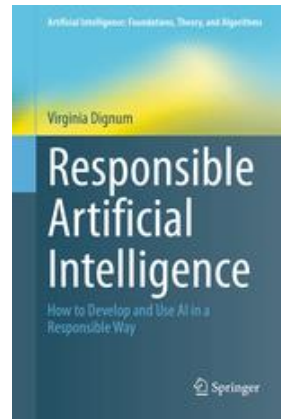# FOR DESIGN: TRUSTWORTHY AI



- Regulation and certification

- Codes of conduct

- Human-centered

- AI as driver for innovation

UMEÅ UNIVERSITY

- Design impacts decisions impacts society impacts design

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

- AI can give answers, but we ask the questions

- AI needs ART (Accountability, Responsibility, Transparency)

# RESPONSIBLE ARTIFICIAL INTELLIGENCE

# WE ARE RESPONSIBLE

Email: virginia@cs.umu.se

Twitter: @vdignum

UMEÅ UNIVERSITY