

GUIDELINES FOR FACILITATING THE LEARNING OF ARTIFICIAL INTELLIGENCE By School students of grades 7-12

FACILITATE – AI: Guidelines for facilitating the learning of Artificial Intelligence (AI) by School Students of Grades 7-12

Project Number: 2021-1-CY01-KA220-SCH-000032567

Co-funded by the European Union



Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Education and Culture Executive Agency (EACEA). Neither the European Union nor EACEA can be held responsible for them.







C1 Training course: Ethics About AI

- Result 1, A1/T1 Module Number and Area/ Topic: Ethics About Al Module owners: UCY, CY
- **Coordinator:** Antonis Kakas







Al Ethics A high-level Overview

AI Ethics is very Important

- Dangers vs Opportunities
- Levels of Dangers

AI Ethics is Hard

- How do we build Ethical Machines?
- Teach or Grow?

AI Ethics in EU

- Guidelines, Regulations and Laws!
 EU AI act
- But







European-approach-artificial-intelligence

EU AI Ethics

European Artificial Intelligence Act: the **right** framework for trustworthy AI and innovation?

Shaping Europe's digital future

Home Policies Activities News Library Funding Calendar Consultations

Home > Policies > A European approach to artificial intelligence

A European approach to artificial intelligence

The EU's approach to artificial intelligence centers on excellence and trust, aiming to boost research and industrial capacity while ensuring safety and fundamental rights.

The way we approach Artificial Intelligence (AI) will define the world we live in the future. To help building a resilient <u>Europe for the Digital Decade</u>, people and businesses should be able to enjoy the benefits of AI while feeling safe and protected.

The <u>European AI Strategy</u> aims at making the EU a world-class hub for AI and ensuring that AI is human contribution tructure the Surphane approach to









Al Ethics Teacher Training Videos

FROM EU ERASMUS+ Trustworthy AI

<u>https://www.trustworthyaiproject.eu/teaching-trustworthy-ai-oers/</u>

Responsible ai 1-0-1: from values to requirements



Dr. Andreas Theodorou andreas.theodorou@umu.se



@recklesscoding



UMEÅ UNIVERSITY

What is AI?



What is AI?

- Simulation of natural intelligence; the field of AI includes the study of theories and methods for adaptability, interaction and autonomy of machines (virtual or embedded)
- A (computational) technology that is able to infer patterns and possibly draw conclusions from data.
- An (autonomous) entity (e.g. when one refers to `an' AI); this is the most usual reference in media and science fiction.

"AS SOON AS IT WORKS, NO ONE CALLS IT AI ANY MORE."

John McCarthy

This Lack of agreement leads to...

- A constant re-writing of similar high-level policy statements:
 - Opens loopholes to be exploited.
 - We will see later more about this!
- Increases misconceptions (true AI, superintelligence, etc):
 - Al is not a magic abstraction.
 - Computation is a **physical process**; it requires *energy, space,* and *time*.

<u>Theodorou, A</u>. and Dignum V. (2020), *Towards ethical and socio-legal governance in AI*. Nature Machine Intelligence, 2(1). Andreas Theodorou | @recklesscoding

Narratives in the news



Andreas Theodorou | t: @recklesscoding

Bots, #StrongerIn, and #Brexit: Computational Propaganda during the UK-EU Referendum

COMPROP RESEARCH NOTE 2016.1

Philip N. Howard	Bence Kollanyi
Oxford University	Corvinus University
philip.howard@oii.ox.ac.uk	kollanyi@gmail.com

DISINFORMATION AND SOCIAL BOT **OPERATIONS IN THE RUN UP TO THE 2017** FRENCH PRESIDENTIAL ELECTION

and th

EMILIO FERRARA **UNIVERSITY OF SOUTHERN CALIFORNIA, INFORMATION SCIENCES INSTITUTE**

ABSTRACT

Recent accounts from researchers, journalists, as well as federal investigators, reached a unanimous conclusion: social media are systematically exploited to manipulate and alter public opinion. Some and have been coordinated by means of bots, social media trolled by

SPOTLIGHT STORY INSIDE THAILAND'S HAZY I EGALIZATION OF WEEK

CHRCCDIR

Millions of Americans Have Lost Jobs in the Pandemic—And Robots and AI Are Replacing Them Faster Than Ever

BUSINESS • COVID-19

2017 Eurobarometer

- **61%** of respondents have a **positive view** of robots
- 84% of respondents agree that robots can do jobs that are too hard/dangerous for people
- 68% agree that robots are a good thing for society because they help people
- 88% of respondents consider robotics a technology that requires careful management
- 72% of respondents think robots steal people's jobs

We need to build trust for our systems

• To perform as we expect them to.

- The implications from their development and deployment fall within:
 - Ethical
 - Legal
 - Social
 - Economic
 - Cultural



(**ESLEC**) specifications and values we want to protect and **communicating that** to the public.

Slide provided by DG Connect

European Union Background on Al





In this context: appointment of Independent High-Level Expert Group on Artificial Intelligence (AI HLEG) in June 2018



Slide provided by DG Connect

Ethics Guidelines for AI – Requirements



Human agency and oversight



Technical Robustness and safety



Privacy and data governance



Transparency



Diversity, nondiscrimination and fairness



Societal & environmental well-being



Accountability

To be continuously implemented & evaluated throughout AI system's life cycle





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

Which Values to implement? How?



EXAMPLE ASSESSMENT





- Did you build in mechanisms for notice and control over personal data depending on the use case (such as valid consent and possibility to revoke, when applicable)?
- Did you assess what level and definition of accuracy would be required in the context of the AI system and use case?
 - Did you assess how accuracy is measured and assured?
 - Did you put in place measures to ensure that the data used is comprehensive and up to date?
 - Did you put in place measures in place to assess whether there is a need for additional data, for example to improve accuracy or to eliminate bias?
- Did you assess:
 - to what extent the decisions and hence the sutcome made by the AI system can be understood?
 - to what degree the system's decision influences the rganisation's decision-making processes?
 - why this particular system was deployed in this specific area?
 - what the system's business model is (for example, how does it create value for the organisation)?
- In case the AI system interacts directly with humans:
 - Did you assess whether the AI system encourages humans to develop attachment and empathy towards the system?
 - Did you ensure that the AI system clearly signals that its social interaction is simulated and that it

Responsible ai & Human Control

WHAT IS AI ETHICS?

- Umbrella term for many things:
 - moral agency;
 - existential crisis/superintelligence;
 - adherence (or not) to human moral values;
 - trustworthiness.

It is not telling people if they are 'morally right' or morally wrong.'

RESPONSIBLE AI

 AI research includes the study of theories and methods for adaptability, interaction and autonomy of machines.

 Responsible AI is the understanding of how we can develop socially-beneficial systems and ensure human control.

Responsible AI = ethical + legal + robust + verifiable AI

Trustworthy AI (HLEG)

ETHICS != LAW != SOCIAL NORMS



Theodorou, A. and Dignum V. (2020), Towards ethical and socio-legal governance in AI. Nature Machine Intelligence, 2(1).

Socio-technical systems



Al involves

• In Design: Development is influenced by ESLEC issues.

• *By* Design: Integration of ethical abilities as part of the behaviour of artificial intelligent systems.

 For Design: Codes of conduct, standards, and certification processes that ensure the integrity of key stakeholders.

Dignum, V (2018). *Ethics in Artificial Intelligence: Introduction to the special issue.* Ethics and Information Technology, 20(1):1–3, 3 2018.

Values integration

WORKING AROUND THE PROBLEM

• Stop thinking of AI as 'magic data' instead of software.

• We do not necessary have to agree what *each value* means.

• We need to make their interpretations **explicit and transparent**.

Theodorou, A. and Dignum V. (2020), Towards ethical and socio-legal governance in AI. Nature Machine Intelligence, 2(1).

Interpretating values

- Structured and explicit process of translating translate abstract values into concrete norms and requirements.
- We aim to not only describe the norms themselves, but also the exact connection between abstract and concrete concepts in each context.
- Fulfilling the norm will be considered as adhering to the value.



Aler Tubella A., **Theodorou A.**, Dignum F., Dignum V. (2019). *Governance by Glass-box: implementing transparent moral bounds for AI behaviour*. IJCAI

Interpreting values



Aler Tubella A., **Theodorou A.**, Dignum F., Dignum V. (2019). *Governance by Glass-box: implementing transparent moral bounds for AI behaviour*. IJCAI

Accountability & responsibility

- Accountability is *backwards thinking*, provides an account of events after they have occurred.
- Responsibility is *forwards thinking*, i.e., acting to deter incidents and violations of our ethical and legal values from occurring.
- To be held accountable, you need to be held responsible...

accountability

• Umbrella term, covering amongst others:

- **1.** legal accountability ;
- 2. professional accountability;
- **3.** political accountability;
- **4.** administrative accountability; and
- 5. social accountability.
- Used not only to indicate punishment but also acceptance of responsibility.

Bovens, M. (2007). Analysing and Assessing Accountability: A Conceptual Framework. Eur. L. J.

accountability

Incidents will occur —and sometimes reoccur.

• It is essential for **maintaining the public's trust** in the technology.

• It is a **mean of control**: that compels adherence to specified regulations and practices to demonstrate *due diligence*.

Transparency in the IEEE 7001 Standard

"To consider an autonomous system transparent to inspection, the stakeholder should have the ability to request meaningful explanations of the system's status either at a specific moment or over a specific period or of the general principles by which decisions are made (as appropriate to the stakeholder) (Theodorou et. al., 2017)"

Winfield AFT, Booth S, Dennis LA, Egawa T, Hastie H, Jacobs N, Muttram RI, Olszewska JI, Rajabiyazdi F, **Theodorou A**, Underwood MA, Wortham RH and Watson E (2021) *IEEE P7001: A Proposed Standard on Transparency*. Front. Robot. AI

Theodorou A., Wortham R.H., and Bryson J. (2017). *Designing transparency for real time inspection of autonomous robots*. Connection Science, Vol. 29, Issue 3

Transparency as <u>Contestability</u>

• Our right to contest decisions made for us is not only protected by the Regulation (EU), 2016 GDPR

 Requires looking beyond why a decision was made: a decision needs to be both *correct* and *permissible*.

Aler Tubella, A., <u>Theodorou, A.</u>, Dignum, V., and Michael, L. (2020). *Contestable Black Boxes*. International Joint Conference on Rules and Reasoning.

Human control



Human control



HUMAN CONTROL OVER ARTIFICIAL INTELLIGENCE

Current frameworks



Meaningful Human control

- Human presence is not sufficient.
- We need to be able to:
 - track relevant reasons behind a decision;
 - verify their compliance to any policy; and
 - trace back to an individual along the chain who is aware and accepting of their responsibility.

Santoni de Sio, F., and Van den Hoven, J. (2018). Meaningful Human Control over Autonomous Systems: A Philosophical Account

Meaningful Human control

- **The world is not static**: contexts that change will demand changing levels of responsiveness.
- Requires us to consider wider issues around Responsible AI.



Methnani L., Aler Tubella A., Dignum V and <u>Theodorou A. (2021</u>). *Let Me Take Over: Variable Autonomy for Meaningful Human Control*. Front. Artif. Intell.











Takes action



You never trained in this environment. I will take over now.



Takes action



I am bored and taking control.





- **Maximises performance**: shifts autonomy both ways.
- **Maximises human control:** human presence adjusted as the environment changes.

• Ensures **ART**.

Methnani L., Aler Tubella A., Dignum V and <u>Theodorou A.</u> (2021). *Let Me Take Over: Variable Autonomy for Meaningful Human Control*. Front. Artif. Intell.

Mental Models

 With poor transparency, robots mislead us due to our lack of theory of mind for them.

 Misunderstanding leads to anxiety, mistrust, fear, and misuse/disuse. Introduces user self-doubt.

Transparency helps us build accurate mental models.

Technical solutions

Keeping the black box

Sometimes black boxes are inevitable.

 Some of the best-performing methods for pattern recognition, e.g. deep learning, are black boxes right now.

• Yet, we still need to audit our systems.

testing black boxes

Black-box testing for Responsible ai

- Al systems are software.
- White-box testing: the focus on testing is in the internal structures an application
- Black-box testing: ignores the 'inside' of a software, the focus on the inputs and outputs.





Black-box testing for Responsible ai

Testing a system adherence to ethical values.



Aler Tubella A., <u>Theodorou A.</u>, Dignum F., Dignum V. (2019). *Governance by Glass-box: implementing transparent moral bounds for AI behaviour*. International Joint Conference on Artificial Intelligence (IJCAI).

Formalising the interpretation stage

- We encode statements of the form: "A counts-as B in context C".
- It allows for verification in reasonable time.



Aler Tubella A., Dignum V. (2019). *The Glass Box Approach: Verifying Contextual Adherence to Values*. Workshop in AI Safety 2019

Formalising the OBSERVATION stage

 Tests can be established metrics in binary format.





But how consistently do we adhere to a policy?

Why not just ask the system itself? Explainability and Dialogue

Transparency through Information-Seeking Dialogues



Aler Tubella A., <u>Theodorou A.</u>, Nieves J.C. (2021). *Interrogating the Black Box: Transparency through Information-Seeking Dialogues*. AAMAS 2021

Agents-modelling-other-agents

- Inspired by agents-modellingagents literature.
- Agent communicate through information-seeking dialogues.
- Aggregation of information with argumentation semantics.



Aler Tubella A., <u>Theodorou A.</u>, Nieves J.C. (2021). *Interrogating the Black Box: Transparency through Information-Seeking Dialogues*. AAMAS 2021

The Awkward architecture

Socially-Aware Agents for Ethics By Design



The social behavior dilemma

- Fulfilment of societal goals requires altruistic or other prosocial behaviour.
- Fulfilment of individual goals may require freeriding, selfish, and anti-social behaviours.
- There is a need to alternate between different behaviours given the situation.

EU AI Ethics

european-approach-artificial-intelligence

European Artificial Intelligence Act: the right framework for trustworthy AI and innovation?



Home > Policies > A European approach to artificial intelligence

A European approach to artificial intelligence

The EU's approach to artificial intelligence centers on excellence and trust, aiming to boost research and industrial capacity while ensuring safety and fundamental rights.

The way we approach Artificial Intelligence (AI) will define the world we live in the future. To help building a resilient <u>Europe for the Digital Decade</u>, people and businesses should be able to enjoy the benefits of AI while feeling safe and protected.

The <u>European AI Strategy</u> aims at making the EU a world-class hub for AI and ensuring that AI is human contribution and tructure the Such an objective translates into the European approach to



Many Thanks to

RAI Group @ UmU



Prof. Dr. Virginia Dignum Professor of Responsible AI @vdignum @rirginia@cs.umu.se



Dr. Andrea Aler Tubella Senior Research Engineer andrea.aler@umu.se



Mattias Brännström Research Engineer mattias.brannstrom@umu.se

Dr. Juan Carlos Nieves Associate Professor ∑@jcnieves ∑juan.carlos.nieves@umu.se



Leila Methnani PhD Student ⊠eila.methnani@cs.umu.se

Other members: A. Brännström, *PhD student* Dr. P. Eriksson, *postdoc* A. Horned, *PhD student* Dr. A. Johnsson, *postdoc* K. Mendez, *PhD student* Dr. L. Vanhée, *Associate Professor*

Other collaborators

Dr. A. Antoniades – Builder. AI (UK) Prof. J.J. Bryson – Hertie Schol (DE) Prof. A. Colman – Princeton University (USA) Dr. A. Cortes – Barcelona Supercomputer Centre (ES) Dr. M. De Vos – Uni. Of Bath (UK) Prof. F. Dignum — UmU (SE) Prof. A. Kakas – Uni. Of Cyprus (CY) Dr. L. Michael – Open Uni. Cyprus (CY) C. Muller – ALLAI (NL) Rotisidis – Uni. Of Bath Α Dr. L. Sartori - Univ. of Bologna (IT) Dr. T. Scantamburlo -Univ. Of Venice (IT) B. Wilder – Harvard University (USA) H. Wilson – Uni. Of Bath (UK) Prof. A. Winfield – Bristol Robotics Lab/UWE (UK) Dr. R. H. Wortham – Uni. Of Bath (UK)

Current Funding Sources:















FACILITATE – AI Partners





Plovdiv University "Paisii Hilendarski"













