

Legal Landscapes of Al-Fairness

D6.1

This content is based on research conducted for the AEQUITAS report "Preliminary Social, Ethical and Legal Landscapes of Al-Fairness"

Social, Ethical and

Crucial to fair and trustworthy AI,

responsiveness by thinking

anticipation, reflexivity, inclusion, and

01. **Landscapes of**

AI-Fairness

www.aequitas-project.eu

Why?

Anticipating unfairnessin **ALL stages of the Al lifecycle**

To align the three

AEQUITAS Engines with the

policy elements and context

social, legal, ethical and

related to Al fairness.

Design phase

Development phase Deployment and use phase Dismantling phase the response (mitigation and AEQUITAS' overarching methodology focuses on the benefits of harm interpretation) stages. AEQUITAS will anticipation. Soliciting and involving the follow this approach to design and develop intended end-users as well as those an anticipatory experimentation eventually impacted by the technology is environment that will enable experimenting critical to anticipate potential harms both with, exploring and adjusting the fairness during the envisioning (model design) and levels of an Al tool.

How are we identifying the social landscape of Al fairness? Al Unfairness Manifestations

Database

Social Landscape

of Al Fairness

02.

are documented ✓ Relevant existing or upcoming policies related to the manifestation at hand Al Stakeholder Identification

We have developed a preliminary methodology for identifying

relevant stakeholders to involve in the design process of

Al-systems. By utilizing a combination of desk research, our

targeted user groups and stakeholders to be involved in the

own expert knowledge, and resources from previous projects, we have created a questionnaire that guides the selection of

We have created a database to track and analyze

✓ Sources of unfairness, such as data,

affected by the Al unfairness

fairness involved

Methodology (AISIM)

by the AI-system

Identification &

negatively and

level of impact

categorization of

positively affected

stakeholders and the

algorithms/models, or interpretations

manifestations of unfairness caused by AI in various domains.

The database allows for the categorization of information

✓ Technical details about the AI techniques used,

→ The ethical, legal, and social implications of AI

✓ Individuals or groups negatively and positively

→ The types of harm resulting from AI unfairness

training data, inputs, outputs, and interpretations

related to each manifestation of unfairness, including:

Affectees Decision makers Stakeholders affected Stakeholders that have Stakeholders that have

power over the

Al-system

development and

deployment of the

Identification &

power over the

deployment and

governance of the

Al-system and their

development,

categorization of

stakeholders with

Domain Experts

and Users

information that would

development of a fair

aid with the

Al-system

Identification &

categorization of

stakeholders that can

aid the development of

the Al-system and their

level of involvement

level of involvement



be found in the 7 Key Requirements

Technical robustness and safety

Helps to prevent errors or unintended

discrimination and other instances of

Includes ensuring that AI-systems are

reliable, accurate, secure, resilient to

attacks, and have a fallback plan.

unfair treatment.

consequences that could result in bias,

(KR) for Trustworthy AI of EGTAI:

Key Requirements

Privacy and Data Governance Helps protect individuals' personal information and prevent discrimination based on sensitive characteristics such as race, gender, or sexual orientation.

Diversity, Non-discrimination and Fairness Transparency Helps ensure that the system is fair, avoid the presence Allows individuals and of any unfair bias but also ensure that the system is

Social and Environmental Well-being

Ensures that the development and

broader social and environmental

development, reducing poverty and

inequality, promoting human rights, and

Legal landscape

of Al-Fairness

goals, such as sustainable

ensuring democracy.

deployment of Al-systems aligns with

organizations to understand

how AI-systems make

issues of bias or

decisions, to identify any

discrimination, and to hold

entities accountable in the

event of Al unfairness.

collected and used in a lawful and ethical manner.

Data should be free of socially constructed biases, inaccuracies, errors or mistakes,

data should be allowed to do so, and that any data used to train Al-systems is

that only duly qualified personnel with the competence and need to access individual's

Fairness related to:

Human agency and oversight

manipulated, deceived, herded or conditioned, or subject to any other

Requires that systems are regularly

reviewed and audited to ensure they are

unfair outcomes.

operating as intended.

Ensures that individuals are not unfairly

gender, abilities, or characteristics. Ensures diversity in the design of the Al-system through the participation of diverse stakeholders, including those who may be directly or indirectly affected by the system.

Ensures that individuals and organizations

can be held accountable for any issues of

bias or discrimination that may arise from

responsibility and liability, minimizing and

reporting negative impacts and trade-offs

Includes establishing clear lines of

accessible and can be used by all regardless of their age,

Accountability

the use of Al-systems.

that influence fairness.

and candidate selection. Non-discrimination, right to work, freedom of occupation, freedom of religion, right to privacy and family life, gender and racial equality, equal employment, data protection, equal treatment and opportunity, harassment prevention, fair and transparent digital HR practices, vigilance against compromising human dignity, human oversight, and compliance with GDPR. ✓ The proposed Al Act classifies Al systems intended to be used for recruitment or selection of natural persons, notably for advertising vacancies, screening or filtering applications, evaluating candidates

categorization as high risk.

often involves a type of scoring.

high risk.

in the course of interviews or tests, as high risk. Separately, the

Moreover, the proposal for the Al Act prohibits social scoring by

public authorities, which is seen as the evaluation or classification

of the trustworthiness of persons based on unrelated or irrelevant

social behaviour or personal characteristics, leading to detrimental

treatment of that person. This prohibition could be relevant when

using AI for HR, recruiting or candidate selection, especially as it

proposal for the Al Act classifies biometric identification and

Rights and the Treaty on the EU) and numerous EU Regulations and Directives hold notions of Al-Fairness relevant for the two use cases that deal with disadvantaged groups: (i) detection of child neglect and abuse, and (ii) access to education for disadvantaged students. Non-discrimination; human dignity; rights of the child; right to private and family life; right to life; right to preventive healthcare and medical treatment; data protection; racial equality; gender equality; freedom of religion Non-discrimination; human dignity; right to private and family life; protection; racial equality; gender equality; freedom of religion; rights of the elderly; rights of people with disability; right to education; freedom to choose an occupation; right to engage in work; equal treatment of qualifications. The proposed Al Act classifies Al systems intended to be used for the purpose of determining access or assigning natural persons to educational and vocational training institutions as high risk.

Development of the AEQUITAS AI Fairness-by-design Methodology Dismantling Scoping Design Development Deployment phase phase phase and use phase phase

We will develop an Al

that covers the

Fairness-by-Design Methodology

developments

Building Blocks Identification and assessment of manifestations of Al unfairness related to the Al-system and application domain at hand. ✓ Identification of the legal landscape of AI-Fairness related to the Al-system and application domain at hand. ✓ Identification of the ethical landscape of Al-Fairness using EGTAI related to the Al-system and the application domain at hand. ✓ Identification of relevant stakeholders using the Stakeholder Identification Methodology.

Legal notions of Al-Fairness in HR, **Recruiting and Candidate Selection** Several European and EU Treaties (including the European Convention on Human Rights, the EU Charter on Fundamental Rights and the Treaty on the EU) and numerous EU Regulations and Directives and self-regulatory instruments (e.g., social partner agreements) hold notions of Al-Fairness relevant for HR, recruiting

Legal notions of Al-Fairness in **Healthcare** Several European and EU Treaties (including the European Convention on Human Rights, the EU Charter on Fundamental Rights and the Treaty on the EU) and numerous EU Regulations and Directives hold notions of Al-Fairness relevant for the Healthcare domain. Non-discrimination; human dignity; right to private and family life; right to life; right to preventive healthcare and medical treatment; data protection; positive discrimination; racial equality; gender equality. The proposed Al Act classifies Al systems intended to be used as

safety components of a medical device, as regulated in the

Legal notions of Al-Fairness

regarding disadvantaged groups

Several European and EU Treaties (including the European

Separately, the proposal for the Al Act classifies biometric

identification and categorization as high risk.

Once the Al Act is adopted,

systems must comply with a

Convention on Human Rights, the EU Charter on Fundamental

Medical Devices Regulation as **high risk**. Separately, the proposal for

the Al Act classifies biometric identification and categorization as

large set of requirements (iii) technical documentation and before they can be put on the record keeping EU internal market. The AI Act (iv) transparency and provision of proposal categorizes the information to users requirements into the (v) measures to ensure overarching categories of: human oversight (vi) accuracy, robustness, and cybersecurity.

(i) risk management

(ii) data and data governance

Al Act

Al Treaty of the Council of Europe

What's next? **EU AI Regulatory Framework Future policy**

Identification and assessment of ethical, legal and social AI-Fairness elements in collaboration with stakeholders



www.aequitas-project.eu

@aequitasEU

Follow us

 \bigcirc in





Funded by

the European Union



Views and opinions expressed are however those of the authors

only and do not necessarily reflect those of the European Union.

Neither the European Union nor the granting authority can be

UCC

Contact us

info@aequitas-project.eu

Be the first to get exclusive project updates

by subscribing tour newsletter!

EUROCADRES

PHILIPS