

# Trustworthy AI And the European AI act

**G.Landais** 

# Trustworthy AI



### How can I provide a trustworthy AI?

- → Am I subject to any rules?
- → Are there ethic ? good practice guideline ?

### **European legislation**

- European AI act https://www.europarl.europa.eu/thinktank/en/document/EPRS\_BRI(2021)698792 https://artificialintelligenceact.eu/
- Ethics Guidelines for Trustworthy AI https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai

# European Artificial intelligence act



EU AI Act: first regulation on artificial intelligence The use of artificial intelligence in the EU is regulated by the AI Act, the world's first comprehensive AI law. Find out how it protects you.

Published: 08-06-2023, Last updated: 19-02-2025



### What is it about?

- Regulation
- Evaluate AI risk
   (Astronomy and ChatGPT are not considered to be a High Risk)
- Encourage AI development with sandbox environment
- · EU act compliance guidance

#### Who is impacted?

Organisation from Europe or outside Europe that deploy, distribute or provide AI, in particular High-risk AI systems

### Calendar

- 2024/08/01 : begining of regulation.
- 2025/02/02 : Prohibition of AI presenting unacceptable risks.
- 2025/08/02 : establish rules for general purposes

### Consequences



- High-risk AI systems requires CE marking
- GPAI (General Purpose AI) needs transparency subject to a good practice guideline

The regulation is not applicable to any activities related to AI research, testing, and development before it is marketed or put into operation.

# European Artificial intelligence act



### Risks levels

#### **Prohibited AI practices**

- · Prohibes subliminal or manipulative techniques
- Protect Fundamental laws
- Biometric categorisation or identification in public spaces
- Al evaluating or classifying individuals or groups based on social
- ...

#### **High-risk AI systems**

- Critical infrastructure
- Education and vocational training
- Access to essential private and public services
- Law enforcement
- Migration and border control management
- Administration of justice and democratic processes
- ..
- → Needs to run a conformity assessment procedure

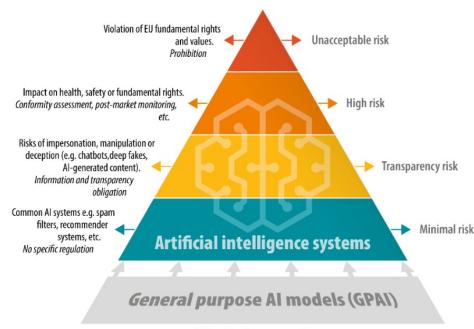
#### **Transparency risk**

- Users must be made aware that they interact with chatbots.
- (deployers and providers of) Al must disclose that the content has been artificially generated

#### Minimal risks

(Example SPAM filter) is out of requirements

113 articles (+annexes) gathered into 13 main topics https://artificialintelligenceact.eu/ai-act-explorer/



GPAI models - Transparency requirements

GPAI with systemic risks - Transparency requirements, risk assessment and mitigation

### **European Artificial intelligence act**



### General-purpose AI (GPAI), what it is?

https://artificialintelligenceact.eu/wp-content/uploads/2022/05/General-Purpose-Al-and-the-Al-Act.pdf

#### **Characteristics:**

- Wide range of possible use, both intented and unintented
- Applied to many different tasks without providing fine tuning
- · Large scale : memory, power and lot of data
- Often GPAI are building-bloks reuse by other AI systems

Large variety of usage: natural language, used in medicine and healthcare, finance, life sciences, programming ...

### **GPAI** obligations

GPAI system transparency requirements:
 maintain tech. Doc. Including content used for training the AI model

(Opensource have less obligations)

- Put a policy that respects the Union copyright law
- Cooperation with authorities keeping track of incident (eg: violation of fundamental rights)
- GPAI model providers will be able to rely on codes of practice

see Article 53: Obligations for Providers of General-Purpose Al Models https://artificialintelligenceact.eu/article/53/

### **Ethics Guidelines for Trustworthy AI**



### https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai

→ Dedicated for any AI provided in Europe or consume in Europe

### The 3 components of a Trustworthy Al

- Lawfull AI: ensuring compliance with all applicable laws and regulations
- Ethical AI: including fundamental rights, ethical principles
- Robust AI: include tech. + societal perpsective robustness technical robustness includes the whole AI life-cycle human respect and representativeness

### The document is made of 3 parts

- Ennounce ethical principles repect of human, acknowledge benefits and AI risks
   "even with good intentions, unintentional harm can occurred"
- List 7 requirements
- Accessment list: Questions to help AI provider and covering the 7 requirements



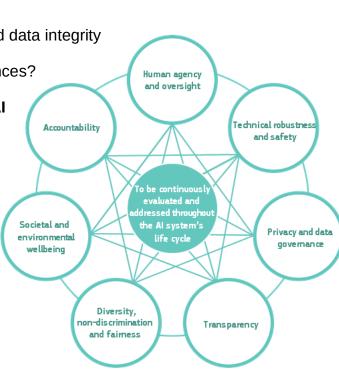
# **Ethics Guidelines for Trustworthy AI**



### 7 requirements for Trustworthy AI

- Human agency and oversight ← technology must be in favor of human (and supervised by human)
  - Can you describe the **level of human control** or involvement?
- Robustness and safety ← safe and reproducibility aspects
  - Did you consider different types and natures of vulnerabilities, such as **data pollution**, physical infrastructure, **cyber-attacks**?
- Privacy and data governance data protection, rights, copyrights and data integrity
  - Did you assess who can access users' data, and under what circumstances?
- Transparency 
   — Human must be aware of their interaction with AI
  - Did you ensure an explanation as to why the system took a certain choice resulting in a certain outcome that all users can understand?
- Diversity, non-discrimination and fairness
  - ← ex: no racial, gender or handicap discrimination
  - Did you consider diversity and representativeness of users in the data?
- Societal and environmental ← ecology aspect
- Accountability ← require responsibility and transparency

IVOA Goerlitz, 2025 - DCP/K



### **Ethics Guidelines for Trustworthy AI**



### **EU AI Act Compliance Checker**

https://artificialintelligenceact.eu/assessment/

Example: I'm a deployer of a AI model in Open source that generates images for educational and vocational training.

### Results:

#### Your results

Providers must submit notification to NCA (National Competent Authority)

If a provider considers their Al system to not pose a significant risk (see <u>Article 6</u> point 2a) they must register their system in the EU database before that system is placed on the market or put into service (see <u>Article 49</u> point 2).

They must also document their assessment and provide this documentation to the National Competent Authorities (NCA) upon request (see Article 6 point 4).

If a market surveillance authority finds that the AI system has been misclassified (see <u>Article 80</u>), your system would be subject to the 'high-risk' obligations described in <u>Chapter III Section 2</u> and you may be subject to fines under Article 99.

# My Al system 'output is in used in the EU

You must take measures to ensure a sufficient level of Al literacy for your staff (and other people dealing with the operation and use of Al systems on their behalf), taking into account their technical knowledge, experience, education and training and the context the Al systems are to be used in, according to <u>Article 4</u>.

# I'm placing on the market General purpose Al models

#### Providers must submit notification to NCA

If a provider considers their AI system to not pose a significant risk (see <u>Article 6</u> point 2a) they must register their system in the EU database before that system is placed on the market or put into service (see <u>Article 49</u> point 2).

They must also document their assessment and provide this documentation to the National Competent Authorities (NCA) upon request (see Article 6 point 4).

If a market surveillance authority finds that the AI system has been misclassified (see <u>Article 80</u>), your system would be subject to the 'high-risk' obligations described in <u>Chapter III Section 2</u> and you may be subject to fines under <u>Article 99</u>.

#### General Purpose AI model obligations

You need to follow these obligations for Providers of General Purpose AI (GPAI) models under <u>Article 53</u>. In summary, you must:

- Create and keep technical documentation for the Al model, and make it available to the Al Office upon request.
- Create and keep documentation for providers integrating Al models, balancing transparency and protection
  of IP
- · Put in place a policy to respect Union copyright law.
- Publish a publicly available summary of Al model training data according to a template provided by the Al
  Office

Also, consider whether the GPAI is used as, or a component of, an AI system. If so, obligations on high risk AI systems may apply directly or indirectly under <u>Recital 85</u>.

#### Al Literacy obligations

You must take measures to ensure a sufficient level of Al literacy for your staff (and other people dealing with the operation and use of Al systems on their behalf), taking into account their technical knowledge, experience, education and training and the context the Al systems are to be used in, according to Article 4.

# Conclusion



### What I retain?

- Respect of fundamental laws
- Recognize the benefit and the risks
  - Risk level, the logo CE
- Different level of risks and obligation
  - Deploying AI may require initiating an endorsement process
- An ethic for trustworthy Al
- Science is not impacted
  - Respects of copyrights (transparency)