# AI Cards: Towards an Applied Framework for Machine-Readable AI and Risk Documentation Inspired by the EU AI Act

Delaram Golpayegani*, Isabelle Hupont^, Cecilia Panigutti^
Harshvardhan J. Pandit**, Sven Schade^, Declan O'Sullivan*, Dave Lewis*

* ADAPT Centre, Trinity College Dublin, Dublin, Ireland
^European Commission, Joint Research Centre (JRC), Ispra, Italy
**ADAPT Centre, Dublin City University, Dublin, Ireland

Annual Privacy Forum 2024
4 September 2024

HOST INSTITUTION
Trinity College Dublin
Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

HOST INSTITUTION
DCU
Ollscoil Chathair
Bhaile Átha Cliath
Dublin City University

PARTNER INSTITUTIONS

T DUBLIN
OLLSCOIL TEICNEOLAÍOCHTA
BHAILE ÁTHA CLIATH
TECHNOLOGICAL
UNIVERSITY DUBLIN

UCD DUBLIN
University College Dublin
An Coláiste Ollscoile, Baile Átha Cliath
Ireland's Global University

MTU
Ollscoil Teicneolaíochta na Mumhan
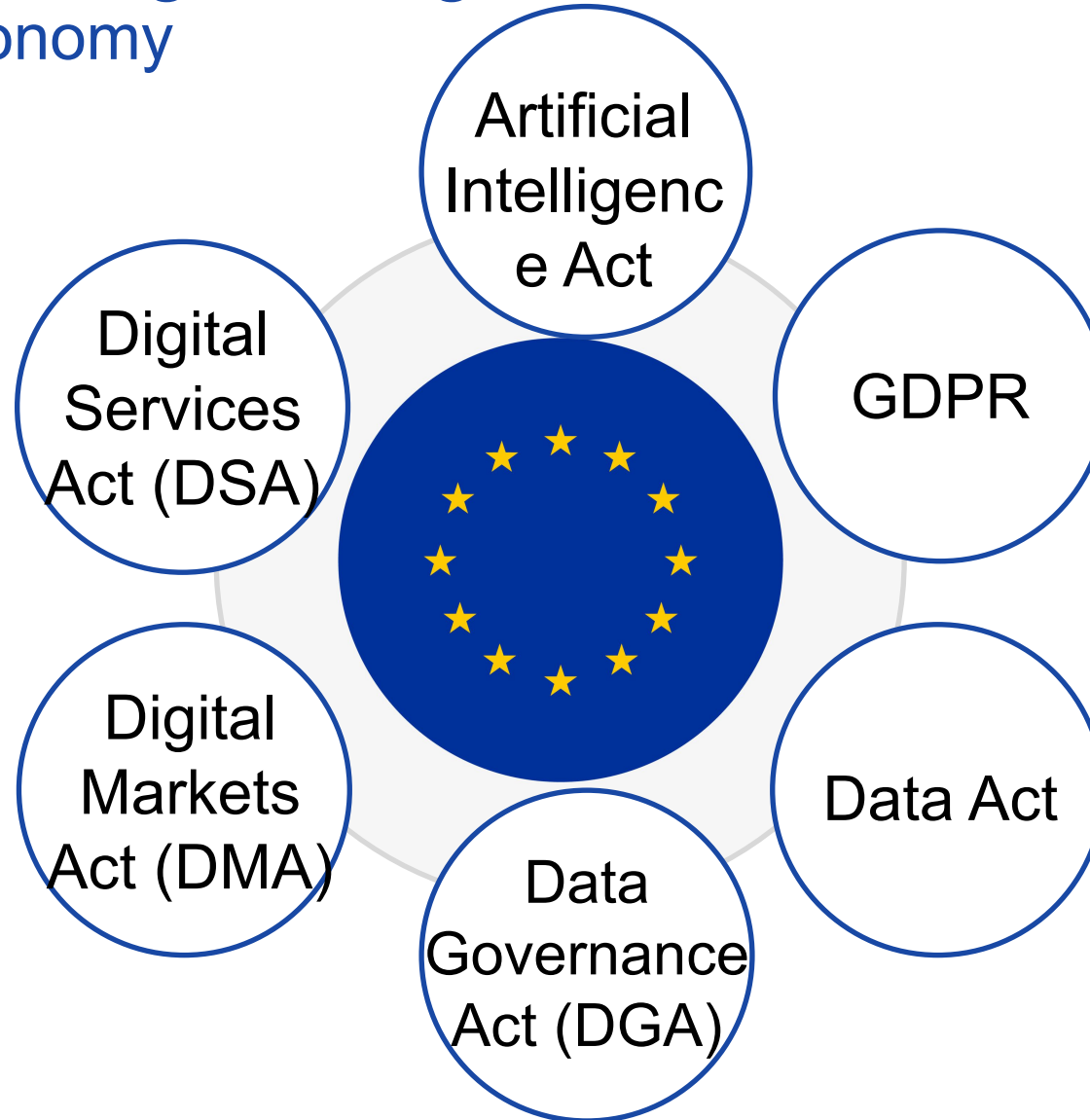Munster Technological University

TUS

Maynooth University
National University
of Ireland Maynooth

OÉ Gaillimh
NUI Galway

# The Big 5+1 EU Digital Regulations
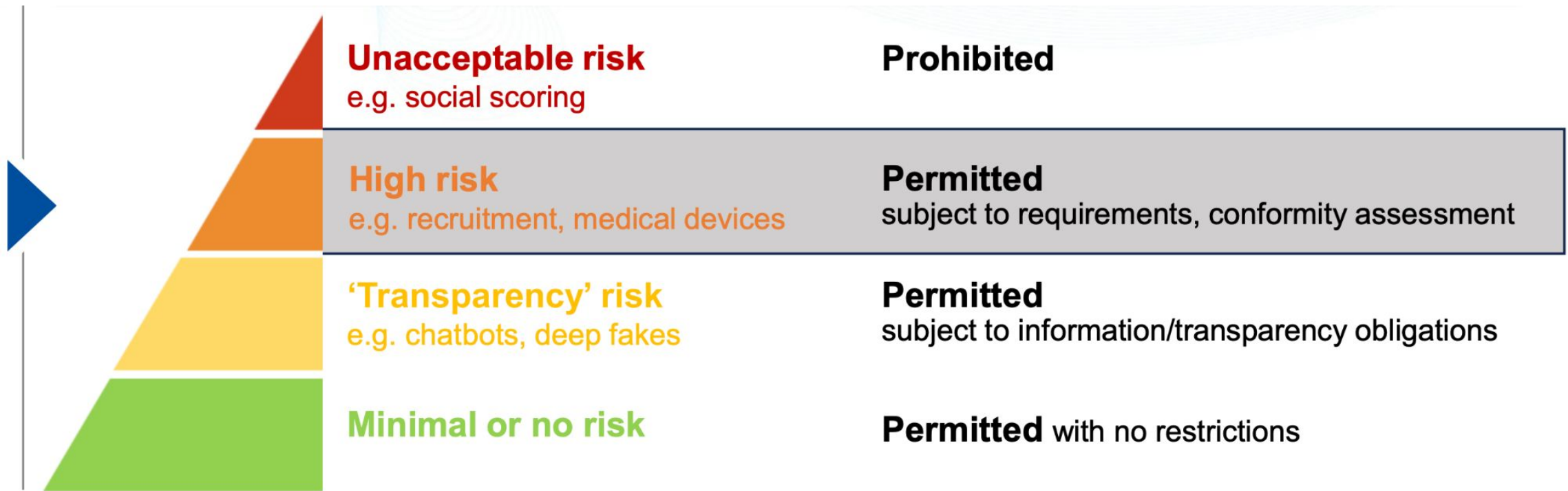## for Data and AI Economy

# The EU AI Act

New Rules for

- **AI Systems**

- **GPAI Models** [General Purpose AI]

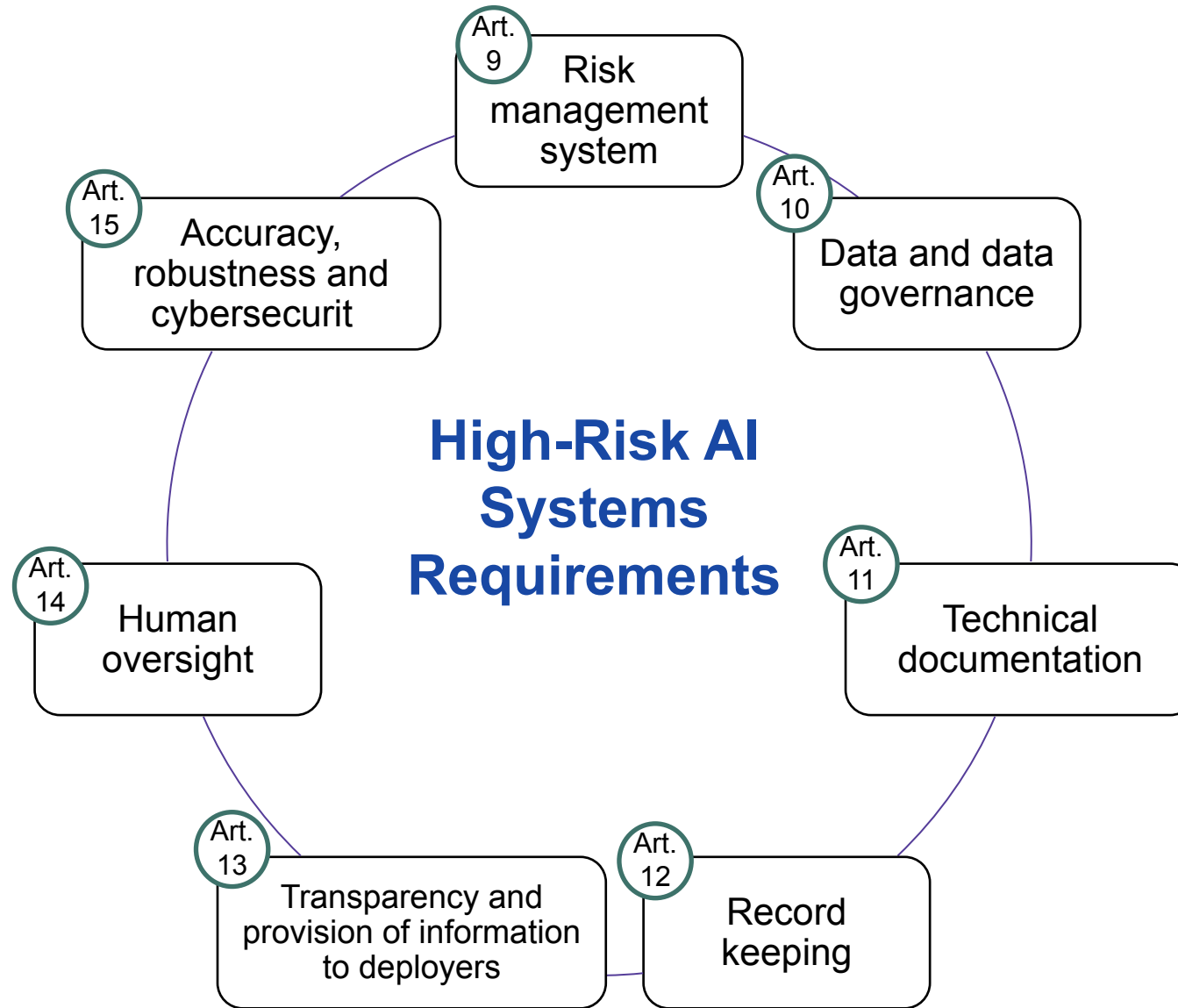Promotes human-centric & trustworthy AI

Protects against harmful effects of AI on

- **Health**

- **Safety**

- **Fundamental Rights**

# AI Systems Risk-Based Classification



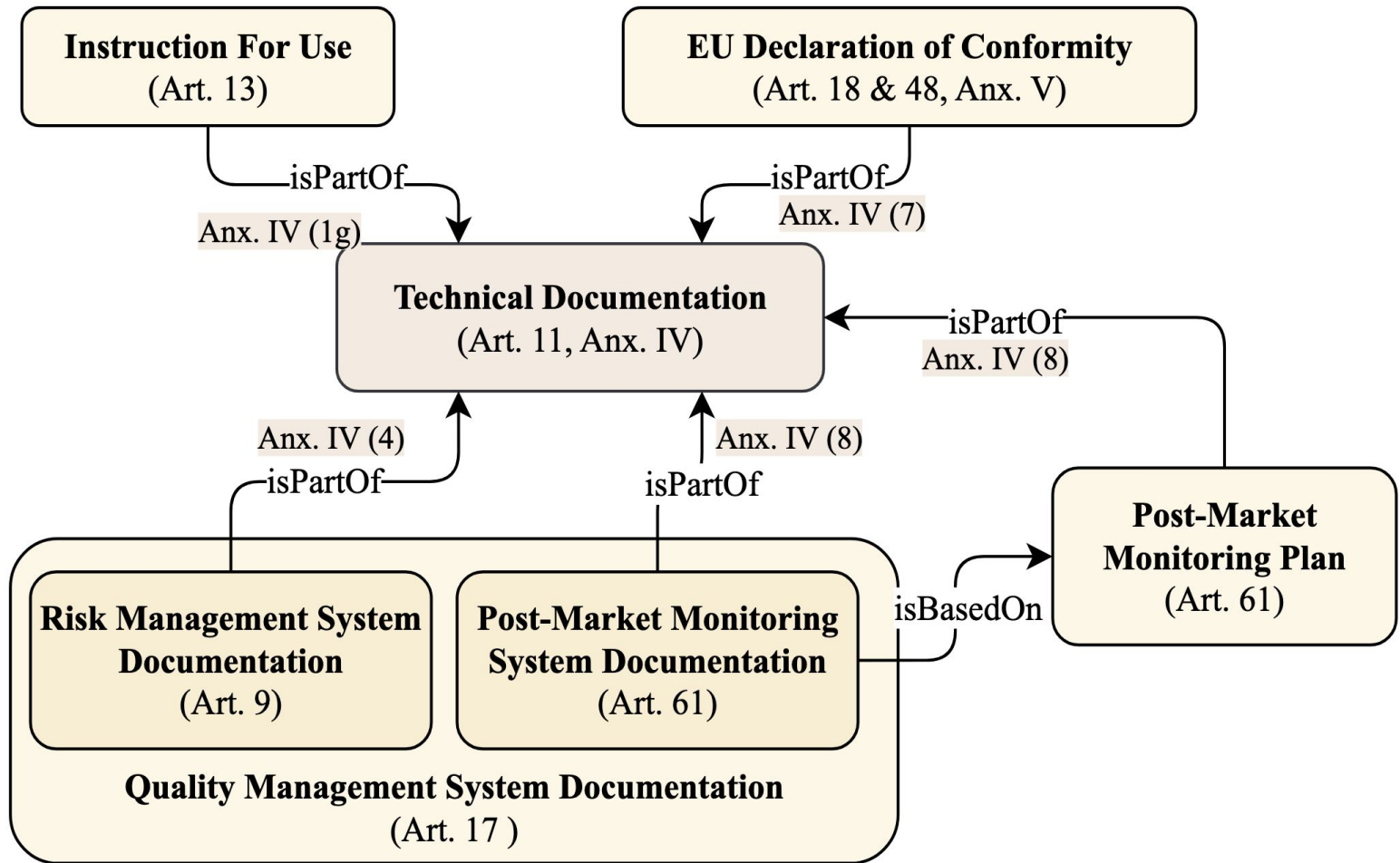| | |
|---|---|
| **Unacceptable risk** <br> e.g. social scoring | **Prohibited** |
| **High risk** <br> e.g. recruitment, medical devices | **Permitted** <br> subject to requirements, conformity assessment |
| **'Transparency' risk** <br> e.g. chatbots, deep fakes | **Permitted** <br> subject to information/transparency obligations |
| **Minimal or no risk** | **Permitted** with no restrictions |

From the EU AI Office webinar on risk management in the AI Act and related standards, 30 May 2024
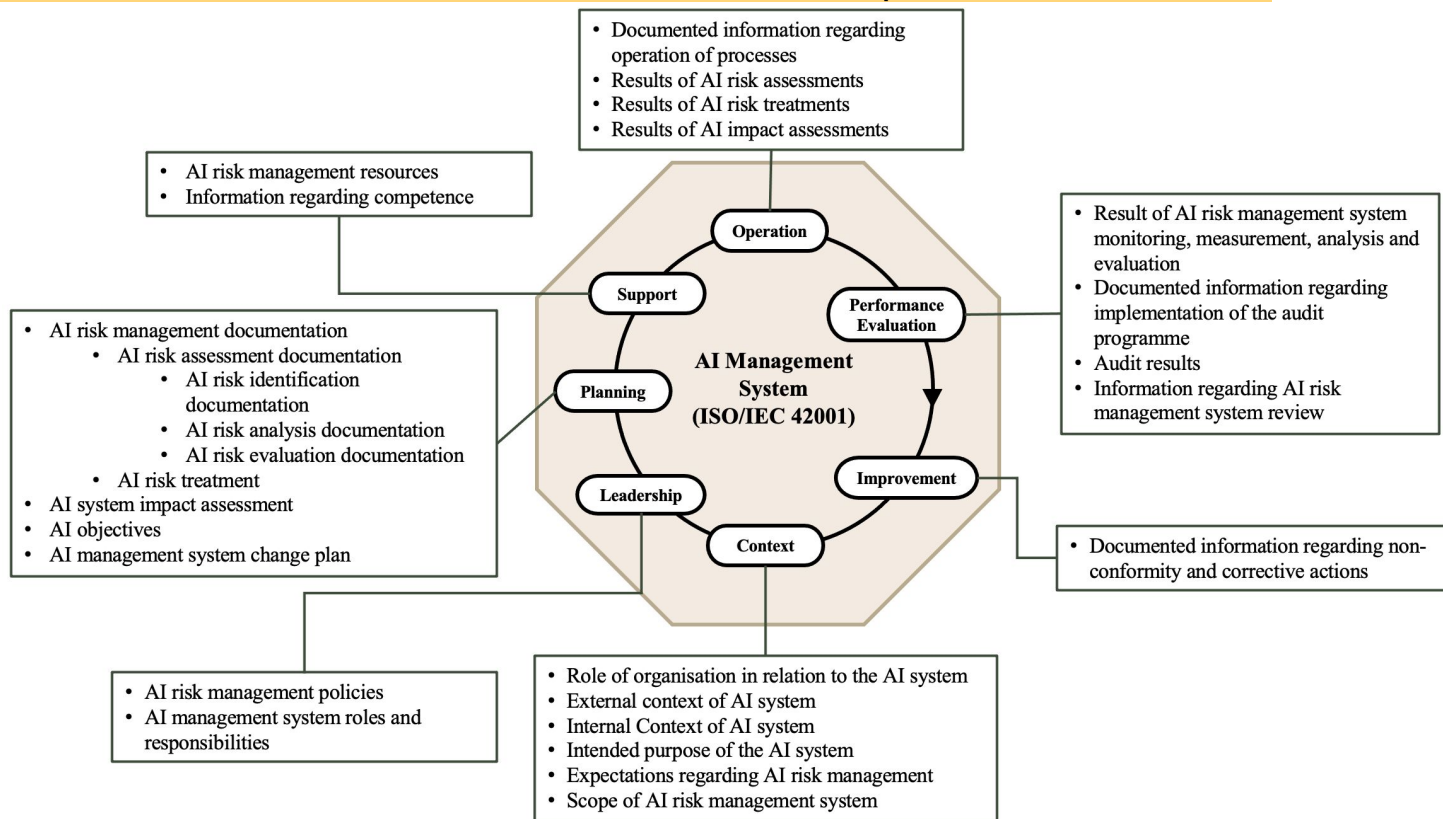
# Documentation Requirements

- Documentation →
  Transparency→
  Trustworthy AI

- "The technical
  documentation shall be
  drawn up in such a way
  as to demonstrate that
  the high-risk AI system
  complies with the
  requirements" (Art. 11)

# Risk Management System Documentation - available today

ISO/IEC 42001 on AI management system

ISO/IEC 23894 on AI risk management

<span style="background-color:#ffd966">Note that these standards are not sufficient for compliance with the AI Act</span>



**Types of information**

1. Information about the <span style="background-color:#ffa500">context of the AI system and the organisation</span>

2. Details of the risk management <span style="background-color:#ffa500">system</span> in place

3. Risk management <span style="background-color:#ffa500">processes</span>

4. <span style="background-color:#ffa500">Results</span> of AI risk management

# AI Cards (I)

# AI Cards (II)

1. General Information about the system

2. Intended use of the AI system using 6 concepts

3. Information about the incorporating components

4. Information about processing of data (including info about legal basis and source of data)

5. Involvement of humans and level of automation



**1. General Information**

Version
Modality
AI Technique(s)
Provider(s)
Developer(s)

**2. Intended Use**

Domain
Purpose
Capability
Deployer
AI Subject
Locality of use

**3. Key Components**

Input (from user)

Component #1
ID            D

Component #2
ID            M

Component #3
ID            D

Component #4
ID            S

Component #5
ID            GPAI

Output (to user)

**D**ataset
**M**odel
**S**ystem
**General P**urpose

Hardware Platform

**4. Data Processing**

| Processing | Legal basis | Data | Data Source |
|---|---|---|---|
| Processing #1 | | Data #1 | |
| | | | |
| | | Data #N | |
| Processing #N | | | |
| | | | |

**5. Human Involvement**

Level of Automation

| Involved Entity | Intended | Active | Informed | Control over output |
|---|---|---|---|---|
| AI Subject#1 | ✓ | ✗ | ✓ | |
| AI Subject#N | ✗ | ✓ | ✗ | |
| End-user#1 | ✗ | ✗ | ✓ | |
| End-user#N | ✗ | ✓ | ✗ | |

# AI Cards (III)

**6. High-level summary of risk management**

**7. Illustration of key qualities of the AI system**

**8. List of pre-determined changes**

**8. Regulation & Certification information**

### 6. Risk Profile

| Impact on ↓ | Risk | | | Measures | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Likeli. | Severity | Residual | Org. | Tech. | Monit. | Secur. | Transp. | Log. |
| Health & Safety | High | V. High | Med. | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Fundamental Rights | V. High | High | High | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ |
| Society | Med. | Med. | Low | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Environment | Low | Low | Low | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ |

### 7. Quality

Quality #1, Quality #2, Quality #3, Quality #4, Quality #5, Quality #6, Quality #7

### 8. Pre-determined Changes

| Changed Entity | Change Frequency | Purpose of Change |
|---|---|---|
| Data | | |
| Model | | |
| ... | | |

### 9. Compliance & Certification

| Regulations | |
|---|---|
| Standards | |
| Codes of conduct | |

# Example:
# An AI-Based Student Proctoring System

Proctify is intended to be used in the education domain, for detecting suspicious behaviour of students during online exams in universities. Facial behaviour analysis and video analysis are used for detecting suspicious behaviour

```
ex:proctify
    airo:isAppliedWithinDomain ex:education ;
    airo:hasPurpose ex:detecting_suspicious_bahviour_during_online_exam
    airo:hasCapability ex:facial_behaviour_analysis ;
    airo:hasCapability ex:video_analysis ;
    airo:isUsedBy ex:university ;
    airo:hasAISubject ex:student ;
```

https://delaramglp.github.io/aicards/example/



AI Cards: Proctify

| Card's Version | 1.2.3 |
| Card's Date (Issued) | 2024-04-23 |
| Card's Language | Eng |
| Card's Publisher | AIEduX |
| Contact Info | proctify@aiedux.org |

https://raw.githubusercontent.com/DelaramGlp/airo/main/usecase/proctify.ttl

**1. General Information**
Version: 1.2
Modality: Software
AI Technique(s): ML>>ANN>>Deep learning
Provider(s): AIEduX
Developer(s): AIEduX

**2. Intended Use**
Domain: Education
Purpose: Detecting suspicious behaviour during online exam
Capability: Facial behaviour analysis, video analysis
Deployer: University
AI Subject: Students
Locality of Use: Educational institution in EU

**3. Key Components**
Facial video → Facial Analysis Toolkit 3.3.2 (tinyurl.com/3wnyxyun) S → SusBehavedModel 1.1.2 (tinyurl.com/2hnth6bb) M → SusBehavedDataset 2.0.1 (tinyurl.com/db4whuw9) D → Suspicious behaviour alarm
Dataset, Model, System, General Purpose

**4. Data Processing**

| Processing | Legal basis | Data | Data Source |
|---|---|---|---|
| Processing of input video | Informed consent | Facial>> Biometrics | User input |
| Behaviour analysis (ML model) | Informed consent | Facial>> Biometrics | SusBehaved dataset contributers |

**5. Human Involvement**
Level of Automation: Partial automation

| Involved Entity | Intended | Active | Informed | Control over output |
|---|---|---|---|---|
| Student | ✓ | ✓ | ✓ | ex-post challenge |
| Occupant (of the room) | ✗ | ✗ | ✗ | No opt-out |
| Instructor | ✓ | ✓ | ✓ | Correct |

**6. Risk Profile**

| Impact on ↓ | Risk | | | Measures | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Likeli. | Severity | Residual | Org. | Tech. | Monit. | Secur. | Transp. | Log. |
| Health & Safety | Med. | V. High | Low | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ |
| Fundamental Rights | High | V. High | Low | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Society | Low | Med. | Med. | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ |
| Environment | Low | Low | Low | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |

**7. Quality**
Accuracy, Cybersecurity, Robustness, Fairness, Functional adaptability, Usability, Explainability

**8. Pre-determined Changes**

| Changed Entity | Frequency | Purpose |
|---|---|---|
| Susbehaved model | 2 Month | Improve performance |
| Mitigation measures | 2 Week | Mitigate newly identified risks |

**9. Compliance & Certification**

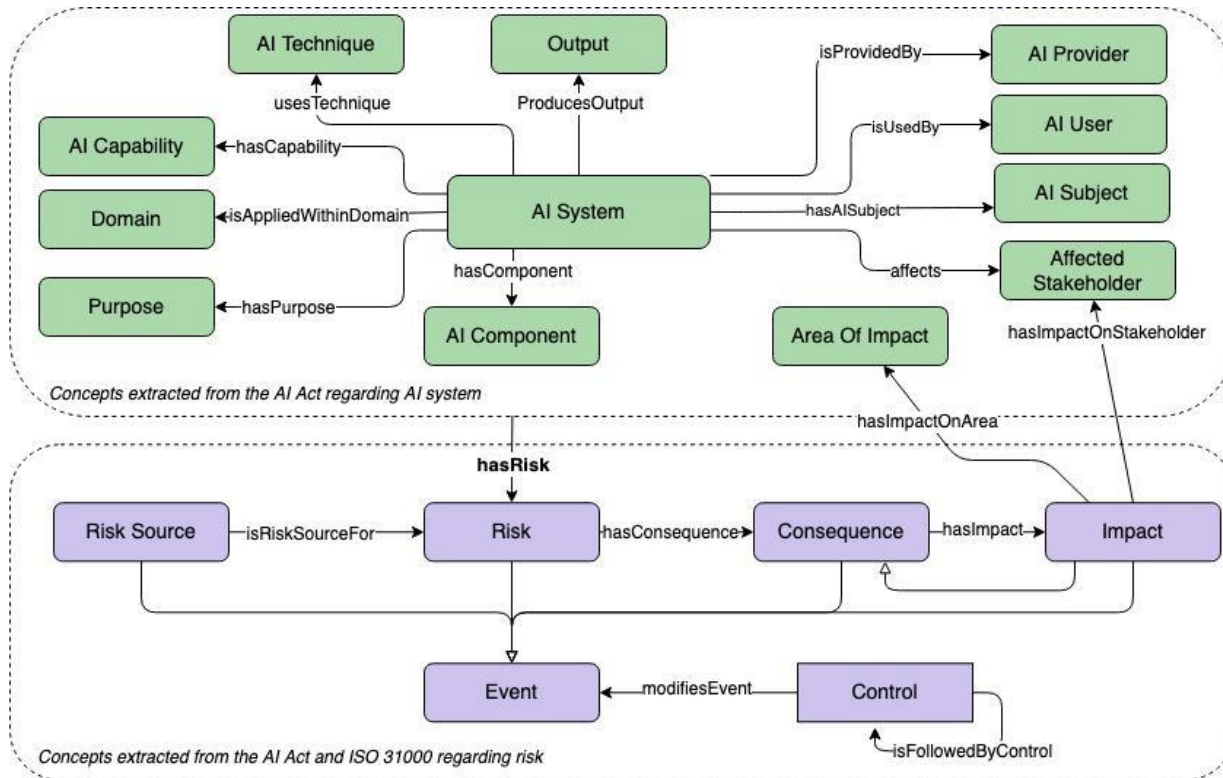| Regulations | [EU, GDPR] [IE, DPA] |
|---|---|
| Standards | [ISO/IEC 27001:2022] |
| Codes of conduct | [EU, use of AI and data in teaching and learning for educators] |

# Benefits of the Machine-Readable Representation

- Consistency

- Interoperability

- Integration with other documents, e.g. HuggingFace's Model Cards

- Automated generation of AI Cards (using SPARQL queries)

- Querying to support compliance checking

# https://w3id.org/airo
## AI Risk Ontology (AIRO)

# https://w3id.org/vair
## Vocabulary of AI Risks (VAIR)



Concepts extracted from the AI Act regarding AI system

Concepts extracted from the AI Act and ISO 31000 regarding risk

| 4. | Purposes | 7. | AI Capabilities |
|---|---|---|---|
| 4.1 | Remote Identification Of People | 7.1 | Biometric Identification |
| 4.2 | Content Generation | 7.2 | RemoteBiometricIdentification |
| 4.3 | Generating Audio Content | 7.3 | Personality Traits Analysis |
| 4.4 | Generating Image Content | 7.4 | Emotion Recognition |
| 4.5 | Generating Video Content | 7.5 | Profiling |
| 4.6 | Knowledge Reasoning | 7.6 | Face Recognition |
| 4.7 | Applying The Law To Facts | 7.7 | Computer Vision |
| 4.8 | Interpreting Law | 7.8 | Image Recognition |
| 4.9 | Interpreting Facts | 7.9 | Automatic Summarisation |
| 4.10 | Decision Making | 7.10 | Dialogue Management |
| 4.11 | Examining Application | 7.11 | Information Retrieval |
| 4.12 | Examining Asylum Application | 7.12 | Machine Translation |
| 4.13 | Examining Migration Related Complaints | 7.13 | Named Entity Recognition |
| 4.14 | Examining Residence Permits Application | 7.14 | Natural Language Generation |
| 4.15 | Examining Visa Application | 7.15 | Part Of Speech Tagging |
| 4.16 | Assessment | 7.16 | Question Answering |
| 4.17 | Assessing Past Criminal Behaviour | 7.17 | Relationship Extraction |
| 4.18 | Assessing Admission Test | 7.18 | Speech Recognition |
| 4.19 | Assigning People To Educational Institutions | 7.19 | Speech Synthesis |
| 4.20 | Determining Access To Education | 7.20 | Pattern Recognition |
| 4.21 | Determining Admission To Educationall nstitutions | 7.21 | Action Recognition |
| 4.22 | Assessing Student | 7.22 | Gesture Recognition |
| 4.23 | Evaluating Learning Outcomes | 7.23 | Object Recognition |
| 4.24 | Recruiting | 7.24 | Music Information Retrieval |
| | | 7.25 | Sound Event Recognition |
| | | 7.26 | Sound Synthesis |
| | | 7.27 | Sound Source Separation |
| | | 7.28 | Speaker Recognition |
| | | 7.29 | Lie Detection |
| | | 7.30 | Sentiment Analysis |

AIRO and VAIR are going to be integrated with DPV (https://w3id.org/dpv)

# Future Work

- Alignment of AI Cards with documentation and reporting requirements of EU digital regulations, including the GDPR, DSA, Interoperability Act, & DGA
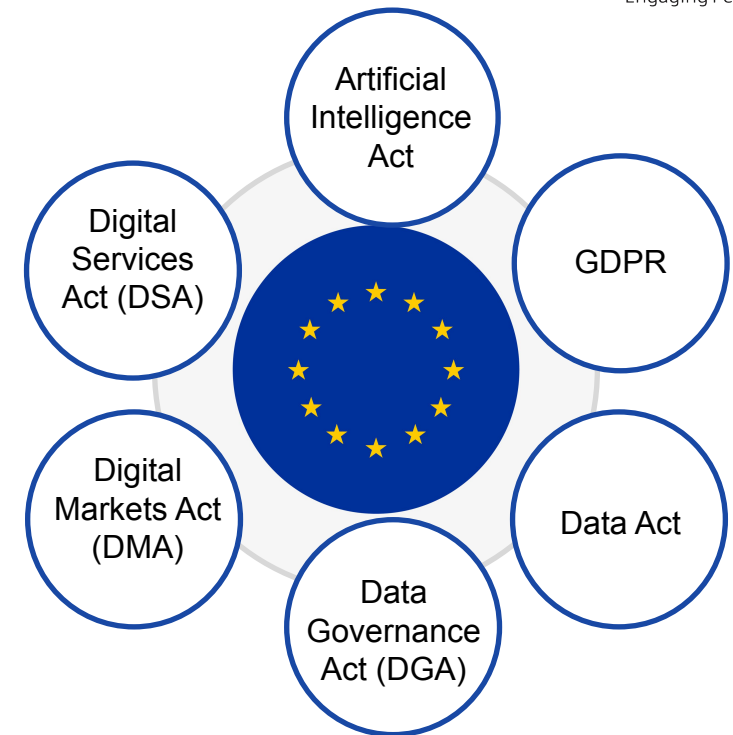
- Providing input to
  - Data Privacy Vocabulary (DPV)
  - ISO/IEC 42005 on AI Impact Assessment
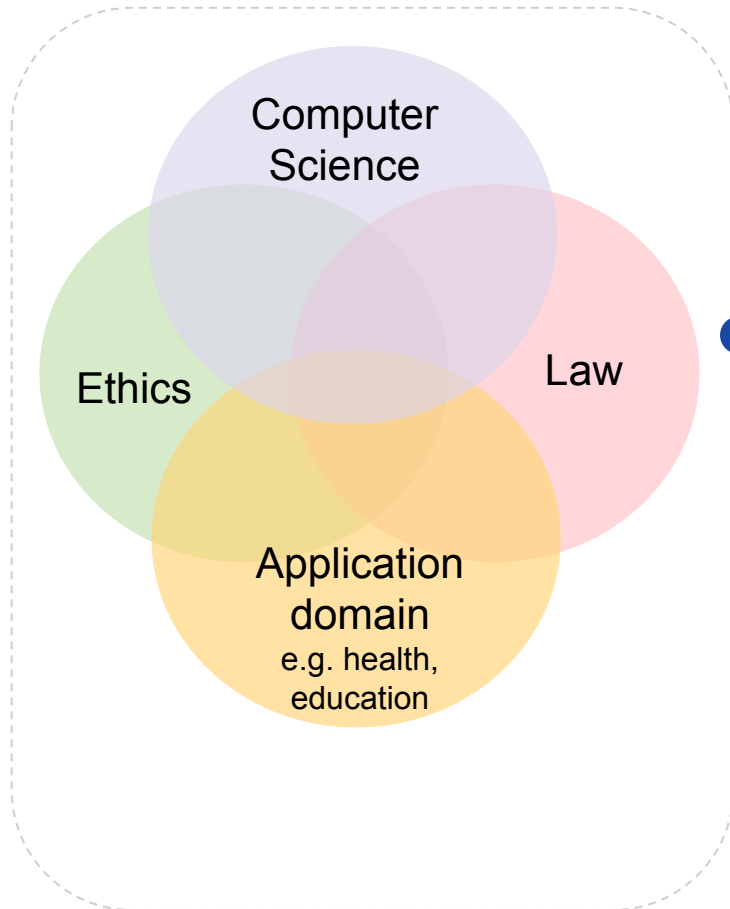  - CEN-CENELEC Catalogue of AI Risks

# Navigating in the New European Regulatory Environment (5+1 regulations)



## RegTech Solutions

- Risk and Operations Tech
- Compliance and Reporting Tech
  - To support compliance and conformity assessment
- Supervisory Tech
- Ethical Tech

## Legal AI ≠ Ethical AI

- Fundamental Rights Impact Assessment  (FRIA)
- General Purpose AI Models obligations
  - Regulatory sandboxes
  - High-risk AI use cases

- Overlaps in regulations (e.g. AI Act & GDPR)
  - Alignments/mappings with standards (e.g. NITS AI RMF)

- AI quality attributes
- AI testing
- AI/Gen-AI risk assessment

- Public awareness and engagement
- AI literacy
- Right to be informed

Computer Science

Ethics

Law

Application domain
e.g. health, education

Impact on

- EC's policies & guidelines
- International and European standards
- Codes of practice in different domains

**Safe, trustworthy, & green AI**

# AI Cards: Towards an Applied Framework for Machine-Readable AI and Risk Documentation Inspired by the EU AI Act

Delaram Golpayegani, Isabelle Hupont, Cecilia Panigutti
Harshvardhan J. Pandit, Sven Schade, Declan O'Sullivan, Dave Lewis

delaram.golpayegani@adaptcentre.ie

16