# Protect

# AIRO

## an Ontology for Representing AI Risks based on the Proposed EU AI Act and ISO Risk Management Standards

Delaram Golpayegani, Harshvardhan J. Pandit, Dave Lewis

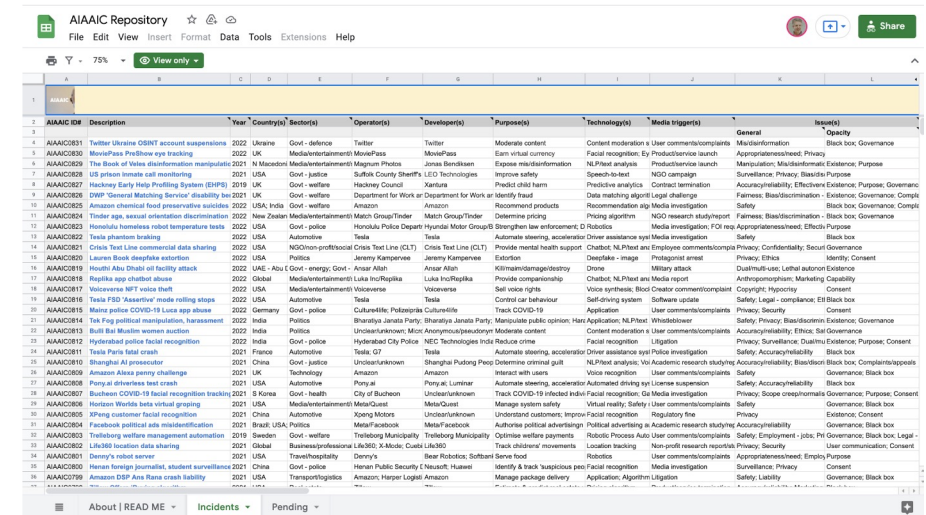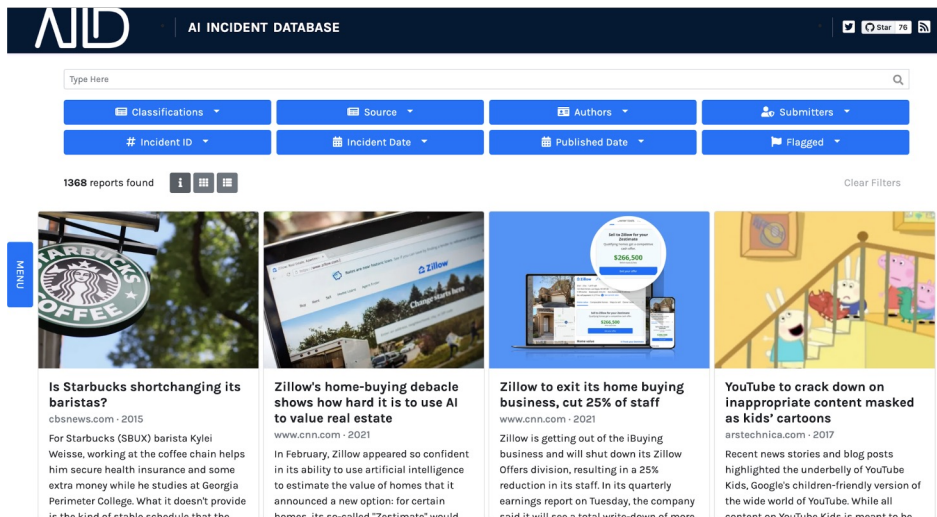ADAPT Centre, Trinity College Dublin, Ireland

sgolpays@tcd.ie

Trinity College Dublin
The University of Dublin

ADAPT
Engaging Content
Engaging People

IRISH RESEARCH COUNCIL
An Chomhairle um Thaighde in Éirinn

# AI Risks



**AIID (AI Incident Database)**
https://incidentdatabase.ai

**AIAAIC Repository**
https://www.aiaaic.org/aiaaic-repository

Icons from https://www.flaticon.com/

# Efforts Addressing AI Risks



**Regulations**

Promote trustworthy AI

Request for harmonised standards
(AI Act, Art. 40)

Highlight the need for AI
regulation

**AI Risks**

Provide technical solutions

**Trustworthy AI Guidelines**

Show gaps in standards

Provide technical solutions

**Standards**

3

AIRO: Ontology for representing AI Risks | Delaram Golpayegani et al. | SEMANTiCS 2022 | contact:sgolpays@tcd.ie | https://w3id.org/AIRO

# AI Act Risk Pyramid

**Unacceptable Risk**

Prohibited

**High Risk**

Requirements for high-risk AI systems

High Risk to health, safety, and fundamental rights of people:
1) Product or safety component of a product covered by Annex II
2) AI system used in Annex III areas

**Limited Risk**

Transparency obligations for certain AI systems

**Minimal Risk**

Codes of conduct

**RQ1**
What is the information required to determine whether an AI system is 'high-risk' as per the AI Act?
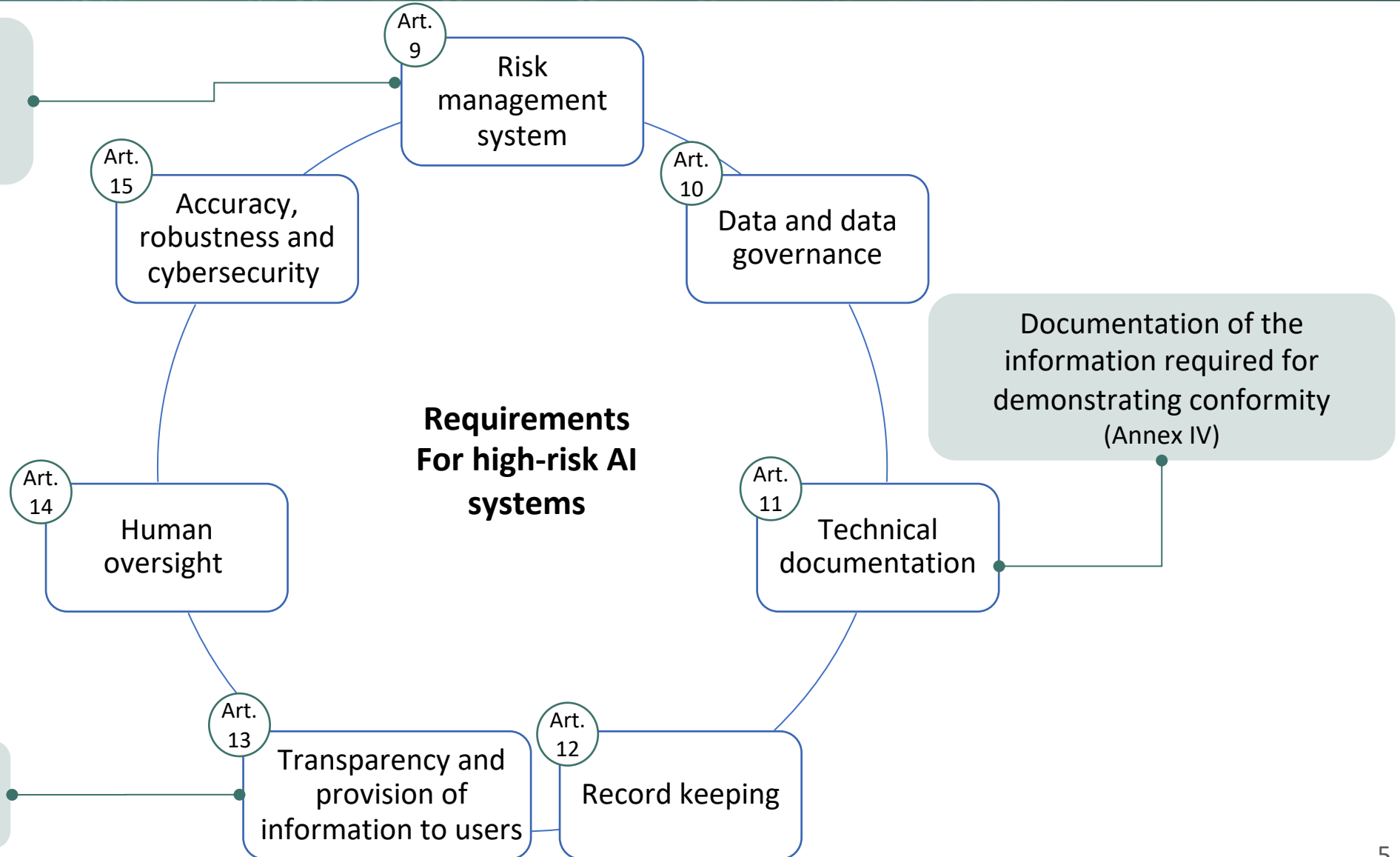
4

# Requirements for High-Risk AI Systems

Identification, assessment and mitigation of risks and impacts

**RQ2**
What information must be maintained regarding risk and impacts of high-risk AI systems according to the AI Act and ISO risk management standards?

Art. 9
Risk management system

Art. 15
Accuracy, robustness and cybersecurity

Art. 10
Data and data governance

Documentation of the information required for demonstrating conformity (Annex IV)

**Requirements For high-risk AI systems**

Art. 14
Human oversight

Art. 11
Technical documentation

Art. 13
Transparency and provision of information to users

Art. 12
Record keeping

Creating instruction of use

# Challenge

- Maintaining, querying, and sharing information associated with risks for compliance checking, demonstrating accountability, and building trust

- Challenges:
  - The pace of changes in AI systems
  - The amount of risk-related information
  - The complexities in the AI value chain

Using semantic web technologies:
- enables automation
- Interoperability

**RQ3**
To what extent can semantic web technologies assist with representing information and generating documentation for high-risk AI systems required by the AI Act?

# Research Questions

**1** What is the information required to determine whether an AI system is 'high-risk' as per the AI Act?

**2** What information must be maintained regarding risk and impacts of high-risk AI systems according to the AI Act and ISO risk management standards?

Identify information requirements from:
- the AI Act
- ISO 31000 family

**3** To what extent can semantic web technologies assist with representing information and generating documentation for high-risk AI systems required by the AI Act?

Create AIRO (AI Risk Ontology) demonstrate its applicability in real-world cases

# State of the Art

| Topic | Summary | Relation to this work |
|---|---|---|
| AI risk management standards | ISO 31000:2018 Risk management– Guidelines<br>ISO 31073:2022 Risk management — Vocabulary | Used for identifying risk concepts |
| AI risk taxonomies | Existing taxonomies of AI risks, harms, risk sources, & mitigation measures | Reusing the taxonomies for populating AIRO |
| Risk models & ontologies | - Generic risk models<br>- Domain-specific risk models | Reusing risk concepts |

# Ontology Development Methodology

AIRO: Ontology for representing AI Risks | Delaram Golpayegani et al. | SEMANTiCS 2022 | contact:sgolpays@tcd.ie | https://w3id.org/AIRO

# Describing High-Risk AI Systems

## Questions to identify whether an AI system is high-risk according to Annex III

| Question | concept | Relation with AISystem |
|---|---|---|
| What techniques are utilised in the system? | AI Technique | usesAITechnique |
| What domain is the system intended to be used in? | Domain | isAppliedWithinDomain |
| What is the intended purpose of the system? | Purpose | hasPurpose |
| What is the application of the system? | AI Application | hasApplication |
| Who is the intended user of the system? | AI User | hasAIUser |
| Who is the subject of the system? | AI Subject | hasAISubject |
| In which environment is the system used? | Environment Of Use | isUsedInEnvironment |

ANNEX I
ARTIFICIAL INTELLIGENCE TECHNIQUES AND APPROACHES
referred to in Article 3, point 1

(a) Machine learning approaches, including supervised, unsupervised and reinforcement learning, using a wide variety of methods including deep learning;

(b) Logic- and knowledge-based approaches, including knowledge representation, inductive (logic) programming, knowledge bases, inference and deductive engines, (symbolic) reasoning and expert systems;

(c) Statistical approaches, Bayesian estimation, search and optimization methods.

ANNEX III
HIGH-RISK AI SYSTEMS REFERRED TO IN ARTICLE 6(2)

High-risk AI systems pursuant to Article 6(2) are the AI systems listed in any of the following areas:

1. Biometric identification and categorisation of natural persons:

    (a) AI systems intended to be used for the 'real-time' and 'post' remote biometric identification of natural persons;

2. Management and operation of critical infrastructure:

    (a) AI systems intended to be used as safety components in the management and operation of road traffic and the supply of water, gas, heating and electricity.

3. Education and vocational training:

    (a) AI systems intended to be used for the purpose of determining access or assigning natural persons to educational and vocational training institutions;

    (b) AI systems intended to be used for the purpose of assessing students in educational and vocational training institutions and for assessing participants in tests commonly required for admission to educational institutions.

4. Employment, workers management and access to self-employment:

    (a) AI systems intended to be used for recruitment or selection of natural persons, notably for advertising vacancies, screening or filtering applications, evaluating candidates in the course of interviews or tests;

    (b) AI intended to be used for making decisions on promotion and termination of work-related contractual relationships, for task allocation and for monitoring and evaluating performance and behavior of persons in such relationships.

5. Access to and enjoyment of essential private services and public services and benefits:

    (a) AI systems intended to be used by public authorities or on behalf of public authorities to evaluate the eligibility of natural persons for public assistance benefits and services, as well as to grant, reduce, revoke, or reclaim such
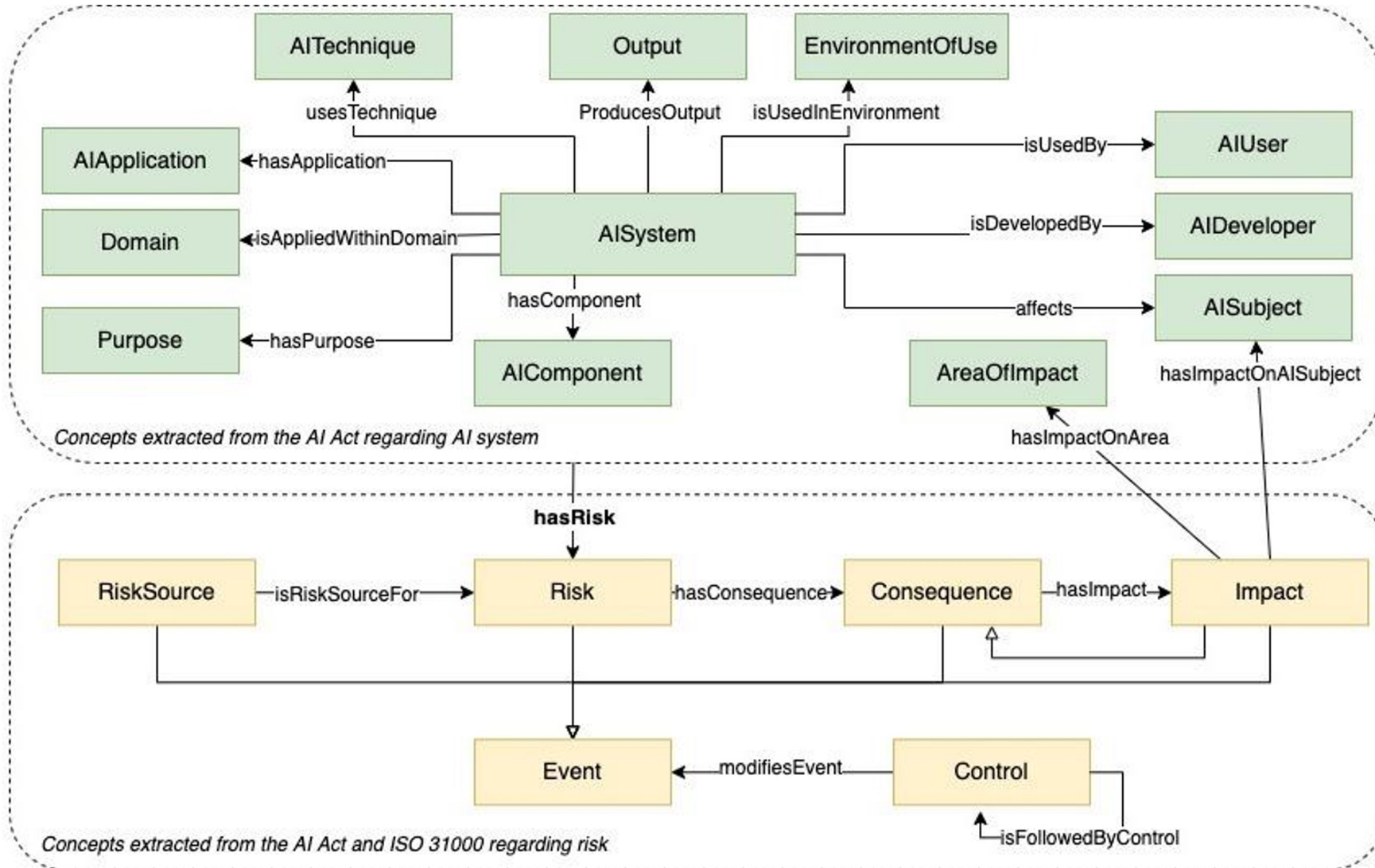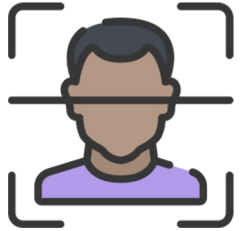
AIRO: Ontology for representing AI Risks | Delaram Golpayegani et al. | SEMANTiCS 2022 | contact:sgolpays@tcd.ie | https://w3id.org/AIRO

# AIRO Requirements
# Technical Documentation

| Annex IV | Required Info | Domain | Relation | Range |
|---|---|---|---|---|
| 1(a) | System's intended purpose | AISystem | hasPurpose | Purpose |
| | System's developers | AISystem | isDevelopedBy | AIDeveloper |
| | System's date | AISystem | dcterms:date | |
| | System's version | AISystem | hasVersion | Version |
| ... | | | | |
| 4 | Risks of AI system | AISystem | hasRisk | Risk |
| | Sources of the risk | RiskSource | isRiskSourceFor | Risk |
| | Consequences of the risk | Risk | hasConsequence | Consequence |
| | Harmful impacts of risk | Consequence | hasImpact | Impact |
| | Probability of risk | Risk | hasLikelihood | Likelihood |
| | Severity of impact | Impact | hasSeverity | Severity |
| | ... | | | |

**ANNEX IV**
**TECHNICAL DOCUMENTATION referred to in Article 11(1)**

The technical documentation referred to in Article 11(1) shall contain at least the following information, as applicable to the relevant AI system:

1. A general description of the AI system including:
   (a) its intended purpose, the person/s developing the system the date and the version of the system;
   (b) how the AI system interacts or can be used to interact with hardware or software that is not part of the AI system itself, where applicable;
   (c) the versions of relevant software or firmware and any requirement related to version update;
   (d) the description of all forms in which the AI system is placed on the market or put into service;
   (e) the description of hardware on which the AI system is intended to run;
   (f) where the AI system is a component of products, photographs or illustrations showing external features, marking and internal layout of those products;
   (g) instructions of use for the user and, where applicable installation instructions;

2. A detailed description of the elements of the AI system and of the process for its development, including:
   (a) the methods and steps performed for the development of the AI system, including, where relevant, recourse to pre-trained systems or tools provided by third parties and how these have been used, integrated or modified by the provider;
   (b) the design specifications of the system, namely the general logic of the AI system and of the algorithms; the key design choices including the rationale and assumptions made, also with regard to persons or groups of persons on which the system is intended to be used; the main classification choices; what the system is designed to optimise for and the relevance of the different parameters; the decisions about any possible trade-off made regarding the technical solutions adopted to comply with the requirements set out in Title III, Chapter 2;

# AIRO https://w3id.org/AIRO



Concepts extracted from the AI Act regarding AI system

Concepts extracted from the AI Act and ISO 31000 regarding risk

12

AIRO: Ontology for representing AI Risks | Delaram Golpayegani et al. | SEMANTiCS 2022 | contact:sgolpays@tcd.ie | https://w3id.org/AIRO

# Use-cases

**Use-case #1: Uber's Real-time ID Check System**
**Purpose**: Ensure the system is used by the registered driver
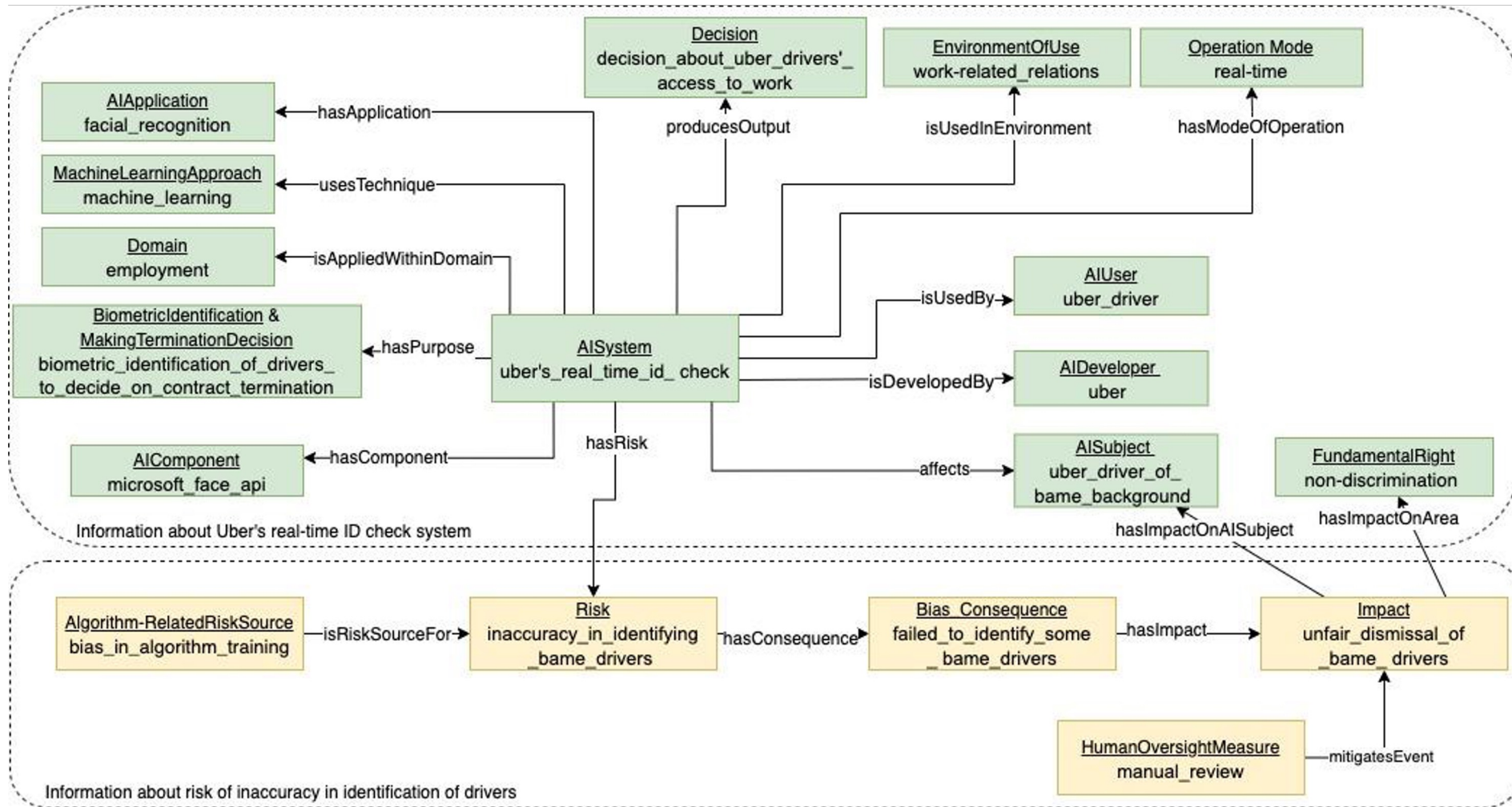**Main issue**: Discrimination against drivers of BAME background

**Use-case #2: VioGen Domestic Violence System**
**Purpose**: Determine the eligibility to access police protection by predicting the likelihood of a victim of gender violence to be assaulted by the same perpetrator again
**Main issue**: Inaccuracy of predictions

AIRO: Ontology for representing AI Risks | Delaram Golpayegani et al. | SEMANTiCS 2022 | contact:sgolpays@tcd.ie | https://w3id.org/AIRO

```
1   PREFIX airo: <https://w3id.org/AIRO#>
2   SELECT  ?system ?technique ?domain ?purpose
3           ?application ?user ?subject ?environment
4   WHERE {
5           ?system a airo:AISystem ;
6                   airo:usesTechnique ?technique ;
7                   airo:isUsedWithinDomain ?domain ;
8                   airo:hasPurpose ?purpose ;
9                   airo:hasApplication ?application ;
10                  airo:isUsedBy ?user ;
11                  airo:affects ?subject ;
12                  airo:isUsedInEnvironment ?environment . }
```

| AIRO concept | |
|---|---|
| AISystem | uber's real time id check |
| AITechnique | machine learning techniques |
| Domain | employment |
| Purpose | biometric identification of drivers to decide on contract termination |
| AIApplication | facial recognition |
| AIUser | uber driver |
| AISubject | uber driver of bame background |
| Environment OfUse | work related relations |

1. Biometric identification and categorisation of natural persons:

   (a) AI systems intended to be used for the 'real-time' and 'post' re... identification of natural persons;

4. Employment, workers management and access to self-employment:

   (a) AI systems intended to be used for recruitment or selection of natural persons, notably for advertising vacancies, screening or filtering applications, evaluating candidates in the course of interviews or tests;

   (b) AI intended to be used for making decisions on promotion and termination of work-related contractual relationships, for task allocation and for monitoring and evaluating performance and behavior of persons in such relationships.

- Manual analysis

**High Risk**

# SHACL Shapes for Automatic Identification of High-Risk AI

- "Rules" to determine whether AI satisfies conditions for being "high-risk"
- Choose your favourite flavour of rule languages & mechanisms

- We chose **SHACL**
- Why:
  - Flexible, Standardised
  - Extensible with plugins/features
  - Built-in documentation of outputs
  - Integrate to instead check outputs e.g. another rule engine

- We implement SHACL shapes for clauses defined in Annex III that determine high-risk

- Validation is to NOT satisfy the expressed criteria

```
1   @prefix dash: <http://datashapes.org/dash#> .
2   @prefix sh: <http://www.w3.org/ns/shacl#> .
3   @prefix airo: <https://w3id.org/AIRO#> .
4   @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
5   :AnnexIII-1
6       a sh:NodeShape ;
7       sh:targetClass airo:AISystem ;
8       sh:message "High-Risk AI System as per AI Act Annex III-1"@en ;
9       sh:description "Biometric Identification of Natural Persons"@en ;
10      sh:not [
11          a sh:PropertyShape ;
12          sh:path airo:hasPurpose ;
13          sh:class airo:BiometricIdentification; ] .
```
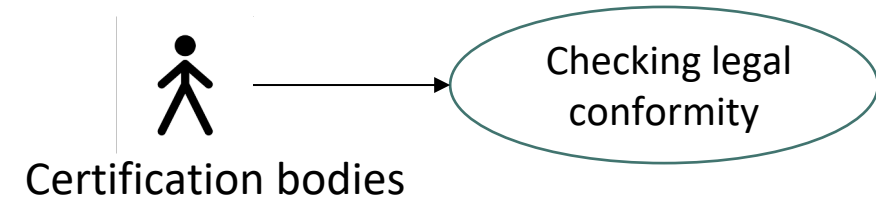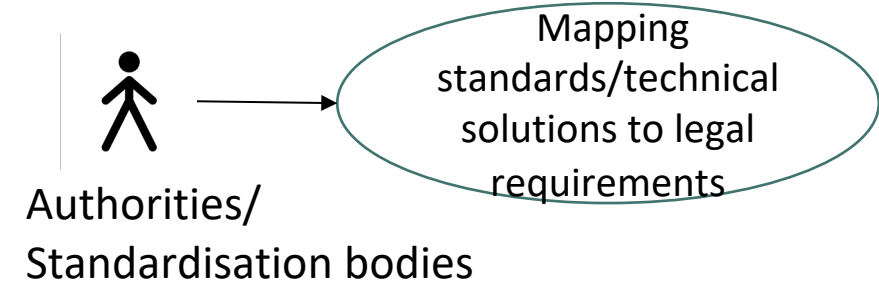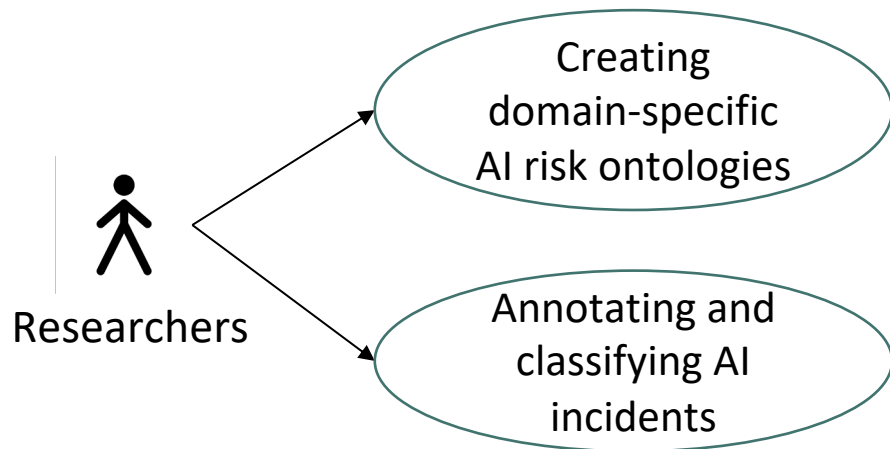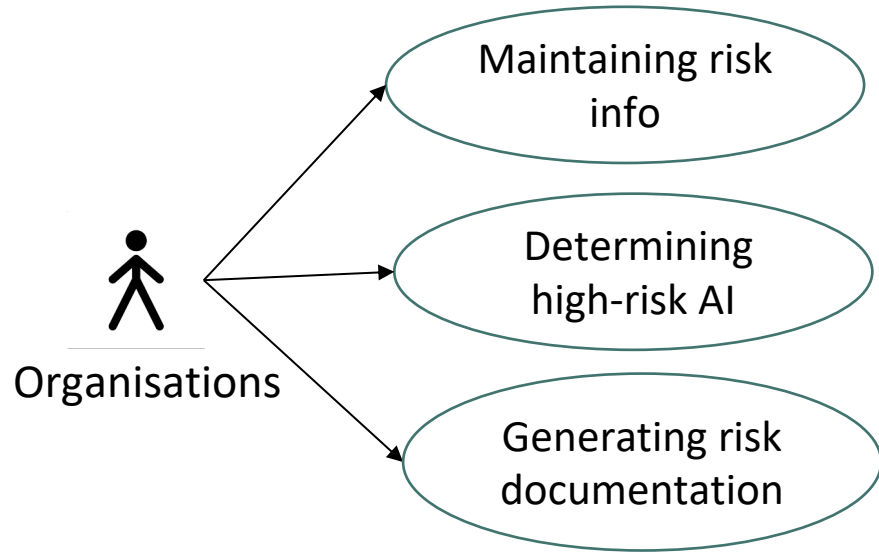
| Anx.IV. Required Information | Concept | Uber's Real-time ID Check |
|---|---|---|
| 1(a). System's intended purpose | Purpose | biometric_identification_of_drivers_to_decide_on_contract_termination |
| 1(a). System's developers | AIDeveloper | uber |
| 1(d). Forms in which AI system is placed on the market or put into service | AISystemForm | service |
| 2(e) & 3. Human oversight measures | HumanOversightControl | manual_review |
| 2(g). Discriminatory impacts of the system | Impact ImpactedArea | unfair_dismissal_of_bame_drivers non-discrimination |
| 3. Expected level of accuracy | AISystemAccuracy | high |
| 3. Foreseeable unintended outcomes of the risk 4. Consequences of the risk | Consequence | failed_to_identify_some_bame_drivers |
| 3 & 4. Sources of the risk | RiskSource | bias_in_algorithm_training |
| 4. Risks associated with the AI system | Risk | inaccuracy_in_identifying_bame_drivers |
| 4. Harmful impacts of the risk | Impact | unfair_dismissal_of_bame_drivers |
| 4. Severity of impact | Severity | N/A |
| 4. Impacted stakeholders | AISubject | uber_driver_of_bame_background |
| 4. Impacted area | AreaOfImpact | non-discrimination |
| 4. Risk management measures applied | Control | manual_review |

Domain Challenge

the incident reports do not provide *detailed information*

Why:
Implementation Details

# Benefit to Stakeholders

**Organisations**
- Maintaining risk info
- Determining high-risk AI
- Generating risk documentation

**Researchers**
- Creating domain-specific AI risk ontologies
- Annotating and classifying AI incidents

**Authorities/ Standardisation bodies**
- Mapping standards/technical solutions to legal requirements

**Certification bodies**
- Checking legal conformity

# Future Work

- Enhance AIRO to:
  - represent known categories of AI risks identified from real-world incidents
  - express provenance of AI risk and impact assessments

- Incorporate changes from the AI Act update and recently developed ISO standards

- Create rules for determining High-Risk AI

- Develop tools for risk documentation generation and sharing

- Apply AIRO's AI impact assessment for the GDPR's DPIA

# AIRO: an Ontology for Representing AI Risks based on the Proposed EU AI Act and ISO Risk Management Standards

Delaram Golpayegani,
Harshvardhan J. Pandit, Dave Lewis

Email: sgolpays@tcd.ie

Ontology: https://w3id.org/AIRO

GitHub:
https://github.com/delaramglp/AIRO