# VDE SPEC



# VCIO based description of systems for AI trustworthiness characterisation

VDE SPEC 90012 V1.0 (en)

**VDE**

## Preface

Publication date of this VDE SPEC draft: 25 April 2022

No draft has been published for the present VDE SPEC.

This VDE SPEC has been developed according to the public VDE SPEC specification. The development is carried out in VDE SPEC consortiums and does not require the participation of all potential stakeholders.

This document has been developed and adopted by the initiator(s) and authors named below:

- Peylo, Christoph / Robert Bosch GmbH (initiator)

- Slama, Dirk / Digital Trust Forum (initiator)

- Hallensleben, Sebastian / VDE e.V. (project lead)

- Hauschke, Andreas / VDE e.V. (author)

- Hildebrandt, Stefanie / VDE e.V. (project manager, co-author)

The experts listed below have been involved in the development process of this VDE SPEC and are part of the project group.

*Adamietz, Peter / BASF*

*Damboldt, Heiko / BASF*

*Fanidakis, Nikolaos / BASF*

*Grünke, Paul / KIT Karlsruher Institut für Technologie*

*Hagendorff, Thilo / Uni Tübingen*

*Hahn, Thomas / Siemens*

*Hapfelmeier, Andreas / Siemens*

*Hubig, Christoph / TU Darmstadt*

*Karls-Eberhard, Andrea / BASF*

*Karnouskos, Stamatis / SAP*

*Lachenmaier, Jens / Steinbeis Institut*

*Loh, Wulf / Uni Tübingen*

*Lohmüller, Stina / iRights.Lab*

*Puntschuh, Michael / iRights.Lab*

*Schlesinger, Dirk / TÜV Süd*

*Wieczorek, Sebastian / SAP*

*Zillner, Sonja / Siemens*

At present, there are no standards covering this topic in any German standard.

Despite great efforts to ensure the correctness, reliability and precision of technical and non-technical descriptions, the VDE SPEC project group cannot give any explicit or implicit warranty with respect to the correctness of the document. The application of this document is made in the knowledge that the VDE SPEC project group cannot be held liable for damage or loss of any kind. The application of this VDE SPEC does not release users from the responsibility for their own actions and therefore at their own risk.

In the course of establishing and/or introducing products into the European market, producers must carry out a risk analysis in order to first determine what risks the product may entail. After carrying out the risk analysis, they evaluate these risks and, if necessary, take appropriate measures to effectively eliminate or minimise the risks (risk mitigation). This VDE SPEC does not absolve the user from this responsibility.

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. VDE shall not be held responsible for identifying any or all such patent rights.

# Content

# List of figures

## List of tables

# 1 Scope

This VDE SPEC provides a way to describe certain socio-technical characteristics of systems and applications that incorporate artificial intelligence techniques and methods. The scope of application refers to products for which a particularly demanding level of trust is desired or required.

By applying the VCIO model explained in this standard, it is possible to describe whether a product adheres to specific values and can be trusted. This standard can therefore e.g., form the basis for attaching a trust label to a product.

The product characterisation according to this standard can be used in a wide variety of contexts. End consumers, companies and government organizations can use the description to define requirements or to compare different products. In doing so, it also becomes possible to assess the compliance with regard to different values (for example, one product might better comply with privacy requirements, while the other might comply better with transparency criteria).

In addition, target requirements can be set during the development of a given product. Those requirements are then considered in the development process in order to achieve a desired value compliance.

The standardised description is independent of the risk posed by the product and does not define any minimum requirements in the context of this. It describes compliance with the specified values in an orthogonal manner. Nevertheless, companies, users or government bodies can themselves set requirements for a minimum level within this framework.

The consortium has worked towards making this standard compatible with the emerging AI Act at the European level. In the case of AI products, the objective is to have a description of trustworthiness aspects that both demonstrate compliance of the product with the AI Act and provide differentiation in the market.

The focus of the standard is on systems and applications that incorporate artificial intelligence techniques and methods. The criteria, indicators and observables therefore aim at characteristics of AI systems, like underlying data sets, the precise definition of the scope, the development, the application, processes, and the clear assignment of responsibilities. In addition, aspects that are not limited to AI systems, but are necessary to demonstrate their trustworthiness, were considered.

# 2 Terms and definitions

The terms and definitions are used differently in different contexts, like AI community, safety community, regulation and legislation. To prevent confusion this section gives an overview of the terms and definitions in the VDE Spec, which aims to be close to the proposed EU AI Act [1].

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

– IEC Electropedia: available at http://www.electropedia.org
– ISO Online browsing platform: available at http://www.iso.org/obp

## 2.1
### accuracy
metric calculated by dividing the number of correctly classified data by the total number of classified data

## 2.2
### adversarial attack
attack on an AI system by an input variation with the aim to manipulate the output of the system

## 2.3
### affected person or entity
people or entities that are directly affected by the system

E.g., because the system classifies or predicts properties of them or collects data from them.

## 2.4
### AI model
knowledge representation using Artificial Intelligence techniques

## 2.5
## AI strategy
processes in the lifecycle of the AI system

## 2.6
## AI system
means software that is developed with one or more of the techniques and approaches listed in Annex I of the proposed EU *ARTIFICIAL INTELLIGENCE ACT* [1] and can, for a given set of human-defined objectives, generate outputs such as content, predictions, recommendations, or decisions influencing the environments they interact with

[SOURCE: [1]]

## 2.7
## Artificial Intelligence techniques
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.

[SOURCE: [1]]

## 2.8
## AI application
Input-output mapping in a given context of use, based on the implemented AI system. Consists of the AI system and optionally the regarding embedding, including additional software components, pre-/postprocessing and an interface for input, output and monitoring.

[SOURCE: adapted from [2]]

## 2.9
## availability
<accessibility/usability> property of being accessible and useable upon demand by an authorized entity

[SOURCE: ISO/TS 21089:2018]

## 2.10
## bias
refers to systematic statistical correlations or distortions that can lead to inaccurate or discriminatory outcomes

E.g., in the data, the AI system or/and its predictions

## 2.11
## clickwork
task of labeling or collecting data for ML-applications accomplished by humans

## 2.12
## confidentiality
<not disclosed> property that information is not made available or disclosed to unauthorized individuals, entities, or processes

[SOURCE: ISO/TS 21089:2018]

## 2.13
## criterion
concretizes values by establishing the reference to the states of affairs to be shaped under the values

## 2.14
## data
reinterpretable representation of information in a formalized manner suitable for communication, interpretation, or processing

[SOURCE: ISO/IEC 25000:2014)

## 2.15
### data centres
physical entity in which data is processed (e.g. servers)

## 2.16
### data poisoning
malicious modification or input of false data to manipulate the output of an AI system

## 2.17
### datasheet
a structured documentation of the characteristics of a dataset

## 2.18
### development data sets
data which is used for developing an AI system, consists of training, validation, and testing data

## 2.19
### environment
<system> context determining the setting and circumstances of all influences upon a system

[SOURCE: ISO/IEC/IEEE 12207:2017]

## 2.20
### explainability
ability or function of an AI system or AI application to explain its outputs, the relation from inputs to outputs or its general behaviour

## 2.21
### global interpretability
possibility to understand the whole logic of an AI system and follow the entire reasoning leading to all the different possible outcomes

[SOURCE: adapted from [3]]

## 2.22
### group
<suffix> indicates that the preceding entity can be divided into meaningful subgroups

## 2.23
### harm
injury or damage to the (physical and mental) health or dignity of people, or impairments to economic or participatory opportunities, or damage to the environment

[SOURCE: adapted from ISO/IEC Guide 51:2014, 3.1]

## 2.24
### hazard
potential source of harm

[SOURCE: ISO/IEC Guide 51:2014, 3.2]

## 2.25
### Human In Command (Control)
### HIC
refers to the necessity for human final decision making, based on the suggestions of the AI system. This includes the ability to decide when and how to use the system in any particular situation. This can include the decision not to use an AI system in a particular situation, to establish levels of human discretion during the use of the system, or to ensure the ability to override a decision made by a system.

[SOURCE: [4]]

## 2.26
### Human In The Loop
### HITL
refers to the capability for human intervention in every decision cycle of the system

[SOURCE: [4]]

## 2.27
## Human On The Loop
## HOTL

refers to the capability for human intervention during the design cycle of the system and monitoring the system's operation

[SOURCE: [4]]

## 2.28
## indicator

an information instance about properties of elements of a state of affairs to be recorded, which are decisive for the qualitative fulfilment/non-fulfilment "(anchor indicators)" or the degree of fulfilment of a criterion

## 2.29
## integrity

designed such that any modification of the electronically stored information, without proper authorization, is not possible

[SOURCE: ISO 17364:2013)

## 2.30
## intended use

reader oriented description of the purpose of the AI system and for what it should and should not be used

## 2.31
## life cycle

evolution of a system, product, service, project, or other human-made entity from conception through retirement

[SOURCE: ISO/IEC/IEEE 12207:2017]

## 2.32
## local interpretability

possibility to understand only the reasons for a specific output or decision. Here only the single prediction/decision is interpretable

[SOURCE: adapted from [3]]

## 2.33
## machine Learning (ML)

process of optimizing model parameters in an AI system through computational techniques, such that the model's behaviour reflects the data or experience

[SOURCE: adapted from ISO/IEC DIS 22989, 3.2.9]

## 2.34
## model

mathematical representation of a physical system or process

[SOURCE: ISO/TS 18166:2016]

## 2.35
## nudging

subtle measures, which remain unconscious to the user, to make certain decision options attractive or unattractive by means of aesthetic/affective impressions and/or the design of the effort required to make and realize respective decisions. This makes use of habits, likes, and dislikes that are known to the systems from the user profile. Here, the decision options themselves remain fundamentally disposable.

## 2.36
## observable

an observable (or measurable) quantity about the state of the element of a situation covered by the indicator

### 2.37
**Operational Design Domain (ODD)**

The set of environments and situations the item is intended to operate within. This includes not only direct environmental conditions and geographic restrictions, but also a characterization of the set of objects, events, and other conditions that will occur within that environment.

Note 1 to entry: A system has a single ODD by definition. Assessment is made with regard to the entire ODD.

[SOURCE: UL4600]

### 2.38
**Operational Domain (OD)**

set of environments and situations the item can reasonably encounter.

[SOURCE: UL4600; ASAM OpenODD: Concept Paper]

### 2.39
**optimization metric**

metric which is optimized by the training process of the AI system. Often a loss function is used and minimized during the training process.

### 2.40
**Out of sample data**

data which was not included in the training process

### 2.41
**performance and evaluation metric**

metric which is used to evaluate the performance or other requirements of an AI system

### 2.42
**persuasive computing**

aims at no longer allowing independent judgments and evaluations of the systemic specifications by the user. This concerns both decisions of the systems for the preselection of the options for action and margins of the decision when informing oneself, evaluating and choosing, as well as systemic information, evaluations and decisions themselves, which are presented as convincing without alternative and/or are made accordingly by the systems themselves (e.g. no more monitoring to be provided or made possible in the case of HOTL).

### 2.43
**reliability**

ability of a device or a system to perform its intended function under given conditions of use for a specified period of time or number of cycles

[SOURCE: ISO/TS 17574:2017]

### 2.44
**system**

set of interrelated or interacting elements

[SOURCE: ISO 9000:2015(en), 3.5.1]

### 2.45
**target group**

meaningful and appropriate classification of persons or entities that interact or are directly affected by an AI System into groups according to different characteristics such as domain knowledge, skill level, etc.

Includes users and affected persons, but can also be used in a limited context, such as target users or target affected persons

### 2.46
**testing data**

means data used for providing an independent evaluation of the trained and validated AI system in order to confirm the expected performance of that system before its placing on the market or putting into service

[SOURCE: [1]]

**2.47**
**testing strategy**
processes to test and assure certain properties of an AI system or AI application

# 3   Abbreviations

| | |
|---|---|
| AI | Artificial Intelligence |
| DKE | Deutsche Kommission Elektrotechnik Elektronik Informationstechnik (www.dke.de) |
| HIC | Human in Command (Control) |
| HITL | Human in the Loop |
| HOTL | Human on the Loop |
| IEC | International Electrotechnical Commission (www.iec.ch) |
| ISO | International Organization for Standardization (www.iso.org) |
| ML | Machine Learning |
| OD | Operational Domain |
| ODD | Operational Design Domain |
| VCIO | Values Criteria Indicators Observables |
| VDE | Verband der Elektrotechnik, Elektronik und Informationstechnik e.V. (www.vde.de) |

# 4 VCIO Model

## 4.1 General

The VCIO model this standard is based on, has previously been demonstrated in the context of AI in "From Principles to Practice – An interdisciplinary framework to operationalise AI ethics" by the AI Ethics Impact Group. The report introduces the approach as follows:

*"The VCIO model distinguishes and combines the four concepts of values, criteria, indicators and observables for the evaluation of AI. […] As values are abstract, often in conflict with each other, and do not include means to evaluate their implementation, it is essential to have other components to fulfil these tasks. This is where the criteria, indicators and observables of the VCIO approach come into play. […]*

*The VCIO approach, therefore, fulfils three tasks:*

1) *It clarifies what is meant by a particular value (value definition).*

2) *It explains in a comprehensible manner how to check or observe whether or to what extent a technical system fulfils or violates a value (measurement).*

3) *It acknowledges the existence of value conflicts and explains how to deal with these conflicts depending on the application context (balancing).*

*To practically implement AI trustworthiness, the VCIO approach operates on four levels:*

*Values formulate a general trustworthiness concern, something that should guide our actions. They are defined at the highest level (as transparency, for example). To verify whether an algorithm fulfils or violates specific values, we must specify Criteria that define the fulfilment or violation of the respective value. Since it is usually not possible to directly observe whether a criterion is met, we need Indicators (as a specific type of sign) to monitor this. Indicators relate criteria on the one hand with Observables on the other.*

*The four hierarchical levels provided by values, criteria, indicators and observables are closely linked, where the fulfilment of the higher level depends on the lower level. However, it is not possible to derive the lower levels from the higher ones in a straightforward, i.e. deductive way. Instead, the normative load runs through all four levels and requires new deliberations at all levels, in the course of which the particular instances must be negotiated in detail.*

*Note that, typically, several indicators are required to evaluate the fulfilment of a criterion; however, we can also use the same indicator as part of different criteria. As there are no deductive relationships between values, criteria, indicators, and observables, at each stage of their determination, normative decisions need to be made in a scientific and technically informed context."* [5]



**Figure 1 – Composition of the VCIO-Model**
(SOURCE: adapted from [5])

## 4.2 Values

### 4.2.1 General

The following sections show the composition of the values:

- Transparency (Section 4.2.2)
- Accountability (Section 4.2.3)
- Privacy (Section 4.2.4)
- Fairness (Section 4.2.5)
- Reliability (Section 4.2.6)

The VCIO´s are represented as tables in the following sections. These tables are constructed as follows:

- The first tables in section 4.2.2 – section 4.2.6 shows an overview of each value and the structure of the underlying criteria and indicators
- The second tables in section 4.2.2 – section 4.2.6 show the corresponding Observables to the indicators subsumed under the criteria.
    - o Here it is also marked if an indicator is skippable (section 5.2.4) and whether it represents, a negative anchor indicator or a positive anchor indicator.
- The criteria, indicators and observables are indexed as follows:
    - o Criteria are indexed with the first letter of the corresponding value and an index number (1, 2, 3, …, n)
    - o Indicators are indexed with the corresponding criterion Index followed by a point and an additional index number.

Observables are indexed by the corresponding indicator index followed by the Level ("A" to "G").

### 4.2.2 Transparency



**TRANSPARENCY**

**T1**
Documentation of data sets

**T1.1** - Is the data's origin documented?

**T1.2** - Are the characteristics of data sets analyzed and documented?

**T2**
Documentation about the AI systems operation

**T2.1 -** Are the characteristics of the AI system(s) documented?

**T2.2 -** Are the characteristics of the AI application documented?

**T3**
Intelligibility

**T3.1** - Have the most intelligible AI models/ systems been selected that can fulfil the application purpose?

**T3.2 -** What degree of explainability including a regarding documentation is provided?

**T3.3** - Are the explanations of the AI system/application outcome designed in a way that adequately informs the affected persons?

**T3.4** - Is the interface of the AI system/application designed in a way that adequately informs the user groups about the outcomes and mechanisms?

**T4**
Accessibility (outside of relevant authorities)

**T4.1** - Who has access to the AI System and the AI application?

**T4.2** - Who has access to the datasets?

**T4.3** - Who has access to the documentation regarding the AI system/application and its data?

**T4.4 -** Who can see which data attributes (including pre-processing) were used as an input for the AI system/application to generate its output?
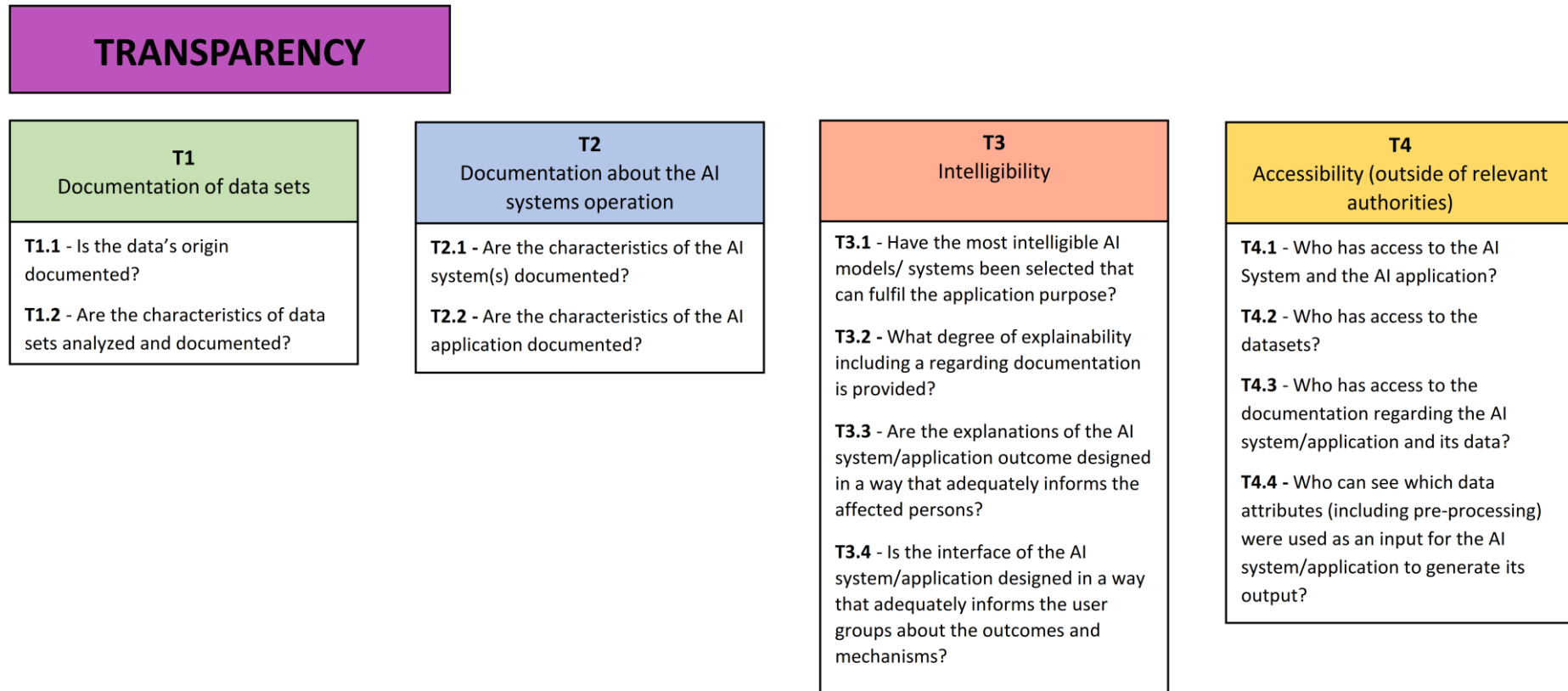
**Figure 2 – Composition of Transparency Criteria and Indicators**

## Table T1 – Documentation of data sets

| T1 | Documentation of data sets | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **T1.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Is the data's origin documented?** | *Established structured notations like "datasheets for datasets" are recommended here.* | Yes, with **structured datasheets**, including **detailed** information on:<br>■ data handling<br>■ data collector<br>■ data collection method | Yes, with **structured datasheets** including **detailed** information on:<br>■ data handling<br>■ data collector<br>■ data collection method<br>containing few (not all) information | Yes, information is collected on:<br>■ data collection method<br>■ data handling<br>■ data collector<br>containing few or missing information **without structured datasheets.** | Only general information on the data's origin is documented. | | | No |
| **T1.2** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Are the characteristics of data sets analysed and documented?** | *Explorative question. Related to R1.2 and F1.7. Characteristics of data sets are:*<br>■ *fit to operational domain*<br>■ *number of data points in relationship to the domain*<br>■ *individual or perturbated data points potential for bias*<br>■ *analysis for potential proxies* | Yes, **structured** information about the characteristics of data sets, including **all** mentioned characteristics are provided. | Yes, **structured** information about the characteristics of data sets, including **all** mentioned characteristics, are provided **only some contain few or missing information.** | Yes, **structured** information about the characteristics of data sets are provided, **but not covering all mentioned characteristics.** | Yes, **some** not-structured information about the characteristics of data sets are provided. | | | None |

## Table T2 – Documentation about the AI systems operation

| T2 | Documentation about the AI systems operation | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **T2.1** | | A | B | C | D | E | F | G |
| **Are the characteristics of the AI system(s) documented?** | *Characteristics of AI system(s) are:*<br>■ *architecture or model graph (Number of layers, Parameters, connectivity input-output dimensions)*<br>■ *expected input data*<br>■ *expected output data*<br>■ *parameter precision (e.g. 8/16/32-bit)* | Yes, characteristics are documented, <u>including:</u><br><br>■ architecture or model graph<br>■ expected input data<br>■ expected output data<br>■ parameter precision<br><br>If there are relevant stakeholders, the documentation is **available** to them. | | Yes, characteristics are documented, <u>including:</u><br><br>■ architecture or model graph<br>■ expected input data<br>■ expected output data<br>■ parameter precision<br><br>The documentation is available for the competent authorities, but nor for all relevant stakeholders. | Yes, some characteristics are documented. | | | No |
| **T2.2** | | A | B | C | D | E | F | G |
| **Are the characteristics of the AI Application documented?** | *Characteristics of AI application are:*<br>■ *Hardware requirements*<br>■ *Training method (e.g. online/ offline/ …)*<br>■ *System architecture*<br>■ *Flow of information* | Yes, characteristics are documented, <u>including:</u><br><br>■ Hardware requirements<br>■ Training method<br>■ System architecture<br>■ Flow of information<br><br><br>If there are relevant stakeholders, the documentation is available to them. | Yes, characteristics are documented, <u>including:</u><br><br>■ Hardware requirements<br>■ Training method<br><br><br><br>If there are relevant stakeholders, the documentation is available to them. | Yes, characteristics are documented, but they are not available to all relevant stakeholders. | | | | No |

## Table T3 – Intelligibility

| T3 | Intelligibility | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **T3.1** | | A | B | C | D | E | F | G |
| **Have the most intelligible AI models/systems been selected that can fulfil the application purpose?** | Aspects *of the justification:*<br>■ *performance*<br>■ *efficiency*<br>■ *simplicity*<br>■ *intelligibility*<br>■ *locally / globally interpretable* | Yes, the AI system and application approach has been evaluated, documented and justified. The most intelligible model from this analysis has been used.<br><br>This evaluation is open to the public | No, but the AI system was evaluated regarding interpretability.<br><br>This evaluation is open to the public. | | No, but the AI system was evaluated regarding interpretability.<br><br>This evaluation is open to the competent authorities. | | | No, the AI system (architecture) has not been evaluated. |
| **T3.2** | | A | B | C | D | E | F | G |
| **What degree of explainability including a regarding documentation is provided?** | *Definition of local and global explainability is in the glossary.* | An interface with details about the AI system/application and the decision-making process is available and the AI application is globally interpretable. | An interface with details about the AI system/application and the decision-making process is available and the AI application is locally interpretable. | An interface with details about the AI system/application and the decision-making process is available. It allows to extract the most relevant features and roughly represent their interrelationships and interactions. | The modes of interpretability are available, but can only be used/understood post hoc by experts. | The modes of interpretability need to be adjusted ex post to the individual model and use by experts. | The model is only theoretically comprehensible. | There are no known modes of interpretability. |

| T3.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is the interface of the AI system/application designed in a way that adequately informs the user groups about the outcomes and mechanisms?** | *User-oriented: The interface / interaction with the system should be designed in such a way that the user-groups understand the procedures and outcomes. The understanding depends on the relevant Information necessary to adequately to fulfil their task.* | Yes, the interface of the system is based on the feedback **of the users-groups and affected persons,** e.g.:<br>■ user and affected person group analysis<br>■ tested *with* the user-group.<br>■ experiences from the analysis or *test* of former products | Yes, the interface of the system is based on the feedback of **specific target users**, e.g.:<br>■ user and *affected* person group analysis<br>■ tested with *the* user-group.<br>■ experiences from the analysis or *test* of former products | Yes, but without participation of the target groups. | Yes, but the modes or interpretability are only specific for one target group. | | | No, the modes of interpretability are not target-group specific. |
| T3.4 | | A | B | C | D | E | F | G |
| **Are the explanations of the AI system/application outcome designed in a way that adequately informs the affected persons?** | *Affected Persons oriented: The explanations of the system should be designed in such a way that the affected persons understand the procedures and outcomes. The understanding depends on the relevant Information necessary to understand the effects to them.* | Yes, the explanation of the system is based on the feedback **of the affected persons,** e.g.:<br>■ affected person groups analysis<br>■ tested with the affected person-groups<br>■ experiences from the analysis or test of former products | Yes, the explanation of the system is based on the feedback **of specific affected person groups,** e.g.:<br>■ affected person-groups analysis<br>■ tested with the affected person groups<br>■ experiences from the analysis or test of former products | Yes, but without participation of the affected persons. | Yes, but the modes or interpretability are only specific for one affected persons group. | | | No, the modes of interpretability are not affected person group specific. |

**Table T4 – Accessibility (outside of relevant authorities)**

| T4 | Accessibility (outside of relevant authorities) | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **T4.1** | | A | B | C | D | E | F | G |
| **Who has access to the AI System and the AI application?** | *If an NDA is used in this context, it must not prevent the publication of conclusions drawn from access and analysis of the data, the right to analyse the data freely and fully as well as partly publication of single data points or database entries as illustrations for conclusions.* | With the possibility of non-disclosure agreement: <br>■ operators of the AI system <br>■ competent authorities <br>■ additional information and trust intermediaries (e.g. regulators, watchdogs, research, courts) | With the possibility of non-disclosure agreement: <br>■ operators of the AI system <br>■ competent authorities | | Only competent authorities. | | Nobody outside of the company. | Nobody outside of the development team, not even inside the company. |
| **T4.2** | | A | B | C | D | E | F | G |
| **Who has access to the datasets?** | *If an NDA is used in this context, It must not prevent the publication of conclusions drawn from access and analysis of the data, the right to analyse the data freely and fully as well as partly publication of single data points or database entries as illustrations for conclusions.* | With the possibility of non-disclosure agreement: <br>■ operators of the AI system <br>■ competent authorities <br>■ additional information and trust intermediaries (e.g. regulators, watchdogs, research, courts) | With the possibility of non-disclosure agreement: <br>■ operators of the AI system <br>■ competent authorities | | Only competent authorities. | | Nobody outside of the company. | Nobody outside of the development team, not even inside the company. |

| T4.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Who has access to the documentation regarding the AI system/application and its data?** | *Documentation from T1.1 and T1.2., including a short description of the operational domain.* | Everyone | With the possibility of non-disclosure agreement:<br><br>▪ operators of the AI system<br>▪ competent authorities<br>▪ additional information and trust intermediaries (e.g. regulators, watchdogs, research, courts) | With the possibility of non-disclosure agreement:<br><br>▪ operators of the AI system<br>▪ competent authorities | Only competent authorities. | | Nobody outside of the company. | Nobody outside of the development team, not even inside the company. |
| T4.4 | | A | B | C | D | E | F | G |
| **Who can see which data attributes (including pre-processing) were used as an input for the AI system/application to generate its output?** | *This refers only to the name/label and not to the individual content of the data attribute.* | Everyone | With the possibility of non-disclosure agreement:<br><br>▪ operators of the AI system<br>▪ competent authorities<br>▪ additional information and trust intermediaries (e.g. regulators, watchdogs, research, courts) | With the possibility of non-disclosure agreement:<br><br>▪ operators of the AI system<br>▪ competent authorities | Only competent authorities. | | Nobody outside of the company. | Nobody outside of the development team, not even inside the company. |

### 4.2.3　Accountability

**ACCOUNTABILITY**

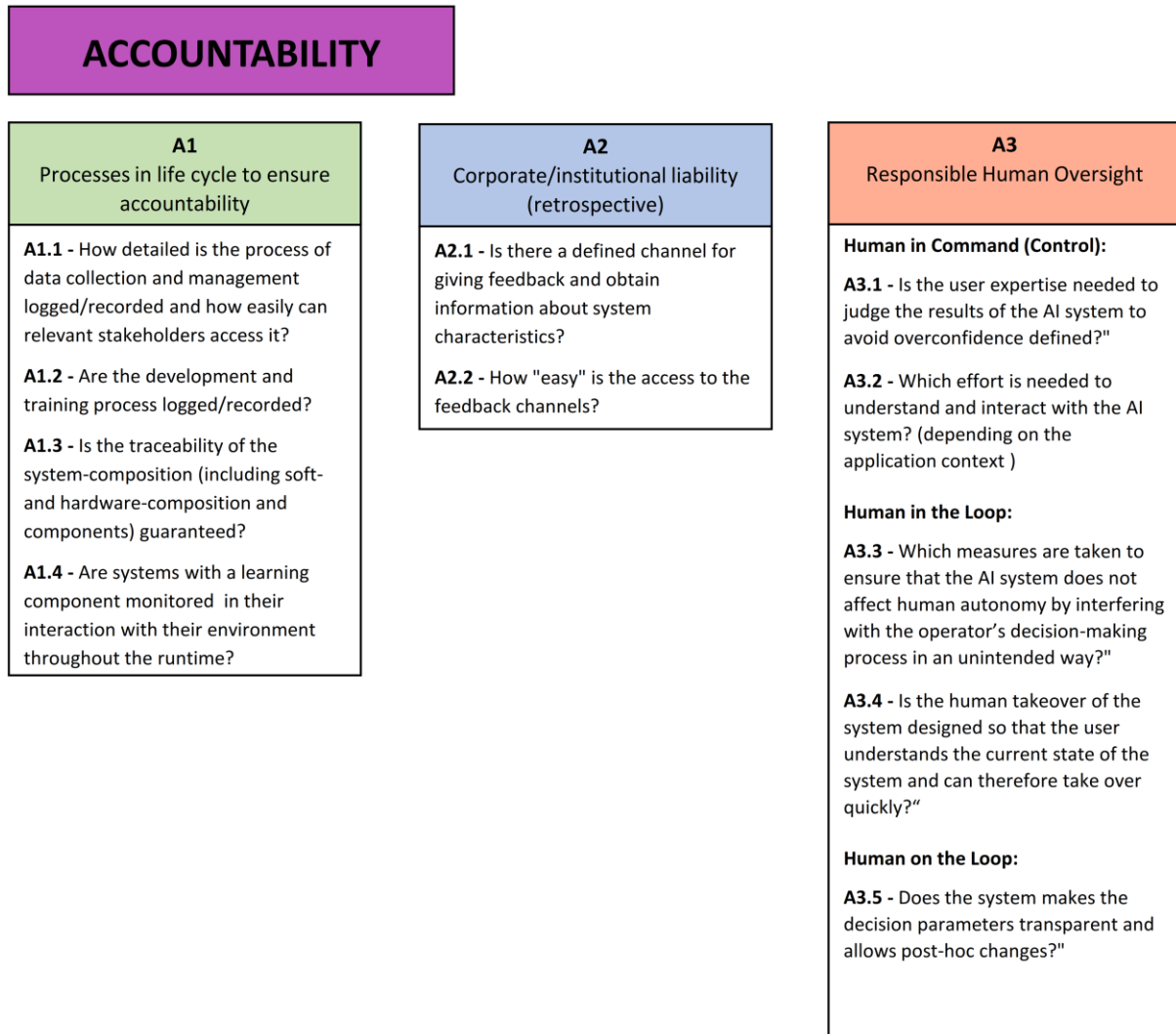| **A1**<br>Processes in life cycle to ensure accountability | **A2**<br>Corporate/institutional liability (retrospective) | **A3**<br>Responsible Human Oversight |
|---|---|---|
| **A1.1 -** How detailed is the process of data collection and management logged/recorded and how easily can relevant stakeholders access it?<br><br>**A1.2 -** Are the development and training process logged/recorded?<br><br>**A1.3 -** Is the traceability of the system-composition (including soft- and hardware-composition and components) guaranteed?<br><br>**A1.4 -** Are systems with a learning component monitored in their interaction with their environment throughout the runtime? | **A2.1 -** Is there a defined channel for giving feedback and obtain information about system characteristics?<br><br>**A2.2 -** How "easy" is the access to the feedback channels? | **Human in Command (Control):**<br><br>**A3.1 -** Is the user expertise needed to judge the results of the AI system to avoid overconfidence defined?"<br><br>**A3.2 -** Which effort is needed to understand and interact with the AI system? (depending on the application context )<br><br>**Human in the Loop:**<br><br>**A3.3 -** Which measures are taken to ensure that the AI system does not affect human autonomy by interfering with the operator's decision-making process in an unintended way?"<br><br>**A3.4 -** Is the human takeover of the system designed so that the user understands the current state of the system and can therefore take over quickly?"<br><br>**Human on the Loop:**<br><br>**A3.5 -** Does the system makes the decision parameters transparent and allows post-hoc changes?" |

**Figure 3 – Composition of Accountability**

## Table A1 – Processes in life cycle to ensure accountability

| A.1 | Processes in life cycle to ensure accountability | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **A1.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **How detailed is the process of data collection and management logged/recorded and how easily can relevant stakeholders access it?** | *Relevant stakeholders can for example be data users, product manager or competent authorities.*<br><br>*Logs/Records have to be stored for a reasonable time to allow delayed analysis.* | Logging/Records includes:<br>1. origin of data<br>2. responsible person<br>3. relevant data preparation processing operations (annotation, labelling, cleaning, enrichment, aggregation, ...)<br>4. Recovery of data in every stage<br><br>Stakeholder access:<br>■ easy, universal format | Logging/Records includes:<br>1. origin of data<br>2. responsible person<br>3. relevant data preparation processing operations (annotation, labelling, cleaning, enrichment, aggregation, ...)<br>4. Recovery of data in every stage<br><br>Stakeholder access:<br>■ In an unstructured and not clearly/prepared form<br>■ only for competent authorities | Logging/Records includes:<br>1. origin of data<br>2. responsible person<br><br>Stakeholder access:<br>■ In an unstructured and not clearly/prepared form<br>■ only for competent authorities | There is no logging/recording, but details about origin of data are documented.<br><br>Stakeholder access:<br>■ In an unstructured and not clearly/prepared form<br>■ only for competent authorities | | | The data collection process is not logged or documented. |

| A1.2 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are the development and training process logged/recorded?** | | Yes, comprehensive logging/recording including the responsibilities is available, including: <br> 1. operation in which the data was used and how they have been modified <br> 2. version control of AI systems and the involved data <br> 3. version restore/recovery of AI systems <br><br> and is available for all relevant stakeholders. | Yes, comprehensive logging/recording including the responsibilities is available, including: <br> 1. operation in which the data was used and how they have been modified <br> 2. version control of AI systems and the involved data <br><br> and is available for all relevant stakeholders. | No logging/recording available, but a general description of the development and training process including responsibilities is provided for all relevant stakeholders. | No logging/recording, but a general description of the development and training process is provided. | | | There is no logging and information about the process. |

| A1.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is the traceability of the system-composition (including soft- and hardware-composition and components) guaranteed?** | *Software-components can be:* <br> *AI model* <br> *AI system* <br> *AI application* <br> *Hardware-components can be:* <br> *…* | There is sufficient information about the system available: <br> ■ to easily reconstruct the composition of the system <br> ■ at every time in its lifecycle | There is sufficient information about the system available: <br> ■ to easily reconstruct the composition of the system <br> ■ at major inflection points (releases, gates) in its lifecycle | There is sufficient information about the system available to: <br> ■ reconstruct the composition of the system with additional efforts | There is not enough information to reconstruct the composition of the system. | | | No |

| A1.4 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are systems with a learning component monitored in their interaction with their environment throughout the runtime?** | *Learning systems, which adapt their behaviour during their use should have additional monitoring applications that track the changes in the system and highlight how the evolving systems differ from the original one.* | Yes, the learning process is monitored:<br><br>■ Logging of input-output behaviour is available for a defined period<br><br>■ Misuse is detected and reported<br><br>■ Concept Drift and Data Drift is detected and reported (e.g. changing operational domain properties)<br><br>If required, the information can be adequately prepared and made available to the relevant stakeholders. | Yes, the learning process is monitored:<br><br>■ Logging of input-output behaviour is available for a defined period<br><br>■ Misuse is detected and reported<br><br>■ Concept Drift and Data Drift is detected and reported (e.g. changing operational domain properties) | The learning process will be logged, but reviewed only at infrequent intervals<br><br>■ Misuse is detected and reported<br><br>■ Concept Drift and Data Drift is detected and reported (e.g. changing operational domain properties) | | | The learning process will be logged, but review is not planned. | There is no logging or monitoring of the learning process. |

| A.2 | Corporate/institutional liability (retrospective) | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **A2.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Is there a defined channel for giving feedback and obtain information about system characteristics?** | *Information and explanations about theses that can be received:*<br>■ *A1.4 - Interactions of learning Systems with the environment*<br>■ *F1.2 - Target Groups*<br>■ *F1.3 - Marginalised Groups*<br>■ *F1.7 - Bias*<br>■ *R1.1 - ODD and Intended Use*<br>■ *R1.5 - Risk and potential harms*<br>■ *T3.3 - User Interface*<br>■ *T3.4 - Explanations for affected persons* | Yes, there is an instance that:<br>■ has enough Information and power to give individualized Information<br>■ can enforce reviews of system characteristics | Yes, there is an instance that:<br>■ has enough Information and power to give individualized Information | Yes, but only standardized Information can be given. | | | | No |
| **A2.2** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **How "easy" is the access to the feedback channels?** | | Everyone with a justified interest can contact the feedback channel and receives direct feedback. | Only users can contact the feedback channels. | Access is only possible when fulfilling certain requirements:<br>e.g. additional payment or only at predefined points during the lifecycle. | | | Only for the competent authorities. | No |

**Table A3 – Responsible Human Oversight**

| A.3 | Responsible Human Oversight | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **A3.1** | Skippable | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| *Skippable if not HIC -* **Human in Command (Control):** **Is the user expertise needed to judge the results of the AI system to avoid overconfidence defined?** | | The level of expertise required for a human user to **understand and judge the system's recommendations**, given the data and body of knowledge in the field, is documented and appropriate with regard to the intended purpose. | The level of expertise that is required by a human operator to check the **plausibility** of the recommendations of the system given the data and the body of knowledge of that domain are documented. | The expertise that is required by a human operator to work successfully with the system are documented. | The level of expertise is not specified but descriptions about it can be requested and obtained. | | | No measures have been taken. |
| **A3.2** | Skippable | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| *Skippable if not HIC -* **Human in Command (Control):** **Which effort is needed to understand and interact with the AI system?** **(depending on the application context)** | | The intended user understands the actions of the AI system/application with no effort and knows how to interact with it **immediately.** | The intedned user understands the actions of the AI system/application with little effort and knows how to interact with it **after a short time period.** | Guided introduction of the AI system / application is needed. | Extended training is needed. | | | Extended training and a prior knowledge in the application context of the AI system is needed. |

| A3.3 | Skippable | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| *Skippable if not HITL -* **Human in the Loop: Which measures are taken to ensure that the AI system does not affect human autonomy by interfering with the operator's decision-making process in an unintended way?** | | The level of human control and involvement in the decision is documented and the interfaces is designed to **allow the operator to easily understand and influence** the decision process. Proposal of the system and human decision are documented. | The level of human control and involvement in the decision is documented and the interfaces is designed to **allow the operator to influence** the decision process. | The level of human control and involvement in the decision is documented in a way that allows auditing. | The level of human control and involvement in the decision is documented. | | | None of the above. |
| A3.4 | Skippable | A | B | C | D | E | F | G |
| *Skippable if not HITL -* **Human in the Loop: Is the human takeover of the system designed so that the user understands the current state of the system and can therefore take over quickly?** | | The system has been designed in a way that a seamless handover to a human operator is always possible. This handover has been tested within the full operations envelope of the system and properly documents | The system has been designed in a way that a seamless handover to a human operator is always possible. This handover has been tested and documented in the main modes of operation. | The system has been designed in a way that a seamless handover to a human operator is always possible. This handover has only been tested sporadically. | The system has been designed in a way that a seamless handover to a human operator is always possible. This handover was not tested. | | | System was not designed and has not been tested regarding the handover to a human operator. |

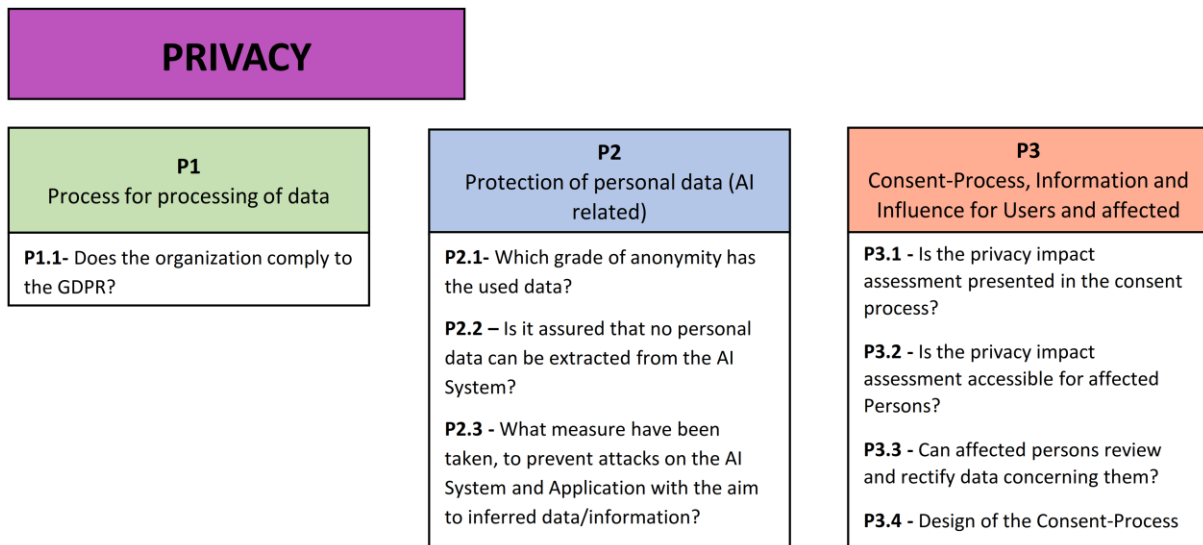| A3.5 | Skippable | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Skippable if not HOL - Human on the Loop: Does the system makes the decision parameters transparent and allows post-hoc changes?** | | Decision parameters are well documented and transparent.<br><br>There is a process and tools in place that allow documentation and ex-post explanation of decisions taken by the system. Process parameters can be changed (System can be retrained by customer). | Decision parameters are well documented and transparent.<br><br>Decisions taken by the system are well documented.<br>Process parameters can be changed (System can be retrained by customer). | Decision parameters are transparent and can be changed (System can be retrained by customer). | Process in place to change decision parameters by developer. | | | No |

## 4.2.4 Privacy

| PRIVACY | | |
|---|---|---|
| **P1**<br>Process for processing of data | **P2**<br>Protection of personal data (AI related) | **P3**<br>Consent-Process, Information and Influence for Users and affected |
| **P1.1-** Does the organization comply to the GDPR? | **P2.1-** Which grade of anonymity has the used data?<br><br>**P2.2 –** Is it assured that no personal data can be extracted from the AI System?<br><br>**P2.3 -** What measure have been taken, to prevent attacks on the AI System and Application with the aim to inferred data/information? | **P3.1 -** Is the privacy impact assessment presented in the consent process?<br><br>**P3.2 -** Is the privacy impact assessment accessible for affected Persons?<br><br>**P3.3 -** Can affected persons review and rectify data concerning them?<br><br>**P3.4 -** Design of the Consent-Process |

**Figure 4 – Composition of Privacy**

## Table P1 – Process for processing of data

| P1 | Process for processing of data | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| P1.1 | | A | B | C | D | E | F | G |
| Does the organization comply to the GDPR? | | Yes | | | | | | No |

## Table P2 – Protection of personal data (AI related)

| P2 | Protection of personal data (AI related) | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| P2.1 | positive Anchor | A | B | C | D | E | F | G |
| Which grade of anonymity has the used data? | *How anonymous is the dataset, i.e. How much personal information can be inferred from the dataset(s).* | It is justified, that no personal data exists in the dataset. | Dataset has been sanitised (removal of all data directly identifing a natural person) and anonymised. | dataset anonymised with state-of-the-art methods (e.g. k--anonymization, differential privacy, etc.). | Pseudonymization | | | No measures taken. |

| P2.2 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is it assured that no personal data can be extracted from the AI System?** | *How does the AI system development and training affect the privacy of the dataset? How much information can be inferred from the AI system / application as a result of the chosen development process? E.g. via model inversion, membership inference attack, etc.* | Yes, it is assured that no personal data can be extracted, and the approach taken is justified in a report. | | It is justified, that personal data can only be extracted with high effort. | | | Personal data can easily be extracted from the AI system. A public report explains to users which kind of data can be extracted and what the risk of extraction is. | No measures taken. |

| P2.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **What measure have been taken, to prevent attacks on the AI System and Application with the aim to inferred data/information?** | *Deployed/Live-system privacy preserving mechanism and privacy attack mitigation process* | Analysis of the potential attacks Evaluation of the possible risk. Countermeasures to mitigate the risks.<br><br>Strict Access Control (no direct access to AI system or AI application for the user)<br><br>Penetration tests have been conducted. | Analysis of the potential attacks Evaluation of the possible risk. Countermeasures to mitigate the risks.<br><br>Access Control<br><br>Penetration tests have been conducted. | Analysis of the potential attacks Evaluation of the possible risk. Countermeasures to mitigate the risks.<br><br>Access Control | Analysis of the potential attacks Evaluation of the possible risk. Countermeasures to mitigate the risks. | Analysis of the potential attacks Evaluation of the possible risk. | | No measures taken. |

## Table P3 – Consent-Process, Information and Influence for Users and affected Persons

| P3 | Consent-Process, information and influence for users and affected persons | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **P3.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Is the privacy impact assessment presented in the consent process?** | *The assessment needs to include the following aspects:*<br>■ *What might be a concrete physical impact when working with the system?*<br>■ *What kind of moral harm can be caused by such a system?*<br>■ *What could be material consequences of such a system?* | Yes, individualized for different user groups. | | Yes, but in a general way. | | | | No |
| **P3.2** | Skippable | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Is the privacy impact assessment accessible for affected Persons?** | *If there is a system that is allowed to take data from non-users due to a permission status (Erlaubnistatbestand), can they still see the impact analysis?* | Yes, individualized for different groups of affected persons. | | Yes, but in a general way. | | | | No |

| P3.3 | Skippable | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Can affected persons review and rectify data concerning them?** | *If there is a system that is allowed to collect data from non-users due to a permission status, can they still see what data has been collected from them and have the opportunity to correct it?* | There is a possibility that allows affected persons to easily review and rectify data concerning them.<br><br>Appropriate information (e.g. sign) that data is collected is available.<br><br>Contact possibilities are indicated. | There is a possibility for affected persons to review and rectify data concerning them on request.<br><br>Appropriate information (e.g. sign) that data is collected is available.<br><br>Contact possibilities are indicated. | | No possibility to review or rectify data<br><br><br>Appropriate information (e.g. sign) that data is collected is available.<br><br>Contact possibilities are indicated. | Appropriate information (e.g. sign) that data is collected is available. | | No |

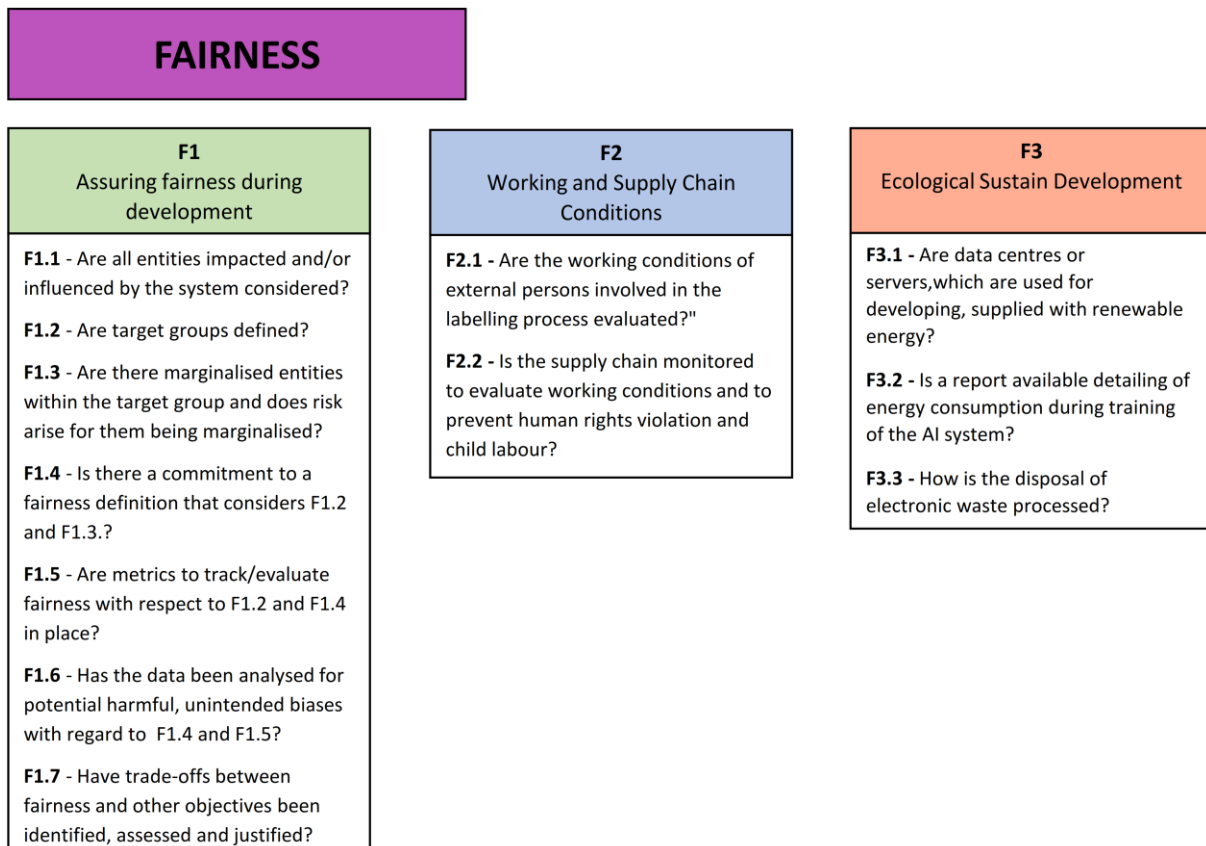| P3.4 | Skippable | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Design of the Consent-Process** | | The consent process was examined and identified for possible nudging or persuasive computing effects from a psychological and sociological perspective.<br><br>The effects were communicated to:<br>■ users<br>■ affected persons<br><br>Appropriate precautionary measures were taken regarding:<br>■ users<br>■ affected persons.<br><br>Privacy by default (opt-in for usage of personal data is needed) for:<br>■ users<br>■ affected persons | The consent process was examined and identified for possible nudging or persuasive computing effects from a psychological and sociological perspective.<br><br>The effects were communicated to:<br>■ users<br>■ affected persons<br><br>Appropriate precautionary measures were taken with regard to:<br>■ users<br><br>Privacy by default (opt-in for usage of personal data is needed) for:<br>■ users | The consent process was examined and identified for possible nudging or persuasive computing effects from a psychological and sociological perspective.<br><br>The effects were communicated to:<br>■ users<br>■ affected persons | The consent process was examined and identified for possible nudging or persuasive computing effects from a psychological and sociological perspective.<br><br>The effects were communicated to:<br>■ users | The consent process was examined and identified for possible nudging or persuasive computing effects from a psychological and sociological perspective. | | An examination has not taken place. |

## 4.2.5    Fairness



**FAIRNESS**

| F1 Assuring fairness during development | F2 Working and Supply Chain Conditions | F3 Ecological Sustain Development |
|---|---|---|
| **F1.1** - Are all entities impacted and/or influenced by the system considered?<br><br>**F1.2** - Are target groups defined?<br><br>**F1.3** - Are there marginalised entities within the target group and does risk arise for them being marginalised?<br><br>**F1.4** - Is there a commitment to a fairness definition that considers F1.2 and F1.3.?<br><br>**F1.5** - Are metrics to track/evaluate fairness with respect to F1.2 and F1.4 in place?<br><br>**F1.6** - Has the data been analysed for potential harmful, unintended biases with regard to  F1.4 and F1.5?<br><br>**F1.7** - Have trade-offs between fairness and other objectives been identified, assessed and justified? | **F2.1 -** Are the working conditions of external persons involved in the labelling process evaluated?"<br><br>**F2.2 -** Is the supply chain monitored to evaluate working conditions and to prevent human rights violation and child labour? | **F3.1 -** Are data centres or servers,which are used for developing, supplied with renewable energy?<br><br>**F3.2 -** Is a report available detailing of energy consumption during training of the AI system?<br><br>**F3.3 -** How is the disposal of electronic waste processed? |

**Figure 5 – Composition of Fairness**

## Table F1 – Assuring fairness during development

| F1 | Assuring fairness during development | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **F1.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Are all entities impacted and/or influenced by the system considered?** | *Considers the impact on the environment in which the system is to be used, including the human agency concept, implementation in the socio-technical system, …* | Yes, entities are considered that:<br>1. directly interface with the AI system **(1st order network effects)**<br>2. affected by the deployment of this system **(2nd order network effects)**<br>3. entities within the broader social-technical system it operates within | Yes, entities are considered that:<br>1. directly interface with the AI system **(1st order network effects)**<br>2. affected by the deployment of this system **(2nd order network effects)** | | Yes, entities are considered that:<br>1. directly interface with the AI system **(1st order network effects)** | | | No |
| **F1.2** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Are target groups defined?** | *Important characteristics to define different target groups:*<br>*Demography (age, income, family size, family status, gender, education, …)*<br>*Geography (residence, origin, …)* | Yes, <u>all of</u><br>1. the target groups, entities and users are identified<br>2. and a justification for the selection is provided<br>3. including additional target groups that arise from reasonably unforeseen misuses | Yes<u>, all of</u><br>1. the target groups, entities and users are identified<br>2. and a justification for the selection is provided | Yes, <u>most of</u><br>1. the target groups, entities and users are identified<br>2. and a justification for the selection is provided | Yes, <u>some of</u><br>1. the target groups, entities and users are identified | | | No |

| F1.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are there marginalised entities within the target group and does risk arise for them being marginalised?** | | Yes,<br>1. A research to identify all marginalised groups/entities is carried out.<br><br>2. All risk that arises from any of the groups being marginalised are detailed and justified.<br><br>3. The marginalised groups/entities are involved in the development process. | Yes,<br>1. A research to identify all marginalised groups/entities is carried out.<br><br>2. All risk that arises from any of the groups being marginalised are detailed and justified. | Yes,<br>1. A research to identify all affected marginalised groups/entities is carried out. | | | | No |

| F1.4 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is there a commitment to a fairness definition that considers F1.2 and F1.3.?** | | A fairness definition is defined and provided in **collaboration with the target group and marginalised entities/groups.**<br><br>There is a commitment to considering it **throughout the lifecycle of the AI system**.<br><br>There is a commitment to a process to validate and ensure the integrity of the fairness criteria throughout the life cycle of the AI system.<br><br>Easy access and transparency to the fairness definition and criteria is provided to the public, including justification for the definition and process. | A fairness definition is provided and considered **throughout the lifecycle of the AI system**.<br><br><br>There is a commitment to a process to validate and ensure the integrity of the fairness criteria throughout the life cycle of the AI system.<br><br>Easy access and transparency to the fairness definition and criteria is provided to the public, including justification for the definition and process. | A fairness definition is provided and considered **throughout the lifecycle of the AI system**.<br><br><br><br>Easy access and transparency to the fairness definition and criteria is provided to the public, including justification for the definition and process. | A fairness definition is provided and considered **within the development process**. | A fairness definition is provided. | | No |

| F1.5 | | A | B | C | D | E | F | G |
|------|---|---|---|---|---|---|---|---|
| **Are metrics to track/evaluate fairness with respect to F1.2 and F1.4 in place?** | *Typical fairness metrics are e.g., statistical parity, equal distribution of false negatives, equal distribution of false positives, decision between group fairness or individual fairness, fairness through unawareness, equality of opportunity.* | Several metrics to measure and track fairness are in place, that completely align with the fairness definition and criteria throughout the whole lifecycle of the AI system.<br><br>These metrics have been developed in collaboration with the target group and the marginalised users/entities.<br><br>The metrics are well documented and transparent to the public. | Several metrics to measure and track fairness are in place, that completely align with the fairness definition and criteria throughout the whole lifecycle of the AI system.<br><br>These metrics have been developed in collaboration with the target group and the marginalised users/entities. | Several metrics to measure and track fairness are in place, that completely align with the fairness definition and criteria throughout the whole lifecycle of the AI system. | Several metrics to measure and track fairness are in place. | A metric to measure and track fairness is in place. | | No |

| F1.6 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Has the data been analysed for potential harmful, unintended biases with regard to F1.4 and F1.5?** | | A datasheet is provided. It documents which data sources have been assessed and with which methods in order to identify biases that might bring harm or risk.<br><br>The documentation covers the objectives and measures taken to avoid harm and risk. It also states why the actions taken are reasonable in relation to the selected fairness metric. The documentation is released to the public.<br><br>The nature of the bias has been ascertained. It was considered and assessed with respect to the fairness definition and criteria. | A datasheet is provided. It documents which data sources have been assessed and with which methods in order to identify biases that might bring harm or risk.<br><br>The documentation covers the objectives and measures taken to avoid harm and risk. The documentation is released to the public.<br><br>The nature of the bias has been ascertained. It was considered and assessed with respect to the fairness definition and criteria. | A datasheet is provided. The documentation covers the objectives and measures taken to avoid harm and risk.<br><br><br><br>The nature of the bias has been ascertained. It was considered and assessed with respect to the fairness definition and criteria. | The nature of the bias has been ascertained. It was considered and assessed with respect to the fairness definition and criteria. | The data has been analysed for the most common and easily identifiable biases. This has been documented and considered with respect to the fairness definition and criteria. | | No |

| F1.7 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Have trade-offs between fairness and other objectives been identified, assessed, and justified?** | *Identification, assessment, and justification according to the target group and marginalised users/entities.*<br><br>*Possible objectives can be performance or privacy.* | Trade-offs have been identified and documented in collaboration with the target group and marginalised groups.<br><br>Consideration of how to balance any trade-off involved collaboration with or are based on the feedback of the target group and marginalised groups.<br><br>This process is well documented and accessible to the target group and marginalised groups. | Trade-offs have been identified and documented in collaboration with the target group and marginalised groups.<br><br>Consideration of how to balance any trade-off involved collaboration with or are based on the feedback of the target group and marginalised groups. | | Trade-offs have been identified and document in collaboration with the target group and marginalised users/entities. | Trade-offs have been identified and document. | | No |

## Table F2 – Working and Supply Chain Conditions

| F2 | Working and Supply Chain Conditions | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **F2.1** | **Skippable** | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| *Skippable if no external participation* **Are the working conditions of external persons involved in the labelling process evaluated?** | *Minimal safety and worker protection standards and standards regarding social security and protection from exploitation in place at the facility providing click work are covered by the Supply Chain Act, for example.* | Yes, the <u>following conditions</u> are evaluated:<br>■ minimal safety and worker protection standards<br>■ minimal standards regarding social security and protection from exploitation<br>■ click work-specific working conditions (diversification of tasks, potential emotional/psychological dangers from explicit material)<br><br>The datasets (e.g. datasheet) contain information about labelling (click working) process.<br><br>It is published, that external persons are involved. | Yes, the following conditions are evaluated:<br>■ minimal safety and worker protection standards<br>■ minimal standards regarding social security and protection from exploitation<br><br><br>The datasets (e.g. datasheet) contain information about labelling (click working) process.<br><br>It is published, that external persons are involved. | Yes, **one of the two** <u>following conditions</u> are evaluated:<br>■ minimal safety and worker protection standards<br>■ minimal standards regarding social security and protection from exploitation<br><br>The datasets (e.g. datasheet) contain information about labelling (click working) process.<br><br>It is published, that external persons are involved. | The datasets (e.g. datasheet) contain information about labelling (click working) process.<br><br>It is published, that external persons are involved. | It is published, that external persons are involved. | | No, there is no evaluation or documentation. |

| F2.2 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is the supply chain monitored to evaluate working conditions and to prevent human rights violation and child labour?** | *Examples for legal requirements:*<br>■ *German Act on Corporate Due Diligence Obligations in Supply Chains*<br>■ *EU Corporate Due Diligence Obligations in Supply Chains* | Yes, the supply chain is consequently monitored due to existing legal requirements or similar obligations. | | The suppliers are reviewed once. There is a company specific policy to handle violations. | The suppliers are reviewed once. | | | No |

### Table F3 – Ecological Sustain Development

| F3 | Ecological Sustain Development | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| F3.1 | | A | B | C | D | E | F | G |
| **Are data centres or servers, which are used for developing, supplied with renewable energy?** | *Renewable energy includes solar, wind, hydro, geothermal, biomass and marine energy. Climate positivity can be reached through e.g., carbon offsetting or reuse of excess energy for heating.* | Yes, at least 99% and the data centres are climate positive. | Yes, at least 99 %. | In part, at least 80%. | In part, at least 60%. | In part, at least 40%. | In part, at least 20%. | **Less than 20 %.** |

| F3.2 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Is a report available detailing of energy consumption during training of the AI system?** | *The report has to include a calculation or estimates of energy consumption and carbon emissions of all system components or the system overall, measures for carbon offsetting or energy reuse, a description of other ecological impact incl. directly resulting waste generation, an explanation of the process of consideration between its ecological impact (incl. energy consumption) and other factors (e.g. accuracy), and an explanation of why the chosen AI model is used with regards to its ecological sustainability.* | Yes, and this report was published before or during launch of the system with estimates and was updated after the launch with actual energy use calculations. | Yes, and this report was published after the system was already in use. | Yes, and this report was published before or during launch of the system with estimate.s | | A report was written but is only available internally. | | No |

| F3.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **How is the disposal of electronic waste processed?** | *Does not include the waste generated by the user.* | It is ensured that electronic waste is recycled as far as possible and not exported to risk areas (including Supply Chain). | It is ensured that electronic waste is recycled as far as possible and not exported to risk areas (just for electronics that are in control of the organisation). | | | | | **There is no detailed knowledge.** |

## 4.2.6    Reliability



**Figure 6 – Composition of Reliability**

## Table R1 – Robustness & Reliability qua Design

| R1 | Robustness & Reliability qua Design | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **R1.1** | | **A** | **B** | **C** | **D** | **E** | **F** | **G** |
| **Is the operational design domain of the AI system / application clearly defined and documented?** | *The Operational Design Domain (ODD) describes the conditions and environment an AI system/application is intended to operate within, and reasonably expected to encounter. This ODD should be described accurately and in enough detail such that the environment and boundaries of operation are clear. The user and stakeholders should be able to easily deduce from this whether the planned/intended use of an AI system is within the scope of the ODD.* | A **onthologically complete, structured and detailed** description of the:<br>■ operational design domain<br>■ and the intended use cases<br><br>These are published and well understood by:<br>■ the users of the AI system<br>■ auditors<br>■ regulatory bodies<br>■ all additional stakeholders | A **onthologically complete, structured and detailed** description of the:<br>■ operational design domain<br>■ and the intended use cases<br><br>These are published and well understood by:<br>■ the users of the AI system<br>■ auditors<br>■ regulatory bodies | A description of the:<br>■ operational design domain<br>■ and the intended use cases<br><br>These are published and well understood by:<br>■ the users of the AI system,<br>■ auditors<br>■ regulatory bodies | A description of the:<br>■ operational design domain<br>■ and the intended use cases<br><br>These are published and well understood by:<br>■ the users of the AI system. | A description of the:<br>■ the intended use cases<br><br>These are published and well understood by:<br>■ the users of the AI system | | No |

| R1.2 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Was ensured, that the quality and quantity of the data fit to the intended purpose and Operational Design Domain?** | | Documentation of which shows, the examination of:<br><br>■ Completeness of the attributes of the data<br>■ Correctness of data<br>■ data format<br>■ the labeling and Annotation Process including quality assurance<br>■ compatibility of Training data with the operational design domain<br>■ relevant data preparation; i.e. raw data pre-processing (e.g. cleaning, enrichment, aggregation) with regard to the intended purpose and Operational Design Domain of the AI System | Documentation of which shows, the examination of:<br><br>■ Completeness of the attributes of the data<br>■ Correctness of data<br>■ data format<br>■ the labeling and Annotation Process including quality assurance<br>■ compatibility of Training data with the operational design domain<br>■ relevant data preparation; i.e. raw data pre-processing (e.g. cleaning, enrichment, aggregation) | | Documentation of/ which shows, the examination of:<br><br>■ Completeness of the attributes of the data<br>■ Correctness of data<br>■ data format<br>■ the Labeling and Annotation Process including quality assurance<br><br>■ relevant data preparation; i.e. raw data pre-processing (e.g. cleaning, enrichment, aggregation) | Documentation of/ which shows, the examination of:<br><br>■ Completeness of the attributes of the data<br>■ Correctness of data<br>■ data format<br>■ the Labeling and Annotation Process including quality assurance | Documentation of/ which shows, the examination of:<br><br>■ Completeness of the attributes of the data<br><br><br>■ data format | No |

| R1.3 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Was the quality of the development of the AI systems ensured?** | | Justification of the approach and models used.<br><br>With documentation and justification of the chosen:<br>■ Performance and Evaluation Metrics<br>■ Optimization metric<br>■ The testing strategy<br><br>Live testing that covers the ODD and reasonably foreseen situations of the OD has been performed. | Justification of the approach and models used.<br><br>With documentation and justification of the chosen:<br>■ Performance and Evaluation Metrics<br>■ Optimization metric<br>■ The testing strategy<br><br>Live testing that covers the ODD has been performed. | Justification of the approach and models used.<br><br>With documentation and justification of the chosen:<br>■ Performance and Evaluation Metrics<br>■ Optimization metric<br>■ The testing strategy<br><br>(Virtual) Testing inside the ODD has been performed. | Justification of the approach and models used.<br><br>With documentation and justification of the chosen:<br>■ Performance and Evaluation Metrics<br>■ Optimization metric<br><br>(Virtual) Testing inside the ODD has been performed. | Justification of the approach and models used.<br><br>With documentation and justification of the chosen:<br>■ Performance and Evaluation Metrics<br>■ Optimization metric | Justification of the approach and models used. | no |
| R1.4 | | A | B | C | D | E | F | .G |
| **Is the system robust against varying environments (i.e. distribution shift) and outliers?** | *Varying environments can influence the Performance of an AI system. The system needs to be able to detect varying environments to adapt his behaviour. flawed data = data that is influenced by a statistical or non-statistical disturbance or malfunction. E.g. rain, dust, lens effects, noise.* | System must be able to gracefully track and monitor changes in the operational environment. It must offer mechanisms to adapt to observed changes in the operational design domain. | Relaxation of Grade A: reasonably adhere to changes in the operational design domain. | Yes, but only in a subdomain. | | | | No |

| R1.5 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are all possible risks assesed and the harms the system could have classified (e.g. life and health, violation of rights etc.)?** | | All risks are transparent, well documented with the product and made available to customers. | All risks are transparent and can be obtained by a defined interface. | Main risks are identified and can be retrieved by a defined process. | | | | None of the above. |

| R1.6 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are measures in place to ensure the integrity, robustness, and overall security of the AI system / application against potential attacks over its life cycle?** | *Implementation of general cybersecurity measures.* | Compliance to cybersecurity standards (e.g. ISO 27k series, IEC 62443, ISO/SAE 21434, ETSI EN 303 645, ... ). Regular review security measures and protocols. Measures (including the ones taken during training of AI system) are defined and transparently documented with the product. | Measures are defined and transparently documented with the product. | Measures defined information can be retrieved by a defined interface. | Measures partly defined and information can be retrieved by a defined process. | | | None of the above. |

| R1.7 | | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|---|
| **Are end-users informed of the duration of security coverage and updates? What length is the expected timeframe within which security updates for the AI system/application will be provided?** | | Information is shipped with the product. | Information can be obtained by a defined interface. | Information partially available and can be retrieved by a defined process. | | | | None of the above. |
| R1.8 | | A | B | C | D | E | F | G |
| **Are technical documentations documented, including standards, that need to be applied by the AI system/application?** | | Yes | | | | | | No |

## Table R2 – Robustness & Reliability in Operation

| R2 | Robustness & Reliability in Operation | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **R2.1** | | A | B | C | D | E | F | G |
| **Is the applied AI lifecycle management robust to changes in the operational domain?** | | Continuous model monitoring and testing (including integrity checks) as a feature of the AI strategy covering the full operational domain. | Continuous model monitoring and testing (including integrity checks) as a feature of the AI strategy covering key/important areas of the operational domain. | Regular model monitoring and testing (including integrity checks) as a feature of the AI strategy covering key/important areas of the operational domain. | Occasional model monitoring and testing (including integrity checks) are carried out. | Occasional model monitoring and testing is carried out. | Occasional testing is carried out. | None of the above. |
| **R2.2** | | A | B | C | D | E | F | G |
| **Is a failure mitigation strategy for the AI-based system in place?** | *Is there a fail-safe strategy for the AI-based system in place? Reaction of the system if parts of it are not working properly (such as sensors malfunctioning) or if the input data is either corrupted or contains noise. Presence of fall-back systems in case the AI-based system cannot work properly anymore e.g., broken/dirty lens/microphone, electromagnetic interference.* | Yes, the following:<br>■ redundancy,<br>■ fall back mechanisms (e.g. defaulting to a safe mode, kill-switch),<br>■ alert system (end-user, provider, competent authority),<br>■ fail-safe logging (i.e. black box),<br>■ secure failure (e.g. tamper protection, safe mode) and system restoration. | Yes, the following:<br>■ redundancy,<br>■ fall back mechanisms,<br>■ alert system,<br>■ fail-safe logging,<br>■ secure failure. | Yes, the following:<br>■ redundancy,<br>■ fall back mechanisms,<br>■ alert system,<br>■ secure failure. | Yes, the following:<br>■ fall back mechanisms,<br>■ alert system,<br>■ secure failure. | Yes, the following:<br>■ fall back mechanisms,<br>■ alert system. | Yes, the following:<br>■ fall back mechanisms,<br>■ basic alert system. | None of them. |

# 5 Determining Trust Classes

## 5.1 General

To receive a rating in the VCIO scheme, an observable for every indicator needs to be determined. The rating of the indicators subsumed under a criterion is then aggregated to a rating for the corresponding criteria. These criteria subsumed under a value are then also aggregated to arrive at the final rating for the corresponding value.

Each indicator has corresponding observables, that have a rating from "A" (best or fully fulfilled) to "G" (worst / not fulfilled at all).

There are three different types of indicators, that are described in the following sections:

- Score indicators (5.2.1)
- Positive anchor indicators (5.2.3)
- Negative anchor indicators (5.2.2)

Each of them has a different impact on the aggregation on the criteria level.

For some criteria, there exist indicators that fulfil the same aim with a different approach. These are marked as alternative indicators and are handled as one indicator. Therefore, the indicator with the higher ranked observable is chosen.

NOTE Since one alternative indicator can be chosen, just one alternative indicator is considered for aggregation.

In specific cases, it is possible that some Indicators do not apply to the rated product. In this case, the indicator and its observables can be omitted and are not counted towards the aggregation.

## 5.2 Types of Indicators

### 5.2.1 Score Indicators

Score indicators contain of observables rating from "A" to "G". Every level of the observables corresponds to a score (see Table 1 – Corresponding scores for the levels), which allows to aggregate them to the criterion level mentioned in section 5.3.1. There it is also explained how the types of indicators have different effects in the aggregation.

**Table 1 – Corresponding scores for the levels**

| Level | A | B | C | D | E | F | G |
|-------|---|---|---|---|---|---|---|
| Score | 0 | 1 | 2 | 3 | 4 | 5 | 6 |

### 5.2.2 Negative Anchor Indicators

Negative anchor indicators are necessary conditions to meet the aim of a criteria. If a negative indicator is not sufficiently fulfilled, the indicator cannot be fulfilled either. Therefore, if a negative indicator is rated with a "G", the corresponding criteria is automatically rated with "G".

If the negative anchor indicator is at least partially fulfilled, the necessary condition is meet and the negative anchor indicator is handled as a Score indicator.

### 5.2.3 Positive Anchor Indicators

Positive anchor indicators are sufficient conditions to fulfil the aspects of a criterion. Therefore, if a positive anchor indicator is rated with an "A", the corresponding criterion is rated with an "A".

### 5.2.4 Skippable Indicators

In some cases, not all indicators are applicable for the intended use of an AI system. These are marked as skippable, which means that they do not need to be answered if not applicable and are not taken into account for aggregation.

For example, it is possible to skip privacy indicators, if no personal data is acquired and therefore no individuals are affected in their privacy. Also, an AI system can just have one underlying Human Agency concept. In this case the other Human Agency indicators can be skipped.

If one or more indicators are skipped, this results in a reduced number of indicators which are considered for the aggregation. This shall be considered for the divider of the rounded average in section 5.3.1.

## 5.3    Aggregation

### 5.3.1    Criteria Level

To reach an aggregation from indicator to criterion, the following steps must be taken, depending on which kind of indicators exist in the considered criterion.

**Step 1: Check for negative anchor indicator**

If a negative anchor indicator exists in the considered criterion this step has to be taken, otherwise it can be skipped.

If the negative anchor indicator is rated with "G" and therefore not fulfilled at all, the whole criterion cannot be fulfilled and is likewise rated with a "G" and the aggregation of the criterion is terminated.

If the level of the negative anchor indicator is "F" or higher, it is at least partially fulfilled and it is treated like a score indicator.

**Step 2: Application of positive anchor indicators**

Positive anchor indicators can fully fulfil a criterion if they exist and are fully fulfilled. If they do not exist, this step can be skipped.

This means if a positive anchor indicator is marked with an "A", the corresponding criterion is also marked with an "A".

If positive anchor indicators are not fully fulfilled, this means level below "A", they are treated as score indicators.

**Step 3: Aggregation with score indicators**

Every level of the observable for an indicator corresponds with a (malus)score. Here "A" gets the lowest score and "G" the highest (see Table 1 – Corresponding scores for the levels).

To aggregate the level of the indicators to the criteria level the rounded average of the scores is taken. This means the scores are added, divided by the number of applied indicators, and are rounded to the next integer. The resulting score can then be translated back to a corresponding level.

If a positive anchor indicator for this criterion exists and is fully fulfilled and rated with an "A" in step 2, the criterion is also marked with an "A".

Example: If an indicator with an "A" gets aggregated with an indicator with a "B" the corresponding scores are "0" and "1". After adding them up this results in "1". Divided by the number of indicators (in this case 2) you receive "0.5", which then is rounded to 1 and corresponds to "B". Therefore "B" would be the result of the aggregation.

### 5.3.2    Value Level

After completing the criteria level, the aggregation on the value level can be performed. Since the criteria all have the same significance, the scores are determined by a rounded average. The scores are equivalent to the one on indicator level, and the aggregation is analogue to Step 3 for the criteria level (see Table 1 – Corresponding scores for the levels) negative or positive anchor criteria do not exist.

Therefore, to aggregate the level of the criteria to the value level the rounded average of the scores is taken. This means the scores get added, divided by the number of applied Indicators, and are rounded to the next integer. The resulting score can then be translated back to a corresponding level.

### 5.3.3    Example for Aggregation

To illustrate the aggregation on criteria and value level an example is shown in Figure 7 – Illustration of the Aggregation.
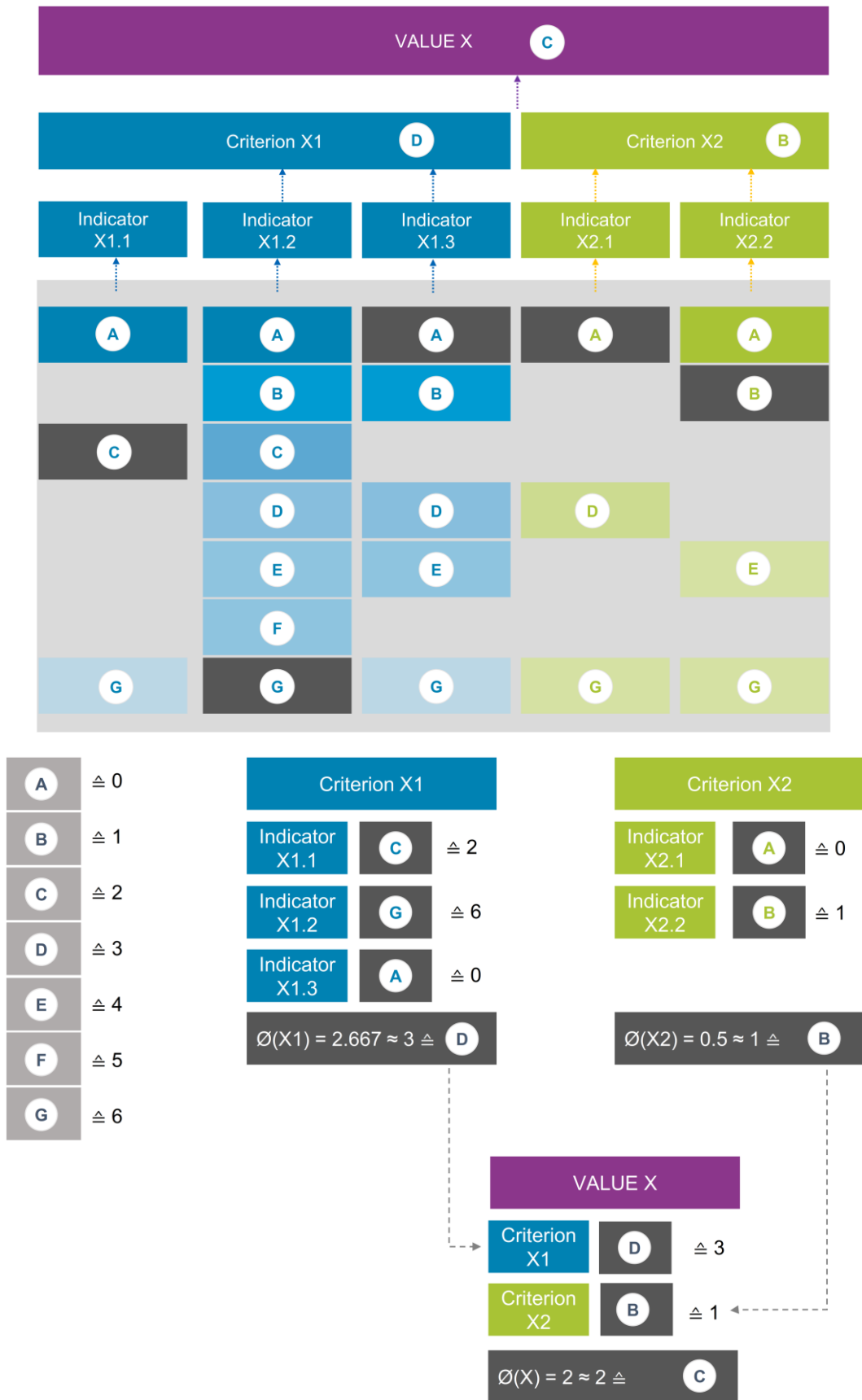
**Figure 7 – Illustration of the Aggregation**

## 5.4 Aggregation in case of multiple sub-systems

AI Systems can consist of multiple sub-items, like multiple AI models or datasets.

These sub-items could have different observables in an indicator. For example, different datasets have differently detailed corresponding information (T1.1). In this case the lowest reached observable is taken for the indicator.

## 5.5 Partial Label

The AI Trust Label is intended to show the trustworthiness and adherence to values of products. In the case of AI, the focus of the standard is on entire products that contain AI technologies.

However, these systems can be composed and integrated from several sub-items within a supply chain. For this reason, these sub-items can also fulfil indicators or criteria if they can be meaningfully assigned.

In this case, it would also be possible for the creator of this sub-item to prove the indicators or criteria including the corresponding level. These can then be used by the integrator without having to perform a further check.

# Bibliography

[1]     ARTIFICIAL INTELLIGENCE ACT – Proposal for a "REGULATION OF THE EUROPEAN
        PARLIAMENT AND OF THE COUNCIL LAYING DOWN HARMONISED RULES ON
        ARTIFICIAL INTELLIGENCE AND AMENDING CERTAIN UNION LEGISLATIVE ACTS" and
        Annexes by the European Commission, 2021-04-21.

[2]     Poretschkin, Maximilian; Schmitz, Anna; Akila, Maram; Adilova, Linara; Becker, Daniel;
        Cremers, Armin B.; Hecker, Dirk; Houben, Sebastian; Mock, Michael; Rosenzweig, Julia;
        Sicking, Joachim; Schulz, Elena; Voss, Angelika; Wrobel, Stefan: Leitfaden zur Gestaltung
        vertrauenswürdiger Künstlicher Intelligenz – „KI-Prüfkatalog", 2021.

[3]     Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino
        Pedreschi. 2018. A Survey of Methods for Explaining Black Box Models. ACM Comput. Surv.
        51, 5, Article 93 (September 2019), 42 pages. DOI: https://doi.org/10.1145/3236009

[4]     ETHICS GUIDELINES FOR TRUSTWORTHY AI by the High-Level Expert Group on Artificial
        Intelligence, Publishing date: 2019-04-08.

[5]     Report of the AIEI Group: "From Principles to Practice – An interdisciplinary framework to
        operationalise AI ethics" from AI Ethics Impact Group (AIEIG), Publishing date: 2020-04-01.

# Other sources

Algorithm Watch. "AI Ethics Guidelines Global Inventory." Accessed August 19, 2020.
https://inventory.algorithmwatch.org/.

ASAM OpenODD: Concept Paper – Version 1.0, 01.10.2021, Available at:
https://www.asam.net/index.php?eID=dumpFile&t=f&f=4544&token=1260ce1c4f0afdbe18261f7137c6
89b1d9c27576. (Accessed: 09.03.2022)

Baer, Tobias: Understand, Manage, and Prevent Algorithmic Bias: A Guide for Business Users and
Data Scientists. Apress, 2019.

Barocas, Solon; Boyd, Danah: "Engaging the Ethics of Data Science in Practice." Communications of
the ACM 60, no. 11 (2017): 23–25.

Barocas, Solon; Hardt, Moritz; Narayanan, Arvind: "Fairness and Machine Learning: Limitations and
Opportunities." Accessed May 19, 2019. https://fairmlbook.org/.

Barocas, Solon; Selbst, Andrew D: "Big Datas Disparate Impact." California Law Review 104 (2016):
671–732.

Behrendt, Hauke; Loh. Wulf: "Informed Consent and Algorithmic Discrimination: Is Giving Away Your
Data the New Vulnerable?" Review of Social Economy, 2022, online first. Algorithm Watch. "AI Ethics
Guidelines Global Inventory." Accessed August 19, 2020. https://inventory.algorithmwatch.org/.

Cave, Stephen; Dihal, Kanta: "Race and AI: The Diversity Dilemma." Philosophy & Technology 34, no.
4 (2021): 1775–79.

Citron, Danielle; Pasquale, Frank: "The Scored Society: Due Process for Automated Predictions."
Washington Law Review 89 (2014): 1–33.

Crawford, Kate; Paglen, Trevor: "Excavating AI: The Politics of Training Sets for Machine Learning."
Accessed March 9, 2021. https://excavating.ai/.

Djeffal, Christian: "Artificial Intelligence and Public Governance: Normative Guidelines for Artificial
Intelligence in Government and Public Administration." In Regulating Artificial Intelligence. Edited by
Thomas Wischmeyer and Timo Rademacher, 277–93. Cham: Springer International Publishing, 2020.

Draft Text of the Recommendation on the Ethics of Artificial Intelligence. SHS/IGM-
AIETHICS/2021/JUN/3 Rev.2. UNESCO. June 25, 2021.

Falco, Gregory; Shneiderman, Ben; Badger, Julia; Carrier, Ryan; Dahbura, Anton; Danks, David;
Eling, Martin et al: "Governing AI Safety Through Independent Audits." Nature Machine Intelligence 3,
no. 7 (2021): 566–71. https://doi.org/10.1038/s42256-021-00370-7.

Fjeld, Jessica; Nagy, Adam: "Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches to Principles for AI." BKC Research Publication 1/2020, January 15, 2020. https://cyber.harvard.edu/publication/2020/principled-ai.

Floridi, Luciano; Cowls, Josh; Beltrametti, Monica; Chatila, Raja; Chazerand, Patrice; Dignum, Virginia; Luetge, Christoph et al: "AI4People — An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations." Minds and Machines 28, no. 4 (2018): 689–707.

General Data Protection Regulation – REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC.

Gebru, Timnit; Morgenstern, Jamie; Vecchione, Briana; Wortman Vaughan, Jennifer; Wallach, Hanna; Daumé III, Hal; Crawford, Kate: Datasheets for datasets. 2018

GESIS. "Guidelines for Using Automated Tools for Gender Inferences." Accessed June 9, 2021. http://193.175.238.89/Gender_Inference/guidelines.html.

Hacker, Peter: "Teaching Fairness to Artificial Intelligence: Existing and Novel Strategies Against Algorithmic Discrimination Under EU Law." Common Market Law Review 55, no. 4 (2018): 1143–85.

Hedden, Brian: "On Statistical Criteria of Algorithmic Fairness." Philosophy & Public Affairs 49, no. 2 (2021): 209–31.

Heesen, Jessica; Müller-Quade, Jörn; Wrobel, Stefan: "Kritikalität von KI-Systemen in ihren jeweiligen Anwendungskontexten: Impulspapier." Plattform Lernende Systeme, 2021.

Heesen, Jessica; Müller-Quade Jörn; Wrobel, Stefan; Poretschkin, Maximilian; Dachsberger, Stephanie; Hösl, Maximilian: "Zertifizierung von KI-Systemen: Impulspapier." Plattform Lernende Systeme, 2020. https://www.plattform-lernende-systeme.de/files/Downloads/Publikationen/AG3_Impulspapier_290420.pdf.

Hellman, Deborah: "Measuring Algorithmic Fairness." Virginia Law Review 106 (2020).

Jobin, Anna; Ienca, Marcello; Vayena, Effy: "The Global Landscape of AI Ethics Guidelines." Nature Machine Intelligence 1, no. 9 (2019): 389–99. https://doi.org/10.1038/s42256-019-0088-2.

Madaio, Michael; Stark, Luke; Wortman Vaughan, Jennifer; Wallach, Hanna: Co-Designing Checklists to Understand Organizational Challenges and Opportunities around Fairness in AI. In Conference on Human Factors in Computing Systems. 2020.

Mann, Monique; Matzner, Tobias: "Challenging Algorithmic Profiling: The Limits of Data Protection and Anti-Discrimination in Responding to Emergent Discrimination." Big Data & Society 6, no. 2 (2019): 205395171989580. https://doi.org/10.1177/2053951719895805.

Mitchell, Margaret; Wu, Simone; Zaldivar, Andrew; Barnes, Parker; Vasserman, Lucy; Hutchinson, Ben; Spitzer, Elena; Raji, Inioluwa Deborah; Gebru, Timnit: Model Cards for Model Reporting. In CoRR. 2018.

Mittelstadt, Brent: "Principles Alone Cannot Guarantee Ethical AI." Nature Machine Intelligence 1, no. 11 (2019): 501–7.

Neyland, Daniel: "Bearing Account-Able Witness to the Ethical Algorithmic System." Science, Technology, & Human Values 41, no. 1 (2016): 50–76. https://doi.org/10.1177/0162243915598056.

Noble, Safiya Umoja: Algorithms of Oppression: How Search Engines Reinforce Racism. New York: New York University Press, 2018.

Sanders, Trooper: "Testing the Black Box: Institutional Investors, Risk Disclosure, and Ethical AI." Philosophy & Technology 34, S1 (2021): 105–9.

Schwartz, Reva; Down, Leann; Jonas, Adam; Tabassi, Elham: "A Proposal for Identifying and Managing Bias in Artificial Intelligence." National Institute of Standards and Technology (NIST), 2021. https://doi.org/10.6028/NIST.SP.1270-draft.

Vollmer, Sebastian; Mateen, Bilal A.; Bohner, Gergo; Király, Franz J.; Ghani, Rayid; Jonsson, Pall; Cumbers, Sarah et al: "Machine Learning and Artificial Intelligence Research for Patient Benefit: 20 Critical Questions on Transparency, Replicability, Ethics, and Effectiveness." BMJ 368 (2020): 1-12.

Wachter, Sandra; Mittelstadt, Brent: "A Right to Reasonable Inferences: Re-Thinking Data Protection Law in the Age of Big Data and AI." Columbia Business Law Review, no. 2 (2019): 494–620.

Yeung, Karen; Howes, Andrew; Pogrebna, Ganna: "AI Governance by Human Rights-Centred Design, Deliberation and Oversight: An End to Ethics Washing." In The Oxford Handbook of AI Ethics. Edited by M. Dubber and Frank Pasquale. Oxford UK: Oxford Univ. Press, 2019.

**VDE**