



FINAL DRAFT Technical Specification

ISO/IEC DTS 6254

Information technology — Artificial intelligence — Objectives and approaches for explainability and interpretability of machine learning (ML) models and artificial intelligence (AI) systems

ISO/IEC JTC 1/SC 42

Secretariat: **ANSI**

Voting begins on:
2025-03-20

Voting terminates on:
2025-05-15

Document Preview

[ISO/IEC DTS 6254](https://standards.iteh.ai/catalog/standards/iso/b027f558-bf17-4bc0-9880-2215d3dabd84/iso-iec-dts-6254)

<https://standards.iteh.ai/catalog/standards/iso/b027f558-bf17-4bc0-9880-2215d3dabd84/iso-iec-dts-6254>

RECIPIENTS OF THIS DRAFT ARE INVITED TO SUBMIT, WITH THEIR COMMENTS, NOTIFICATION OF ANY RELEVANT PATENT RIGHTS OF WHICH THEY ARE AWARE AND TO PROVIDE SUPPORTING DOCUMENTATION.

IN ADDITION TO THEIR EVALUATION AS BEING ACCEPTABLE FOR INDUSTRIAL, TECHNOLOGICAL, COMMERCIAL AND USER PURPOSES, DRAFT INTERNATIONAL STANDARDS MAY ON OCCASION HAVE TO BE CONSIDERED IN THE LIGHT OF THEIR POTENTIAL TO BECOME STANDARDS TO WHICH REFERENCE MAY BE MADE IN NATIONAL REGULATIONS.

iTeh Standards
(<https://standards.iteh.ai>)
Document Preview

ISO/IEC DTS 6254

<https://standards.iteh.ai/catalog/standards/iso/b027f558-bf17-4be0-9880-2215d3dabd84/iso-iec-dts-6254>



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2025

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

Page

Foreword	v
Introduction	vi
1 Scope	1
2 Normative references	1
3 Terms and definitions	1
4 Symbols and abbreviated terms	5
5 Overview	6
6 Stakeholders' objectives	6
6.1 General	6
6.2 AI user	7
6.3 AI developer	7
6.4 AI product or service provider	7
6.5 AI platform provider	8
6.6 AI system integrator	8
6.7 Data provider	8
6.8 AI evaluator	8
6.9 AI auditor	8
6.10 AI subject	8
6.11 Relevant authorities	8
6.11.1 Policy makers	8
6.11.2 Regulators	8
6.11.3 Other authorities	9
7 Explainability considerations throughout the AI system life cycle	9
7.1 General	9
7.2 Inception	10
7.3 Design and development	10
7.3.1 General	10
7.3.2 Development of the explainability component	10
7.3.3 Explainability's contribution to development	11
7.4 Verification and validation	11
7.4.1 General	11
7.4.2 Evaluation of the explainability component	11
7.4.3 Explainability's contribution to evaluation	13
7.5 Deployment	14
7.5.1 General	14
7.5.2 Deployment of the explainability component	14
7.5.3 Explainability's contribution to deployment	14
7.6 Operation and monitoring	14
7.7 Continuous validation	14
7.8 Re-evaluation	14
7.9 Retirement	15
8 Property taxonomy of explainability methods and approaches	15
8.1 General	15
8.2 Properties of explanation needs	16
8.2.1 General	16
8.2.2 Expertise profile of the targeted audience	16
8.2.3 Frame activity of interpretation or explanation	17
8.2.4 Scope of information	17
8.2.5 Completeness	17
8.2.6 Depth	18
8.2.7 Reasoning path	18
8.2.8 Implicit and explicit explanations	19

8.3	Forms of explanation	19
8.3.1	General	19
8.3.2	Numeric	19
8.3.3	Visual	19
8.3.4	Textual	20
8.3.5	Structured	20
8.3.6	Example-based	20
8.3.7	Interactive exploration tools	20
8.4	Technical approaches towards explainability	20
8.4.1	General	20
8.4.2	Empirical analysis	21
8.4.3	Post-hoc interpretation	21
8.4.4	Inherently interpretable components	21
8.4.5	Architecture- and task-driven explainability	22
8.5	Technical constraints of the explainability method	22
8.5.1	General	22
8.5.2	Genericity of the method	22
8.5.3	Transparency requirements	23
8.5.4	Display requirements	23
9	Approaches and methods to explainability	23
9.1	General	23
9.2	Empirical analysis methods	24
9.2.1	General	24
9.2.2	Fine-grained evaluation	25
9.2.3	Error analysis	25
9.2.4	Analysis-oriented datasets	25
9.2.5	Ablation	26
9.2.6	Known trends	26
9.3	Post hoc methods	27
9.3.1	Local	27
9.3.2	Global	32
9.4	Inherently interpretable components	36
9.4.1	General	36
9.4.2	Legible models	37
9.4.3	Meaningful models	39
9.4.4	Models with explicit knowledge	41
9.5	Architecture- and task-driven methods	43
9.5.1	General	43
9.5.2	Informative features	43
9.5.3	Rich and auxiliary inputs	44
9.5.4	Multi-step processing	44
9.5.5	Rich outputs	45
9.5.6	Rationale-based processing	46
9.5.7	Rationale generation as auxiliary output	46
9.6	Data explanation	47
	Annex A (informative) Extent of explainability and interaction with related concepts	48
	Annex B (informative) Illustration of methods' properties	51
	Annex C (informative) Concerns and limitations	61
	Bibliography	65

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives or www.iec.ch/members_experts/refdocs).

ISO and IEC draw attention to the possibility that the implementation of this document may involve the use of (a) patent(s). ISO and IEC take no position concerning the evidence, validity or applicability of any claimed patent rights in respect thereof. As of the date of publication of this document, ISO and IEC had not received notice of (a) patent(s) which may be required to implement this document. However, implementers are cautioned that this may not represent the latest information, which may be obtained from the patent database available at www.iso.org/patents and <https://patents.iec.ch>. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html. In the IEC, see www.iec.ch/understanding-standards.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 42, *Artificial intelligence*.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html and www.iec.ch/national-committees.

Introduction

When AI systems are used to help make decisions that affect people's lives, it is important that people understand how those decisions are made. Achieving useful explanations of the behaviour of AI systems and their components is a complex task. Industry and academia are actively exploring emerging methods for enabling explainability, as well as scenarios and reasons why explainability can be required.

Due to the multitude of stakeholders and communities contributing to this effort, the field is suffering from a certain terminological inconsistency. Most notably, the methods to provide such explanations of the behaviour of an AI system are discussed under the banner of "explainability", "interpretability", (sometimes even other terms like "transparency"), raising the question of how these terms relate to each other. This document aims to provide practical guidance for stakeholders regarding compliance with regulatory requirements labelled one way or another. With this goal in mind, it uses the umbrella term "explainability" and provides a non-exhaustive taxonomy and list of approaches that stakeholders can use to comply with regulatory requirements.

While the overarching goal of explainability is to evaluate the trustworthiness of AI systems, at different stages of the AI system life cycle, diverse stakeholders can have more specific objectives in support of the goal. To illustrate this point, several examples are provided. For developers, the goal can be improving the safety, reliability and robustness of an AI system by making it easier to identify and fix bugs. For users, explainability can help to decide how much to rely on an AI system by uncovering potential sources or existence of unwanted bias or unfairness. For service providers, explainability can be essential for demonstrating compliance with legal requirements. For policy makers, understanding the capabilities and limitations of different explainability methods can help to develop effective policy frameworks that best address societal needs while promoting innovation. Explanations can also help to design interventions to improve business outcomes.

This document describes the applicability and the properties of existing approaches and methods for improving explainability of ML models and AI systems. This document guides stakeholders through the important considerations involved with selection and application of such approaches and methods.

While methods for explainability of ML models can play a central role in achieving the explainability of AI systems, other methods such as data analytics tools and fairness frameworks can contribute to the understanding of AI systems' behaviour and outputs. The description and classification of such complementary methods are out of scope for this document.

Information technology — Artificial intelligence — Objectives and approaches for explainability and interpretability of machine learning (ML) models and artificial intelligence (AI) systems

1 Scope

This document describes approaches and methods that can be used to achieve explainability objectives of stakeholders with regard to machine learning (ML) models and artificial intelligence (AI) systems' behaviours, outputs and results. Stakeholders include but are not limited to, academia, industry, policy makers and end users. It provides guidance concerning the applicability of the described approaches and methods to the identified objectives throughout the AI system's life cycle, as defined in ISO/IEC 22989.

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 22989:2022, *Information technology — Artificial intelligence — Artificial intelligence concepts and terminology*
(<https://standards.iteh.ai>)

3 Terms and definitions

For the purposes of this document, the following terms and definitions apply.

ISO and IEC maintain terminology databases for use in standardization at the following addresses: <https://www.iso.org/obp>

— ISO Online browsing platform: available at <https://www.iso.org/obp>

— IEC Electropedia: available at <https://www.electropedia.org/>

3.1

stakeholder

any individual, group, or organization that can affect, be affected by or perceive itself to be affected by a decision or activity

[SOURCE: ISO/IEC 22989:2022, 3.5.13]

3.2

explainability

property of an *AI system* (3.4) that enables a given human audience to comprehend the reasons for the system's *behaviour* (3.22)

Note 1 to entry: Explainability methods are not limited to the production of explanations, but also include the enabling of interpretations.

3.3

transparency

<system> property of a system that appropriate information about the system is communicated to relevant *stakeholders* (3.1)

Note 1 to entry: Appropriate information for system transparency can include aspects such as features, performance, limitations, components, procedures, measures, design goals, design choices and assumptions, data sources and labelling protocols.

Note 2 to entry: Inappropriate disclosure of some aspects of a system can violate security, privacy, or confidentiality requirements.

[SOURCE: ISO/IEC 22989:2022, 3.5.15]

3.4

artificial intelligence system

AI system

engineered system that generates outputs such as content, forecasts, recommendations or decisions for a given set of human-defined objectives

Note 1 to entry: The engineered system can use various techniques and approaches related to artificial intelligence to develop a model to represent data, knowledge, processes, etc. which can be used to conduct tasks.

[SOURCE: ISO/IEC 22989:2022, 3.1.4]

3.5

machine learning

ML

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

[SOURCE: ISO/IEC 22989:2022, 3.3.5]

3.6

trustworthiness

ability to meet *stakeholder* (3.1) expectations in a verifiable way

Note 1 to entry: Depending on the context or sector and also on the specific product or service, data and technology used, different characteristics apply and need verification to ensure *stakeholders'* (3.1) expectations are met.

Note 2 to entry: Characteristics of trustworthiness include, for instance, reliability, availability, resilience, security, privacy, safety, accountability, transparency, integrity, authenticity, quality and usability.

Note 3 to entry: Trustworthiness is an attribute that can be applied to services, products, technology, data and information as well as, in the context of governance, to organizations.

[SOURCE: ISO/IEC TR 24028:2020, 3.42, modified — Stakeholders' expectations replaced by stakeholder expectations; comma between quality and usability replaced by "and".]

3.7

feature

measurable property of an object or event with respect to a set of characteristics

Note 1 to entry: Features play a role in training and prediction.

Note 2 to entry: Features provide a machine-readable way to describe the relevant objects. As the algorithm will not go back to the objects or events themselves, feature representations are designed to contain all useful information.

[SOURCE: ISO/IEC 23053:2022, 3.3.3]

3.8**global**

property of an *explanation* (3.27) or an *interpretation* (3.28) that describes how model predictions are determined

Note 1 to entry: A global explanation provides an overall understanding of the model's typical operation. For instance, a list of rules or *features* (3.7) that determine the model outputs is an example of global explanation.

3.9**local**

property of an *explanation* (3.27) or an *interpretation* (3.28) that describes how a single model prediction was determined

Note 1 to entry: Compared to a global explanation, a local explanation does not try to explain the whole model.

3.10**post-hoc explanation**

explanation (3.27) built by applying analysis on the model after it has been trained or developed

Note 1 to entry: Post-hoc explanations are often used with *opaque box* (3.16) models, but they are not limited to opaque box models.

3.11**feature-based explanation**

explanation (3.27) of model *behaviour* (3.22) based on input *features* (3.7)

Note 1 to entry: For instance, a measure of how much each input feature contributes to a model's output for a given data point is an example of feature-based explanation.

Note 2 to entry: An input feature for a model does not necessarily correspond to the inputs a user gives as entry as several layers of processing can be applied.

3.12**application programming interface****API**

boundary across which a software application uses facilities of programming languages to invoke software services

[SOURCE: ISO/IEC 13522-6:1998, 3.3]

3.13**backpropagation**

neural network training method that uses the error at the output layer to adjust and optimise the weights for the connections from the successive previous layers

[SOURCE: ISO/IEC 23053:2022, 3.2.1]

3.14**classification model**

<machine learning> machine learning model whose expected output for a given input is one or more classes

[SOURCE: ISO/IEC 23053:2022, 3.1.1]

3.15**closed box****black box**

<access> property of an *AI system* (3.4) or a model within an AI system, whereby only its outputs can be obtained programmatically

3.16**opaque box****black box**

<explainability> property of an *AI system* (3.4) or a model within an AI system, whereby it does not offer *intrinsic interpretability* (3.17)

3.17**intrinsic interpretability****inherent interpretability**

property of an AI model that holds its criteria and *decision process* (3.21) in an intelligible way in its structure or content

Note 1 to entry: Intrinsic interpretability is not limited to access only, but also implies an ability to understand the provided information. For instance, a structure of millions of parameters does not usually constitute an intelligible way of holding it.

Note 2 to entry: Intrinsic interpretability is opposed to *opaque box* (3.16).

3.18**decision**

content or item produced by the *AI system* (3.4) as a fulfilment of its task, based on a given input

Note 1 to entry: The decision can be a class, but also any other form of structured or unstructured data (e.g. a sentence, an image).

3.19**outcome**

one of the various options that the *AI system* (3.4) considers when choosing a given *decision* (3.18)

Note 1 to entry: Outcomes are candidate decisions.

3.20**output**

any data or information returned by the *AI system* (3.4) when processing a given input

Note 1 to entry: Outputs encompass decisions but also any additional data or information that is returned together with a decision, e.g. to contextualize or explain it.

3.21**decision process**

set of steps and criteria used by the *AI system* (3.4) to analyse an input and choose the *decision* (3.18) among the possible *outcomes* (3.20)

Note 1 to entry: Depending on the design of the AI system, that decision process can be embedded in part or in whole, implicitly or explicitly, into the AI system's models.

3.22**behaviour**

<AI system> any observable effect of a given *decision process* (3.21), such as a particular *decision* (3.18), the preferences made among different *outcomes* (3.19), a relationship among multiple decisions or a statistical property of the complete set of decisions made by the *AI system* (3.4) (including future decisions)

Note 1 to entry: Depending on the design of the AI system, the behaviour of the AI system can be attributed to the behaviour of the AI system's models or to their interplay.

3.23**factor**

element, property or other characteristic that is considered during the *decision process* (3.21) and can have an effect on the chosen *decision* (3.18)

3.24**cause**

any type of circumstance that can lead to a given *decision* (3.18), including for instance the presence, absence or value of a *factor* (3.23), but also the analysis made of that factor, its similarity or interaction with other factors, or the presence or absence of a given step or criterion in the *decision process* (3.21)

3.25**rationale**

piece of information or the analysis made of that information, based on which *decisions* (3.18) are made

Note 1 to entry: A rationale provided for a single decision identifies one or more causes as having affected the *decision process* (3.21) of the *AI system* (3.4) when choosing that particular decision. A rationale provided without the context of a specific decision identifies a set of *causes* (3.24) that can affect the *behaviour* (3.22) of the AI system during past or future decisions.

3.26**justification**

piece of information or the analysis made of that information, that is sufficient to choose a given *decision* (3.18) among the possible *outcomes* (3.19)

Note 1 to entry: A justification identifies causes relevant to a given decision, without assumption on the set of causes that have affected the decision process of the *AI system* (3.4).

3.27**explanation**

result of expressing a given *rationale* (3.25) or *justification* (3.26) in a way that humans can understand

Note 1 to entry: Explanations can pertain to a *decision* (3.18) or to an *AI system* (3.4).

3.28**interpretation**

result of understanding (by a human) a given *rationale* (3.25) or *justification* (3.26)

Note 1 to entry: Interpretations can pertain to a *decision* (3.18) or to an *AI system* (3.4).

Note 2 to entry: Interpretations can be produced either based on a received explanation, or directly from observation without an explicit act of expression.

3.29**behavioural accuracy**

adequacy between the *outcomes* (3.19) to which the *explanation* (3.27) leads and the actual *decisions* (3.18) made by the *AI system* (3.4)

3.30**simulatability**

ability of humans to process the provided information and apply the corresponding criteria mentally to obtain the output

4 Symbols and abbreviated terms

CEM	contrastive explanations method
CEM-MAF	contrastive explanations method with monotonic attribute functions
CNN	convolutional neural networks
CV	computer vision
ML	machine learning
XAI	explainable artificial intelligence

5 Overview

Explainability is the property of an AI system that enables a given human audience to comprehend the reasons for the system's behaviour. Reasons are rationales or justifications, as defined in this document with respect to the system's behaviour. The appropriate way of achieving explainability depends on the context and stakeholder characteristics. Stakeholder-appropriate explainability helps to achieve concrete objectives such as:

- identifying the causes of an incorrect decision;
- ensuring that a decision was taken for the right reasons;
- strengthening the confidence in the system.

Users of this document are advised that this concept of explainability (and thus the corresponding set of methods) is more encompassing than some existing uses of the term “explainability”, while more restrictive than others. For instance, it includes some analysis and visualization methods, and not only explanations given by the system itself, but it does not include transparency or AI literacy. Some circles call this concept ‘interpretability’ and use “explainability” in a different way, but this document does not make this kind of distinction. See [Annex A](#) for further explanations on the exact technical scope targeted in this document.

The relevance of specific objectives depends on the stakeholders and what they are trying to achieve. The stakeholders can be interested in achieving one or several of the objectives. Stakeholders' objectives are discussed in more details in [Clause 6](#), subject to the limitations and concerns discussed in [Annex C](#).

Achieving explainability in a system warrants specific methodological considerations throughout the whole life cycle. Guidance on that process is offered in [Clauses 7, 8](#) and [9](#) provide further technical material (taxonomy of properties for needs assessment and corresponding landscape of methods) to support that methodology.

6 Stakeholders' objectives

6.1 General

XAI is a broad field, and stakeholders can have very different reasons to seek explainability. As a result, there is a large variety of corresponding expectations. It is important to consider them and ensure clarity on the stakeholders' objectives, because the utility of an explanation can depend a lot on the stakeholder receiving it, or on the action taken based on it.

Stakeholders can be characterized into various different types as defined in ISO/IEC 22989. This characterization is neither unique, comprehensive, nor non-overlapping, but serves to make a point about different types of explanations. The stakeholder can be someone who participates in developing the AI system (e.g. a developer) or someone who uses the AI system (e.g. an end user applying for a loan).

Each stakeholder can perform different actions. The developer, for instance, is trying to improve the system or decide whether to deploy it and can modify the system based on the explanation. It is therefore more useful if the explanation refers to system characteristics. Such explanations are concerned with the inner workings of a system or are associated with the system's development process. Such explanations can be incomprehensible to end users who are unfamiliar with the system's inner workings.

Alternatively, the system can make decisions that the end user wants to address. Therefore, it is more useful for explanations to emphasize aspects that the end user can comprehend and control. Often there are several influential paths from input to output and the explanation can pick aspects of the input under the end user's control. In the process, the explanation can optimize for comprehension (and control) over faithfulness to the inner-workings of the system.

[Figure 1](#) pictures the various AI stakeholder roles as they are defined in ISO/IEC 22989. This [Clause 6](#) discusses (and illustrates on practical scenarios) the various objectives that stakeholders can have depending on their role. See ISO/IEC 22989 for further information on those stakeholder roles.

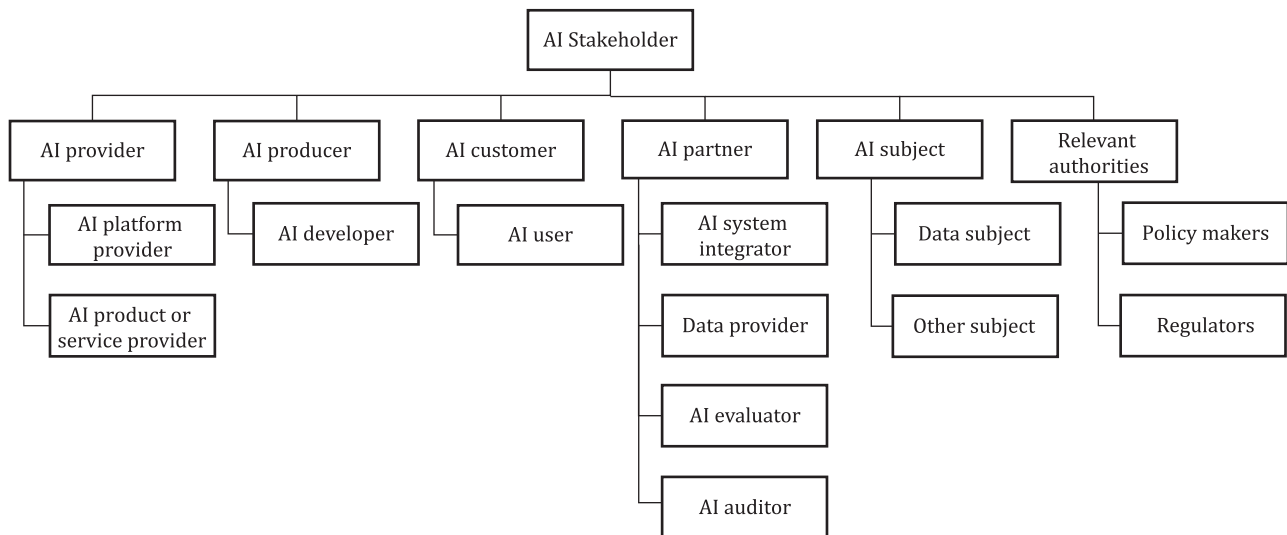


Figure 1 — AI stakeholder roles and their sub-roles

6.2 AI user

A typical objective for AI users seeking explainability is to be reassured on the appropriateness of decisions that the AI system will induce, and that they will not cause harm to others or themselves by using and being influenced by that AI system. This can be linked with a reluctance to interact with automated systems that are perceived as having lower cognitive capabilities than humans, and a desire to maintain human control and enable interventions.

Another common objective is the desire to be able to justify their own actions (driven by the AI system's decisions) when interacting with other stakeholders.

It also happens that AI users seek explainability for their own purposes, in order to gain a better understanding of a given domain, topic or item.

Example scenarios where AI users seek explainability include: a judge that needs to decide whether to use an automated pretrial bail risk assessment system for upcoming cases; or the same judge who for a given case wants to know the extent to which the generated risk score is reliable; a medical doctor considering to overrule the recommendation of an automated diagnosis tool because the proposed diagnostic is atypical; the medical doctor who announces a serious illness to a patient; or a biologist who can discover new causes for known illnesses thanks to “explainable” analysis of health data.

Another scenario is to learn new information (e.g. motivations for customer churn) that can help routing the business logic (e.g. churn mitigation strategy, such as offering discounts to paid content subscription).

6.3 AI developer

AI developers usually seek explainability as a means to acquire a better understanding of the shortcomings of the AI system, in order to improve it. This notably includes detecting biases, spurious correlations, and other failure modes that are unexpected and hard to delineate based on accuracy measures alone. Explainability can also provide them with insights on possible mitigations.

6.4 AI product or service provider

AI product or service providers can be interested in explainability in order to augment the trustworthiness of their AI systems, so that they are more easily adopted by users. It also happens that this desire is driven by the need to comply with an applicable policy, law or regulation.

Another possible objective of AI product or service providers is to gain more information on possible issues caused by their system, so that they can mitigate them, either to augment the system's adoption or to avoid

causing harm. The objectives can be linked with monitoring processes, including drift analysis and churn mitigation.

An example scenario is to ensure that a given product or service does not discriminate AI users or AI subjects based on their gender, race, or any attribute that would be legally punishable and hold the provider accountable.

6.5 AI platform provider

No explainability objectives have been identified.

6.6 AI system integrator

No explainability objectives have been identified.

6.7 Data provider

No explainability objectives have been identified.

6.8 AI evaluator

No explainability objectives have been identified.

6.9 AI auditor

The most typical objective of AI auditors regarding explainability is to gain a better understanding of the limitations and biases of the AI system. They can make judgements on whether it is safe to use that system, by checking that its underlying logic matches domain knowledge, and that it is conceptually sound with principles laid out by the regulations.

An example scenario is ensuring that the AI system is not racist or sexist for instance.

6.10 AI subject

AI subjects' typical interest in explainability is as a means to provide actionable recourse.

An example scenario is for a loan applicant whose application was denied due to automated credit risk assessment, to understand what they should change to get the loan.

It also happens that explainability is desired by AI subjects for purely informative reasons, such as getting solace, or gaining self-awareness. They can want to know whether and how their actions, beliefs and life events have been influenced by AI systems.

6.11 Relevant authorities

6.11.1 Policy makers

No explainability objectives have been identified

6.11.2 Regulators

Regulators are often interested in "explainable" AI systems for purposes such as ensuring that the uses of AI are compatible with preserving human agency, fundamental rights or citizens' rights. This includes a willingness that AI subjects' objectives are met, but also the ability to investigate broader effects on society, such as helping to detect whether an AI system can have a biased behaviour towards certain groups.