



**Publicly  
Available  
Specification**

**ISO/PAS 8800**

**Road vehicles — Safety and artificial  
intelligence**

*Véhicules routiers — Sécurité et intelligence artificielle*

**First edition  
2024-12**

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024



**COPYRIGHT PROTECTED DOCUMENT**

© ISO 2024

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](mailto:copyright@iso.org)  
Website: [www.iso.org](http://www.iso.org)

Published in Switzerland

# Contents

	Page
<b>Foreword</b> .....	<b>vi</b>
<b>Introduction</b> .....	<b>vii</b>
<b>1 Scope</b> .....	<b>1</b>
<b>2 Normative references</b> .....	<b>1</b>
<b>3 Terms and definitions</b> .....	<b>2</b>
3.1 General AI-related definitions.....	2
3.2 Data-related definitions.....	7
3.3 General safety-related definitions.....	9
3.4 Safety: Root cause-, error-and failure-related definitions.....	11
3.5 Miscellaneous definitions.....	12
<b>4 Abbreviated terms</b> .....	<b>14</b>
<b>5 Requirements for conformity</b> .....	<b>15</b>
5.1 Purpose.....	15
5.2 General requirements.....	15
<b>6 AI within the context of road vehicles system safety engineering and basic concepts</b> .....	<b>16</b>
6.1 Application of the ISO 26262 series for the development of AI systems.....	16
6.2 Interactions with encompassing system-level safety activities.....	17
6.3 Mapping of abstraction layers between the ISO 26262 series, ISO/IEC 22989 and this document.....	20
6.4 Example architecture for an AI system.....	22
6.5 Types of AI models.....	23
6.6 AI technologies of a ML model.....	23
6.7 Error concepts, fault models and causal models.....	24
6.7.1 Cause-and-effect chain.....	24
6.7.2 Root cause classes.....	26
6.7.3 Error classification based on the safety impact.....	27
<b>7 AI safety management</b> .....	<b>28</b>
7.1 Objectives.....	28
7.2 Prerequisites and supporting information.....	28
7.3 General requirements.....	28
7.4 Reference AI safety life cycle.....	31
7.5 Iterative development paradigms for AI systems.....	33
7.6 Work products.....	34
<b>8 Assurance arguments for AI systems</b> .....	<b>35</b>
8.1 Objectives.....	35
8.2 Prerequisites and supporting information.....	35
8.3 General requirements.....	36
8.4 AI system-specific considerations in assurance arguments.....	36
8.5 Structuring assurance arguments for AI systems.....	37
8.5.1 Context of the assurance argument.....	37
8.5.2 Categories of evidence.....	38
8.6 The role of quantitative targets and qualitative arguments.....	39
8.7 Evaluation of the assurance argument.....	40
8.8 Work products.....	41
<b>9 Derivation of AI safety requirements</b> .....	<b>41</b>
9.1 Objectives.....	41
9.2 Prerequisites and supporting information.....	42
9.3 General requirements.....	42
9.4 General workflow for deriving safety requirements.....	43
9.5 Deriving AI safety requirements on supervised machine learning.....	46
9.5.1 The need for refined AI safety requirements.....	46

# ISO/PAS 8800:2024(en)

9.5.2	Derivation of refined AI safety requirements to manage uncertainty	47
9.5.3	Refinement of the input space definition for AI safety lifecycle	50
9.5.4	Restricting the occurrence of AI output insufficiencies	50
9.5.5	Metrics, measurements and threshold design	54
9.5.6	Considerations for deriving safety requirements	55
9.6	Work products	56
<b>10</b>	<b>Selection of AI technologies, architectural and development measures</b>	<b>56</b>
10.1	Objectives	56
10.2	Prerequisites	56
10.3	General requirements	56
10.4	Architecture and development process design or refinement	57
10.5	Examples of architectural and development measures for AI systems	58
10.6	Work products	62
<b>11</b>	<b>Data-related considerations</b>	<b>62</b>
11.1	Objectives	62
11.2	Prerequisites and supporting information	62
11.3	General requirements	62
11.4	Dataset life cycle	63
11.4.1	Datasets and the AI safety lifecycle	63
11.4.2	Reference dataset lifecycle	64
11.4.3	Dataset safety analysis	65
11.4.4	Dataset requirements development	71
11.4.5	Dataset design	74
11.4.6	Dataset implementation	75
11.4.7	Dataset verification	75
11.4.8	Dataset validation	76
11.4.9	Dataset maintenance	77
11.5	Work products	77
<b>12</b>	<b>Verification and validation of the AI system</b>	<b>78</b>
12.1	Objectives	78
12.2	Prerequisites and supporting information	78
12.3	General requirements	78
12.4	AI/ML specific challenges to verification and validation	80
12.5	Verification and validation of the AI system	81
12.5.1	Scope of verification and validation of the AI system	81
12.5.2	AI component testing	84
12.5.3	Methods for testing the AI component	86
12.5.4	AI system integration and verification	88
12.5.5	Virtual testing vs physical testing	88
12.5.6	Evaluation of the safety-related performance of the AI system	89
12.5.7	AI system safety validation	90
12.6	Work products	91
<b>13</b>	<b>Safety analysis of AI systems</b>	<b>91</b>
13.1	Objectives	91
13.2	Prerequisites and supporting information	92
13.3	General requirements	92
13.4	Safety analysis of the AI system	93
13.4.1	Scope of the AI safety analysis	93
13.4.2	Safety analysis based on the results of testing	95
13.4.3	Safety analysis techniques	95
13.5	Work products	97
<b>14</b>	<b>Measures during operation</b>	<b>97</b>
14.1	Objectives	97
14.2	Prerequisites and supporting information	98
14.3	General requirements	98
14.4	Planning for operation and continuous assurance	99

# ISO/PAS 8800:2024(en)

14.4.1	Safety risk of the AI system during operation phase.....	99
14.4.2	Safety activities during the operation phase.....	99
14.5	Continual, periodic re-evaluation of the assurance argument.....	100
14.6	Measures to assure safety of the AI system during operation.....	101
14.6.1	General.....	101
14.6.2	Technical safety measures.....	101
14.6.3	Safe operation guidance and misuse prevention in the field.....	102
14.7	Field data collection.....	103
14.8	Evaluation and continuous development.....	104
14.8.1	Field risk evaluation.....	104
14.8.2	Countermeasures addressing field risk.....	105
14.8.3	AI re-training, re-validation, re-approval and re-deployment.....	105
14.9	Work products.....	106
<b>15</b>	<b>Confidence in use of AI development frameworks and software tools used for AI model development.....</b>	<b>106</b>
15.1	Objectives.....	106
15.2	Prerequisites and supporting information.....	107
15.3	General requirements.....	107
15.4	Confidence in the use of AI development frameworks.....	107
15.5	Confidence in the use of tools used to support the AI-safety lifecycle.....	109
15.6	Principles for data-driven AI model training and evaluation.....	110
15.7	Work products.....	110
<b>Annex A</b>	<b>(informative) Overview and workflow of this document.....</b>	<b>111</b>
<b>Annex B</b>	<b>(informative) Example assurance argument structure for an AI system.....</b>	<b>116</b>
<b>Annex C</b>	<b>(informative) ISO 26262 gap analysis for ML.....</b>	<b>130</b>
<b>Annex D</b>	<b>(informative) Detailed considerations on safety-related properties of AI systems.....</b>	<b>137</b>
<b>Annex E</b>	<b>(informative) STAMP/STPA example.....</b>	<b>139</b>
<b>Annex F</b>	<b>(informative) Identification of software units within NN-based systems.....</b>	<b>144</b>
<b>Annex G</b>	<b>(informative) Architectural and development measures for AI systems.....</b>	<b>147</b>
<b>Annex H</b>	<b>(informative) Typical performance metrics for machine learning.....</b>	<b>162</b>
<b>Bibliography</b>	.....	<b>167</b>

## Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

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 ISO 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](http://www.iso.org/directives)).

ISO draws attention to the possibility that the implementation of this document may involve the use of (a) patent(s). ISO takes 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 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](http://www.iso.org/patents). ISO 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](http://www.iso.org/iso/foreword.html).

This document was prepared by Technical Committee ISO/TC 22, *Road vehicles*, Subcommittee SC 32, *Electrical and electronic components and general system aspects*.

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](http://www.iso.org/members.html).

## Introduction

The purpose of this document is to provide industry-specific guidance on the use of AI systems in safety-related functions. It is not restricted to specific AI methods or specific vehicle functions.

This document defines a framework for managing AI safety that tailors or extends existing approaches currently defined in the ISO 26262 series and in ISO 21448.

Functional safety-related risks associated with malfunctioning behaviour of an AI system are addressed by tailoring or extending relevant clauses from ISO 26262-series.

The risks related to functional insufficiencies in the AI system are addressed by extending the concepts and guidance provided by ISO 21448. A causal model for understanding the sources of functional insufficiencies in the AI system is proposed. The model is used to derive a set of safety requirements on the AI system as well as a set of risk reduction measures.

**NOTE 1** ISO 21448 is applicable to intended functionalities where proper situational awareness is essential to safety and where such situational awareness is derived from sensors and processing algorithms, especially functionalities of emergency intervention systems and systems with ISO/SAE PAS 22736 levels 1 to 5 for driving automation. It is therefore possible that systems utilize AI technologies that do not fall within the scope of ISO 21448.

**EXAMPLE 1** ISO 21448 does not apply to the development of an engine control unit that uses AI to optimize its performance whereas this document does.

This document recognizes that due to the wide range of applications of AI and associated safety requirements, as well as the rapidly evolving state-of-the-art, it is not possible to provide detailed requirements on the process or product characteristics required to achieve an acceptably low level of residual risk associated with the use of AI systems. Therefore, in addition to providing guidance for tailoring or extending the ISO 26262 series and ISO 21448, this document focuses on the principles that support the creation of a project-specific assurance argument for the safety of the AI elements within on-board vehicle systems. This includes proposing risk reduction measures during the design and operation phases using an iterative approach to reducing risk as outlined in ISO/IEC Guide 51.

Hazard analysis and risk analysis are beyond the scope of this document. These are considered a part of the vehicle level systems safety engineering activities described in the ISO 26262 series and ISO 21448, or in application of specific standards such as ISO TS 5083.

ISO/IEC TR 5469 provides generic guidance for the application of AI technologies as part of safety functions, independent of specific industry sectors. Many of the concepts outlined in ISO/IEC TR 5469 can be applied in the context of road vehicles. There is therefore a close relationship to concepts described within this document and ISO/IEC TR 5469.

ISO/IEC TR 5469 provides classification schemes to determine the safety requirements on the AI/ML function. These include the usage level and AI technology class.

The usage level is related to the nature of the task being performed by the engineered AI system.

**NOTE 2** The usage levels are described in ISO/IEC TR 5469:2024, 6.2.

The technology class is related to the problem complexity and the transferability of existing standards to demonstrating an adequate level of safety based on properties of the target function and the AI technology used.

**NOTE 3** For the technology classes, see ISO/IEC TR 5469:2024, 6.2.

This document does not explicitly call out the classes and usage levels of ISO/IEC TR 5469.

**EXAMPLE 2** For some AI technology, the application of ISO 26262 is deemed to be sufficient. This corresponds to Class I of ISO/IEC TR 5469.

The guidance outlined within this document is relevant for all usage of AI for which safety requirements can foreseeably be allocated either through:

- a) the use of AI for the functionality itself;

b) the use of AI as a safety mechanism.

NOTE 4 These usages correspond to the usage levels A1, A2, C of ISO/IEC TR 5469. In all cases, the applicability of the guidance provided within this document can be determined by the allocation of safety requirements to the AI technology, whereas the usage levels of ISO/IEC TR 5469 can be used to support the requirements elicitation process.

This document is aligned with standards and documents developed by ISO/IEC JTC1/SC42. AI-specific definitions are used from ISO/IEC 22989, unless in conflict with safety-specific definitions.

Other documents developed within ISO/IEC JTC1/SC42 can be used to provide additional guidance on specific aspects of AI that are relevant to safety-related properties. Examples of such documents include ISO/IEC TR 24027 and ISO/IEC TR 24029-1.

This document harmonizes the concepts already described in ISO 21448:2022, Annex D.2 and ISO/TS 5083:20—<sup>1)</sup>, Annex B whilst extending these with specific guidance regarding the definition of safety requirements of machine learning (ML), ML safety analyses and the creation of associated safety evidence during the development and deployment lifecycle.

ISO/TS 5083:20—, Annex B is an application of this document to automated driving systems (ADS).

The relationship with the above-mentioned documents is summarized in [Table 1-1](#).

**Table 1-1 — How this document relates to other publications on AI safety**

Publication	Relationship with this document
ISO/IEC 22989	AI-specific definitions are used from ISO/IEC 22989, unless in conflict with safety-specific definitions. Safety-related properties are a subset of generic AI properties described in ISO/IEC 22989.
ISO/IEC TR 5469	This document does not explicitly call out the classes and usage levels of ISO/IEC TR 5469. This document considers and adapts to road vehicles the general framework described in ISO/IEC TR 5469 on safety properties, virtual testing and physical testing, confidence in use of AI development frameworks and architectural redundancy patterns.
ISO 26262	This document is a tailoring or extension of ISO 26262 for AI elements of the system. See <a href="#">Clause 5</a> for details.
ISO 21448	This document is a tailoring or extension of ISO 21448 for AI elements of the system. See <a href="#">Clause 5</a> for details.
ISO TS 5083:20—	ISO TS 5083:20—, Annex B is an application of this document to automated driving systems (ADS).

This document adds the following contents with respect to the documents listed in [Table 1-1](#):

- tailoring or extensions of ISO 26262 and ISO 21448 required specifically for AI elements of the system (referred to as AI systems);
- a conceptual model for reasoning about errors and their causes specific to AI systems;
- a reference AI safety lifecycle;
- the safety assurance argument for AI systems;
- a method for deriving AI safety requirements for AI systems;
- considerations for the design of safe AI systems;
- considerations on data management for the AI systems;

1) Under preparation. Stage at the time of publication: ISO/DTS 5083.

## ISO/PAS 8800:2024(en)

- a verification and validation strategy for AI systems;
- a safety analysis approach for AI systems (focused on insufficiencies);
- activities during operation required to ensure the continuous AI safety.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

[STANDARDSISO.COM](https://standardsiso.com) : Click to view the full PDF of ISO/PAS 8800:2024

# Road vehicles — Safety and artificial intelligence

## 1 Scope

This document applies to safety-related systems that include one or more electrical and/or electronic (E/E) systems that use AI technology and that is installed in series production road vehicles, excluding mopeds. It does not address unique E/E systems in special vehicles, such as E/E systems designed for drivers with disabilities.

This document addresses the risk of undesired safety-related behaviour at the vehicle level due to output insufficiencies, systematic errors and random hardware errors of AI elements within the vehicle. This includes interactions with AI elements that are not part of the vehicle itself but that can have a direct or indirect impact on vehicle safety.

EXAMPLE 1 Examples of AI elements within the vehicle include the trained AI model and AI system.

EXAMPLE 2 Direct impact on safety can be due to object detection by elements external to the vehicle.

EXAMPLE 3 Indirect impact on safety can be due to field monitoring by elements external to the vehicle.

The development of AI elements that are not part of the vehicle is not within the scope of this document. These elements can conform to domain-specific safety guidance. This document can be used as a reference where such domain-specific guidance does not exist.

This document describes safety-related properties of AI systems that can be used to construct a convincing safety assurance claim for the absence of unreasonable risk.

This document does not provide specific guidelines for software tools that use AI methods.

This document focuses primarily on a subclass of AI methods defined as machine learning (ML). Although it covers the principles of established and well-understood classes of ML, it does not focus on the details of any specific AI methods e.g. deep neural networks.

## 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 21448:2022, *Road vehicles — Safety of the intended functionality*

ISO 26262-1:2018, *Road vehicles — Functional safety — Part 1: Vocabulary*

ISO 26262-2:2018, *Road vehicles — Functional safety — Part 2: Management of functional safety*

ISO 26262-6:2018, *Road vehicles — Functional safety — Part 6: Product development at the software level*

ISO 26262-8:2018, *Road vehicles — Functional safety — Part 8: Supporting processes*

ISO/IEC 22989:2022, *Information technology — Artificial intelligence — Artificial intelligence concepts and terminology*

### 3 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO 26262-1, ISO 21448, ISO/IEC 22989 and the following apply.

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

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <https://www.electropedia.org/>

#### 3.1 General AI-related definitions

##### 3.1.1

##### **AI component**

element of an *AI system* (3.1.17)

EXAMPLE 1 An *AI pre-processing* (3.1.11) component.

EXAMPLE 2 An *AI post-processing* (3.1.9) component.

EXAMPLE 3 An *AI model* (3.1.7).

EXAMPLE 4 A conventional software component inside an AI system.

Note 1 to entry: AI components that are not AI models or that do not contain AI models are not developed according to this document. The integration of these components with AI components that are AI models or that contain AI models is performed according to this document.

Note 2 to entry: See 6.3 for an elaboration of the relationship of the different abstraction layers of the ISO 26262 series, ISO/IEC 22989 and this document with each other.

[SOURCE: ISO/IEC 22989:2022, 3.1.2, modified to be consistent with ISO 26262-1 definitions — "Functional element" was replaced with "element", reworded to not use "construct", examples and Notes to entry were added.]

##### 3.1.2

##### **AI controllability**

ability of an external agent to control the *AI element* (3.1.3), its output or the behaviour of the item influenced by the AI output in order to prevent harm

EXAMPLE Before setting a pulse-width modulation (PWM) signal of an actor determined by an *AI model* (3.1.7), the PWM output is limited by a simple threshold or the consumer substitutes the PWM signal with an approximate physical model.

Note 1 to entry: An external agent is a person or an element not belonging to the *AI system* (3.1.17).

##### 3.1.3

##### **AI element**

*AI component* (3.1.1) or *AI system* (3.1.17)

Note 1 to entry: An AI element can refer to a subset of *components* (3.5.2) within an AI system that provide related functionality.

Note 2 to entry: See 6.3 for an elaboration of the relationship of the different abstraction layers of the ISO 26262 series ISO/IEC 22989 and this document with each other.

##### 3.1.4

##### **AI explainability**

property of an *AI system* (3.1.17) to express important factors influencing the AI system's outputs in a way that humans can understand

EXAMPLE The AI system can be explainable by natural language or by visualizing feature attribution methods like gradient-based heat/saliency maps.

### 3.1.5

#### AI generalization

ability of an *AI model* (3.1.7) to adapt and perform well on previously unseen data during inference

### 3.1.6

#### AI method

type of *AI model* (3.1.7)

EXAMPLE 1 Deep neural network.

EXAMPLE 2 K-nearest neighbour.

EXAMPLE 3 Support vector machine.

### 3.1.7

#### AI model

construct containing logical operations, arithmetical operations or a combination of both to generate an inference or prediction based on input data or information without being completely defined by human knowledge

Note 1 to entry: Inference is using a model to understand the relation between predictors and a target. Prediction is using a model to generate a prediction (values close to the real seen or unseen targets) based on the inputs.

### 3.1.8

#### AI model validation

evaluation of the performance of different *AI model* (3.1.7) candidates through testing

Note 1 to entry: There are three terms, "AI model validation", "validation" and "safety validation", that are distinguished in this document. AI model validation originates from the validation data used by the AI community, validation originates from classic system development and safety validation originates from the ISO 26262 series.

Note 2 to entry: The AI model validation is executed using the AI validation dataset.

### 3.1.9

#### AI post-processing

any processing that is applied to the output of an *AI model* (3.1.7) for the purpose of mapping the raw output/s to a more contextually relevant and consumable format

EXAMPLE 1 A non-maximum suppression and thresholding for a bounding-box generation that serves to remove bounding boxes of low relevance and duplicates.

EXAMPLE 2 The outputs of a mixture density network are combined with a physical model (a hybrid model).

Note 1 to entry: AI post-processing also includes any data conversion that is used to bring the output into a common format for better comparability.

Note 2 to entry: AI post-processing can have a positive or a negative impact on the safety-related properties of the output of the *AI system* (3.1.17).

### 3.1.10

#### AI predictability

ability of the *AI system* (3.1.17) to produce trusted predictions

Note 1 to entry: Trusted predictions means that the predications are accurate and that this claim is supported by statistical evidence.

### 3.1.11

#### AI pre-processing

any processing that is applied to the input of an *AI model* (3.1.7)

### 3.1.12

#### AI reliability

ability of the *AI element* (3.1.3) to perform the *AI task* (3.1.18) without *AI error* (3.4.1) under stated conditions and for a specified period of time

### 3.1.13

#### AI resilience

ability of the *AI element* (3.1.3) to recover and continue performing the *AI task* (3.1.18) after the occurrence of an *AI error* (3.4.1).

### 3.1.14

#### AI robustness

ability to maintain an acceptable level of performance under the presence of semantically insignificant but reasonably expected changes to the input

EXAMPLE In image data these insignificant input changes can stem from naturally-induced image corruptions or sensor noise.

### 3.1.15

#### AI safety

absence of unreasonable *risk* (3.3.10) due to *AI errors* (3.4.1) caused by faults and functional insufficiencies

Note 1 to entry: This definition only applies in the context of this document. The term "AI safety" is commonly understood to have a broader meaning which includes ethics, value alignment, long-term considerations, etc.

### 3.1.16

#### AI safety requirement

*safety requirement* (3.3.14) of an *AI element* (3.1.3)

### 3.1.17

#### AI system

item or element that utilises one or more *AI models* (3.1.7)

EXAMPLE An AI system consisting of the *AI component* (3.1.1) "deep neural network for bounding box generation (AI model)" and of the AI component "non-maximum suppression algorithm (*AI post-processing* (3.1.9) AI component)".

Note 1 to entry: The AI system can use various *AI methods* (3.1.6) and can utilize different *AI technologies* (3.1.19).

Note 2 to entry: The boundaries of the AI system are determined during the definition of AI system architecture.

Note 3 to entry: The AI system can contain one or more AI components.

Note 4 to entry: The term "AI system" serves in this document as the top level of abstraction of the content to be developed in conformity to the corresponding standard. As such it is possible in a distributed development that what one party considers to be an AI component, the other party considers to be an AI system, as for the latter it represents the top level of the content they develop.

Note 5 to entry: See 6.3 for an elaboration of the relationship of the different abstraction layers of the ISO 26262 series, ISO/IEC 22989 and this document with each other.

### 3.1.18

#### AI task

action required by the *AI element* (3.1.3) to achieve a specific goal

Note 1 to entry: Examples of AI tasks include classification, regression, ranking, clustering and dimensionality reduction.

Note 2 to entry: The AI task can be seen as a semantic description of the *AI model* (3.1.7).

[SOURCE: ISO/IEC 22989:2022, 3.1.35, modified — "task" has been replaced with "AI task", "by the AI element" has been added, "<artificial intelligence>" has been removed; and the Notes to entries have been modified.]

### 3.1.19

#### AI technology

any technology used within the lifecycle of an *AI system* (3.1.17) to design, develop, train, test, validate and implement the *AI model* (3.1.7)

EXAMPLE Examples of AI technologies are provided in 6.6

### 3.1.20

#### AI testing

testing the *AI system* (3.1.17) or *AI model* (3.1.7) to estimate the expected performance and generalization capability in the field

Note 1 to entry: The AI testing is executed using an AI test dataset.

Note 2 to entry: See also ISO 26262-1:2018, 3.169.

### 3.1.21

#### AI system safety validation

confirmation that an *AI safety requirement* (3.1.16) allocated to the *AI system* (3.1.17) is fulfilled

Note 1 to entry: In other standards, validation indicates that requirements are suitable for the intended use. In this document, the term is intentionally used in a different way that is common in the ML community, i.e. to verify the implementation of the requirement.

### 3.1.22

#### bias

undesired, systematic difference in the *AI systems* (3.1.17) predictions with respect to particular classes of inputs in comparison with others due to potential incorrect learning process

EXAMPLE The classes of inputs can refer to images of objects and people in the context of computer vision.

Note 1 to entry: Bias can arise from an undesired systematic difference within the dataset, from limitations within the training process, or from limitations within the *AI model* (3.1.7) capability itself to accurately reflect the dataset.

[SOURCE: ISO/IEC 22989:2022, 3.5.4, modified —the definition was adapted to the AI context, Note 1 to entry was replaced, the EXAMPLE was added.]

### 3.1.23

#### control element

*element* (3.5.4) controlling the execution of the *AI task* (3.1.18) by the *AI element* (3.1.3) and other AI element-related operations like updates

Note 1 to entry: The control element can control non-AI elements as well.

### 3.1.24

#### data pre-processing

part of the AI workflow that transforms raw data so they are usable as the input to create the *AI model* (3.1.7)

Note 1 to entry: Pre-processing can include reformatting, removal of outliers and duplicates, and ensuring the completeness of the dataset.

### 3.1.25

#### encompassing system

item which contains the *AI system* (3.1.17)

### 3.1.26

#### ground truth

set of dataset annotations that are taken to be correct

Note 1 to entry: Individual annotations are derived from information external to the dataset.

Note 2 to entry: Individual annotations may be refined as new information becomes available.

[SOURCE: ISO/IEC 2382-37:2022, 37.09.34]

### 3.1.27

#### hyperparameter

parameters of the used *AI technologies* (3.1.19) that affect both the performance of the *AI model* (3.1.7) and its learning process

Note 1 to entry: Hyperparameters are selected prior to training and can be used to help estimate *model parameters* (3.1.35).

Note 2 to entry: Examples of hyperparameters include the number of network layers, the width of each layer, the type of activation function, the optimization method, the learning rate for neural networks, the choice of kernel function in a support vector machine, the number of leaves or the depth of a tree, the number of clusters in K-means clustering, the maximum number of iterations of the expectation maximization algorithm and the number of Gaussians in a Gaussian mixture.

[SOURCE: ISO/IEC 22989:2022, 3.3.4, modified — The term has been redefined to be applicable to all kinds of AI methods, not only machine learning.]

### 3.1.28

#### input space

set of possible input values

Note 1 to entry: See *semantic input space* (3.1.34) and *syntactic input space* (3.1.36) for ways an input space can be specified.

### 3.1.29

#### machine learning

##### ML

process of optimizing model parameters through computational techniques, such that the model's behaviour aligns with data or experience and enables prediction beyond the training set

EXAMPLE Learning from experience can mean trying to represent non-static data like simulation, reinforcement learning environment, etc.

[SOURCE: ISO/IEC 22989:2022, 3.3.5, modified — "Reflects the data or experience" was replaced with "aligns with data or experience and enables prediction beyond the training set", the EXAMPLE was added.]

### 3.1.30

#### ML algorithm

algorithm to optimize parameters of a *ML model* (3.1.31) from data according to given criteria

EXAMPLE Consider solving a univariate linear function  $y = \theta_0 + \theta_1 x$  where  $y$  is an output or result,  $x$  is an input,  $\theta_0$  is an intercept (the value of  $y$  where  $x=0$ ) and  $\theta_1$  is a weight. In ML, the process of determining the intercept and weights for a linear function is known as linear regression.

[SOURCE: ISO/IEC 22989:2022, 3.3.6, modified — "determine" was replaced with "optimize".]

### 3.1.31

#### ML model

mathematical construct that generates an inference or prediction based on input data or information and comprises a functionality that is created by *machine learning* (3.1.29)

EXAMPLE If a univariate linear function ( $y = \theta_0 + \theta_1 x$ ) has been trained using linear regression, the resulting model can be  $y = 3 + 7x$ .

Note 1 to entry: A ML model results from training based on a *ML algorithm* (3.1.30).

[SOURCE: ISO/IEC 22989:2022, 3.3.7, modified — "And comprises a functionality that is created by *machine learning*" was added to distinguish from other mathematical constructs.]

### 3.1.32

#### **ML model training**

iterative process to optimize a *ML model's* (3.1.31) input and output behaviour on a given training dataset with the intention to improve its quality (e.g. AI accuracy, *AI robustness* (3.1.14), generalization capability, run time), based on a *ML algorithm* (3.1.30) that can adapt ML model parameters, *hyperparameters* (3.1.27), cost function or the model structures itself

[SOURCE: ISO/IEC 22989:2022, 3.3.15, modified to elaborate the procedure and intention.]

### 3.1.33

#### **safety-related AI element**

*AI element* (3.1.3) that contributes to the achievement of an *AI safety requirement* (3.1.16) allocated to the AI system or can contribute to the violation of an *AI safety requirement* allocated to the AI system or can contribute to both

### 3.1.34

#### **semantic input space**

set of possible input values on a semantic level

EXAMPLE The semantic input space acquired by a camera sensor can be described as consisting of street images containing lane markers of different colours, orientation and degradations that appear in different lighting and weather conditions.

Note 1 to entry: The semantic values correspond and conform to abstract semantic concepts expected within the input space.

### 3.1.35

#### **semantic output space**

set of possible output values on a semantic level

### 3.1.36

#### **syntactic input space**

set of possible input values on a syntactic level

EXAMPLE The syntactic input space acquired by a camera sensor can be described as an RGB image array of integers.

Note 1 to entry: The syntactic values can correspond and conform to the output values from a low-level sensor.

### 3.1.37

#### **syntactic output space**

set of possible output values on a syntactic level

### 3.1.38

#### **trained ML model**

*ML model* (3.1.31) with a set of model parameters as result of *ML model training* (3.1.32)

[SOURCE: ISO/IEC 22989:2022, 3.3.14, modified — "ML model with a set of model parameters as" added as part of the term.]

## 3.2 Data-related definitions

### 3.2.1

#### **AI test dataset**

dataset used to estimate the performance and generalization capability of an AI model or an AI system

Note 1 to entry: See [Clause 11](#) for more details.

### 3.2.2

#### **AI validation dataset**

dataset used to compare the performance of different candidate *AI models* (3.1.7)

### 3.2.3

#### **dataset insufficiency**

insufficiency of the dataset regarding data-related safety properties under consideration

Note 1 to entry: Dataset insufficiency includes data integrity errors and data distribution errors.

### 3.2.4

#### **field monitoring dataset**

dataset collected after the release of the *AI system* (3.1.17) while the product is in operation and used specifically for field monitoring of the performance of the AI system

### 3.2.5

#### **hybrid dataset**

dataset comprising data elements that are both real-world data elements and synthetic data elements

### 3.2.6

#### **in distribution data**

data whose features relevant to the *AI task* (3.1.18) are present and sufficiently well represented in the training dataset

Note 1 to entry: In distribution, input does not guarantee correctness of AI model output.

### 3.2.7

#### **metadata**

data that provides additional information about the data element or dataset but is usually not directly involved in the training process

Note 1 to entry: Some metadata (e.g. ground truth) is also used for training.

### 3.2.8

#### **out of distribution data**

data containing features relevant for the *AI task* (3.1.18), either absent or not sufficiently well represented in the *training dataset* (3.2.12), that can result in an *AI error* (3.4.1)

Note 1 to entry: Out of distribution (OOD) refers to data or inputs that fall outside the scope of what an AI or ML model was trained on or is designed to handle. When an AI system encounters OOD data, it can struggle to make accurate predictions or decisions because it lacks the necessary knowledge and experience to handle such inputs effectively. OOD data can lead to unexpected or unreliable model behaviour.

### 3.2.9

#### **real-world dataset**

dataset comprising data elements that have been created by real world acquisitions

### 3.2.10

#### **safety-related KPI**

key performance indicator relevant for the achievement of *AI safety* (3.1.15)

### 3.2.11

#### **synthetic dataset**

dataset comprising data elements that have been created artificially

Note 1 to entry: "Created artificially" implies that the data was not directly collected from something that happened in the real world. Additionally, the data does not necessarily represent something that already happened in the real world.

### 3.2.12

#### **training dataset**

dataset used to train an *ML model* (3.1.31)

[SOURCE: ISO/IEC 22989:2022, 3.3.16, modified — definition reworked to contain the term "dataset".]

### 3.3 General safety-related definitions

#### 3.3.1

##### **assurance**

grounds for justified confidence that a claim has been or will be achieved

[SOURCE: ISO/IEC/IEEE 15026-1:2019]

#### 3.3.2

##### **assurance argument**

reasoned, auditable artefact created supporting the contention that its top-level claim (or set of claims) is satisfied, including systematic arguments, its underlying evidence and explicit assumptions that support the claim(s)

Note 1 to entry: An assurance argument contains the following and their relationships:

- one or more claims about properties;
- arguments that logically link the evidence and any assumptions to the claim(s);
- a body of evidence and possible assumptions supporting these arguments for the claim(s);
- justification of the choice of top-level claim and the method of reasoning.

[SOURCE: ISO/IEC/IEEE 15026-1:2019, modified — "argumentation" replaced with "arguments".]

#### 3.3.3

##### **claim**

true-false statement about the limitations on the values of an unambiguously defined property — called the claim's property — and limitations on the uncertainty of the property's values falling within these limitations during the claim's duration of applicability under stated conditions

Note 1 to entry: Uncertainties may also be associated with the duration of applicability and the stated conditions.

Note 2 to entry: A claim can contain the following:

- property of the system-of-interest;
- limitations on the value of the property associated with the claim (e.g. on its range);
- limitations on the uncertainty of the property value meeting its limitations;
- limitations on the duration of the claim's applicability;
- duration-related uncertainty;
- limitations on conditions associated with the claim;
- condition-related uncertainty.

Note 3 to entry: The term "limitations" is used to fit the many situations that can exist. Values can be a single value or multiple single values, a range of values or multiple ranges of values, and can be multi-dimensional. The boundaries of these limitations are sometimes not sharp, e.g. they can involve probability distributions and can be incremental.

[SOURCE: ISO/IEC/IEEE 15026-1:2019]

#### 3.3.4

##### **undesired safety-related behaviour at the vehicle level**

hazardous behaviour, *RFIM prevention issue* (3.3.9) or malfunctioning behaviour that can cause a hazard

### 3.3.5

#### **hazard**

potential source of harm

[SOURCE: ISO 26262-1:2018, 3.75, modified — "caused by malfunctioning behaviour of the item" and Note 1 to entry deleted.]

### 3.3.6

#### **influencing factor**

factor contributing to the achievement or the absence of a *safety-related property* (3.3.13)

### 3.3.7

#### **misuse**

usage in a way not intended by the manufacturer or the service provider

[SOURCE: ISO 21448:2022, 3.17 modified — The Note to entry and EXAMPLEs were deleted.]

### 3.3.8

#### **reasonably foreseeable indirect misuse**

RFIM

reasonably foreseeable misuse which leads to a reduced controllability of the hazardous behaviour, to a potentially increased severity of an occurring accident or a combination of both

[SOURCE: ISO 21448:2022, 3.17 modified — The term was taken from Note 5 to entry of 3.17 and is now explicitly defined.]

### 3.3.9

#### **RFIM prevention issue**

inability to prevent or detect and mitigate a *reasonably foreseeable indirect misuse* (3.3.8)

### 3.3.10

#### **risk**

combination of the probability of occurrence of harm and the severity of that harm

Note 1 to entry: Other forms of risk definitions exist, e.g. risk for other topics like the risk of a project to fail, etc. This document focuses on the risk regarding safety. Hence this definition was chosen.

Note 2 to entry: The resulting risk evaluation of an error of an AI component is typically equivalent to the evaluation of the potential to lead to a violation of a safety requirement allocated to the AI system. The evaluation can be quantitative as well as qualitative, depending on the safety requirement.

[SOURCE: ISO 26262-1:2018, 3.128, modified — Notes 1 and 2 to entry were added.]

### 3.3.11

#### **AI safety measure**

activity or technical solution to avoid, detect or control *AI errors* (3.4.1), to mitigate their harmful effects or a combination thereof

EXAMPLE AI safety analysis.

Note 1 to entry: Safety measures include architectural measures.

Note 2 to entry: The AI safety measures include safety measures of *AI elements* (3.1.3) as defined in the ISO 26262 series as well as measures to address functional insufficiencies in conformity to ISO 21448 (e.g. functional modifications addressing SOTIF-related risks).

### 3.3.12

#### **safety validation**

assurance, based on examination and tests, that the safety goals are adequate and have been achieved with a sufficient level of integrity

Note 1 to entry: There are three terms, "AI model validation", "validation" and "safety validation", that are distinguished in this document. AI model validation originates from the validation data used by the AI community, validation originates from classic system development and safety validation originates from the ISO 26262 series. The three validation meanings are not the same.

[SOURCE: ISO 26262-1:2018, modified — Note 1 to entry was replaced.]

### 3.3.13

#### **safety-related property**

property impacting safety

### 3.3.14

#### **safety requirement**

requirement related to safety

EXAMPLE 1 SOTIF requirement.

EXAMPLE 2 Functional safety requirement.

EXAMPLE 3 Technical safety requirement.

Note 1 to entry: This includes, but is not limited to, safety requirements motivated by functional safety as well as SOTIF.

### 3.3.15

#### **work product**

work product of the safety lifecycle that can be used as evidence within a safety assurance argument

## 3.4 Safety: Root cause-, error-and failure-related definitions

### 3.4.1

#### **AI error**

one or more discrepancies between computed, observed or measured values or conditions of the *AI element* (3.1.3) and the true, specified or theoretically correct values or conditions of the AI element

Note 1 to entry: An AI error can be a single discrepancy or a sequence of discrepancies.

Note 2 to entry: An AI error can be an error caused by a fault. Faults are typically addressed by the ISO 26262 series.

Note 3 to entry: An AI error can be an output insufficiency caused by a functional insufficiency.

### 3.4.2

#### **AI triggering condition**

specific conditions of a scenario that serve as an initiator for a subsequent *AI error* (3.4.1)

Note 1 to entry: Functional insufficiencies or faults are themselves not AI triggering conditions but are potentially activated by them thus leading to the occurrence of an AI error.

### 3.4.3

#### **contributing AI error**

*AI error* (3.4.1) which can lead to a violation of an *AI safety requirement* (3.1.16) allocated to the *AI system* (3.1.17), either by itself or in combination with one or more other AI errors

### 3.4.4

#### **AI error rate**

probability density of *AI error* (3.4.1) occurrence divided by probability of no AI error occurring until the measuring point

Note 1 to entry: Measurement units can include errors per h, errors per km, etc.

Note 2 to entry: This is an analogue definition to the failure rate.

### 3.4.5

#### **functional insufficiency**

*insufficiency of specification* (3.4.6) or performance insufficiency

Note 1 to entry: A functional insufficiency activated by a triggering condition leads, per definition, to either an output insufficiency, a hazardous behaviour, a RFIM prevention issue or a combination of these.

[SOURCE: ISO 21448:2022, 3.8, modified – The EXAMPLEs, Figures and Notes to entry have been removed. Note to entry 1 has been added.]

### 3.4.6

#### **insufficiency of specification**

specification, possibly incomplete, contributing to either a hazardous behaviour or an *RFIM prevention issue* (3.3.9) when activated by one or more triggering conditions

Note 1 to entry: An insufficiency of specification activated by a triggering condition leads per definition to either an output insufficiency, a hazardous behaviour, an RFIM prevention issue or a combination of these.

Note 2 to entry: More details can be found in 6.7.1.

[SOURCE: ISO 21448:2022, 3.12, modified – The EXAMPLEs and Notes to entry have been removed. A new Note to entry has been added for clarification.]

### 3.4.7

#### **output insufficiency**

incorrect output of an element as a result of a triggering condition activating a *functional insufficiency* (3.4.5) of the element, contributing to either a hazardous behaviour, a *RFIM prevention issue* (3.3.9) or both

[SOURCE: ISO 21448:2022, 3.8, modified – The term was taken from Note 6 to entry and explicitly defined.]

### 3.4.8

#### **safety-related AI error**

*AI error* (3.4.1) of a *safety-related AI element* (3.1.33)

### 3.4.9

#### **safety-related fault**

fault of a *safety-related AI element* (3.1.3)

### 3.4.10

#### **systematic error**

error due to a systematic fault

## 3.5 Miscellaneous definitions

### 3.5.1

#### **architectural measure**

technical solution implemented by the *AI element* (3.1.3) to detect and mitigate or tolerate *AI errors* (3.4.1) in order to uphold the ability to execute the *AI task* (3.1.18) in a safe manner or to achieve or maintain a dedicated operating mode in case of AI errors without unreasonable *risk* (3.3.10)

EXAMPLE 1 Addition of output layers in the AI model for classification. AI models can make incorrect predictions that can lead to hazardous behaviour. Therefore, it would be beneficial for a model to be cautious in situations where it is uncertain about its predictions. One way to accomplish this is to design AI models by adding output layer(s) to represent reject class or reject option. Such models assess their confidence in each prediction and have the option to abstain from making a prediction when they are likely to make incorrect predictions.

EXAMPLE 2 Addition of redundant *AI components* (3.1.1).

Note 1 to entry: Architectural measures have a tangible impact on the *AI system* (3.1.17) or AI component and can enhance or modify the architecture of AI system or AI component.

Note 2 to entry: If an architectural artefact is used during development (e.g. using a saliency map to argue the explainability of the system) but is removed for system deployment, this is not considered an architectural measure.

### 3.5.2

#### **component**

non-system level element that is logically or technically separable and is comprised of more than one hardware part or one or more software units or a combination of hardware part(s) and software unit(s)

EXAMPLE A microcontroller.

Note 1 to entry: A component is a part of a system.

[SOURCE: ISO 26262-1:2018, modified — "Or a combination of hardware part(s) and software unit(s)" was added.]

### 3.5.3

#### **development measure**

appropriate process step(s) (activity) for the development of an *AI system* (3.1.17) or an *AI component* (3.1.1) that facilitates fulfilling *AI safety requirements* (3.1.16) and/or enhancing the AI properties

Note 1 to entry: When analysing an AI system or AI component, specific activity used during or after the training of the AI component can be a development measure. See 10.4 for more details.

### 3.5.4

#### **element**

system, components (system, hardware or software), hardware parts or software units

[SOURCE: ISO 26262-1:2018, modified — "System" added to the components and both Note 1 and Note 2 to entry removed.]

### 3.5.5

#### **item**

system or combination of systems, to which ISO 26262 is applied, that implements a function or part of a function at the vehicle level

[SOURCE: ISO 26262-1:2018, modified — Note 1 to entry removed.]

### 3.5.6

#### **off-board**

property indicating that a given task is external to the vehicle system

### 3.5.7

#### **on-board**

property indicating that a given task is internal to the vehicle system

### 3.5.8

#### **testing**

process of planning, preparing and operating or exercising an item or element to verify that it satisfies specified requirements, to detect safety anomalies, to validate that requirements are suitable in the given context and to create confidence in its behaviour

Note 1 to entry: "To create confidence in its behaviour" also includes evaluating the performance of the element or item.

[SOURCE: ISO 26262-1:2018, modified — Note 1 to entry was added]

**3.5.9  
validation**

confirmation, through the provision of objective evidence, that the requirements for a specific intended use or application have been fulfilled

Note 1 to entry: There are three terms, "AI model validation", "validation" and "safety validation", that are distinguished in this document. AI model validation originates from the validation data used by the AI community, validation originates from classic system development and safety validation originates from ISO 26262. The three validation meanings are not the same.

[SOURCE: ISO/IEC 22989:2022, 3.5.18, modified — Note 1 to entry was added]

**3.5.10  
verification**

confirmation, through the provision of objective evidence, that specified requirements have been fulfilled

EXAMPLE The typical verification activities can be classified as follows:

- verification review, walk-through, inspection;
- verification testing;
- simulation;
- prototyping;
- analysis (safety analysis, control flow analysis, data flow analysis, etc.)

Note 1 to entry: Verification only provides assurance that a product conforms to its specification.

[SOURCE: ISO/IEC 22989:2022, 3.5.18, modified — EXAMPLE was added]

**4 Abbreviated terms**

ACP	assurance claim point
ADS	automated driving system
AI	artificial intelligence
ASIL	automotive safety integrity level
DFA	dependent failure analysis
DLC	dataset lifecycle
DNN	deep neural network
E/E	electrical/electronic
FMEA	failure mode and effects analysis
FN	false negative
FP	false positive
FPS	frames per second
GSN	goal structuring notation
HARA	hazard analysis and risk assessment

HAZOP	hazard and operability study
HMI	human machine interface
KPI	key performance indicator
ID	in distribution
ML	machine learning
NN	neural network
ODD	operational design domain
OOD	out of distribution
OTA	over the air
PFD	probability of failure on demand
PWM	pulse width modulation
RFDM	reasonably foreseeable direct misuse
RFIM	reasonably foreseeable indirect misuse
SOTIF	safety of the intended functionality
TN	true negative
TP	true positive

## 5 Requirements for conformity

### 5.1 Purpose

This clause describes how to:

- achieve conformity to this document;
- interpret the applicability of each clause;
- interpret the tables and figures used in this document.

### 5.2 General requirements

When claiming conformity to this document, each requirement shall be met, unless one of the following applies:

- tailoring of the safety activities as defined in [Clause 6](#) or in accordance with ISO 26262-2 has been performed that shows that the requirement does not apply;
- a rationale is available that the non-conformity is acceptable, and the rationale has been evaluated in accordance with this document and ISO 26262-2, when applicable.

The results of safety activities are given as work products. "Prerequisites" are information which shall be available as work products of a previous phase or from an external source. Given that certain requirements of a clause depend on the automotive safety integrity level (ASIL) or may be tailored, certain work products

may not be needed as prerequisites. A summary of the normative parts of this document is provided in [Annex A](#).

NOTE External sources are used in this document to refer to work products that are the result of activities outside of the defined AI safety lifecycle, such as input received from the development process of the encompassing system and interpretations of tables and figures.

Tables and figures can be normative or informative depending on their context. Tables and figures that are referenced by normative requirements are considered normative unless it is explicitly specified otherwise. Other tables and figures are only informative.

If it is possible to fulfil a requirement with a different combination of methods, a rationale is provided that the chosen combination of methods fulfil the requirement.

## 6 AI within the context of road vehicles system safety engineering and basic concepts

### 6.1 Application of the ISO 26262 series for the development of AI systems

This document is intended to be applied in combination with the ISO 26262 series to specifically address the safety of AI systems.

- For AI components that are not AI models, or do not contain AI models, the ISO 26262 series can be applied by itself.
- For AI components that are AI models, or that contain AI models, the ISO 26262 series can be tailored and applied in combination with this document (see [Figure 6-1](#)).

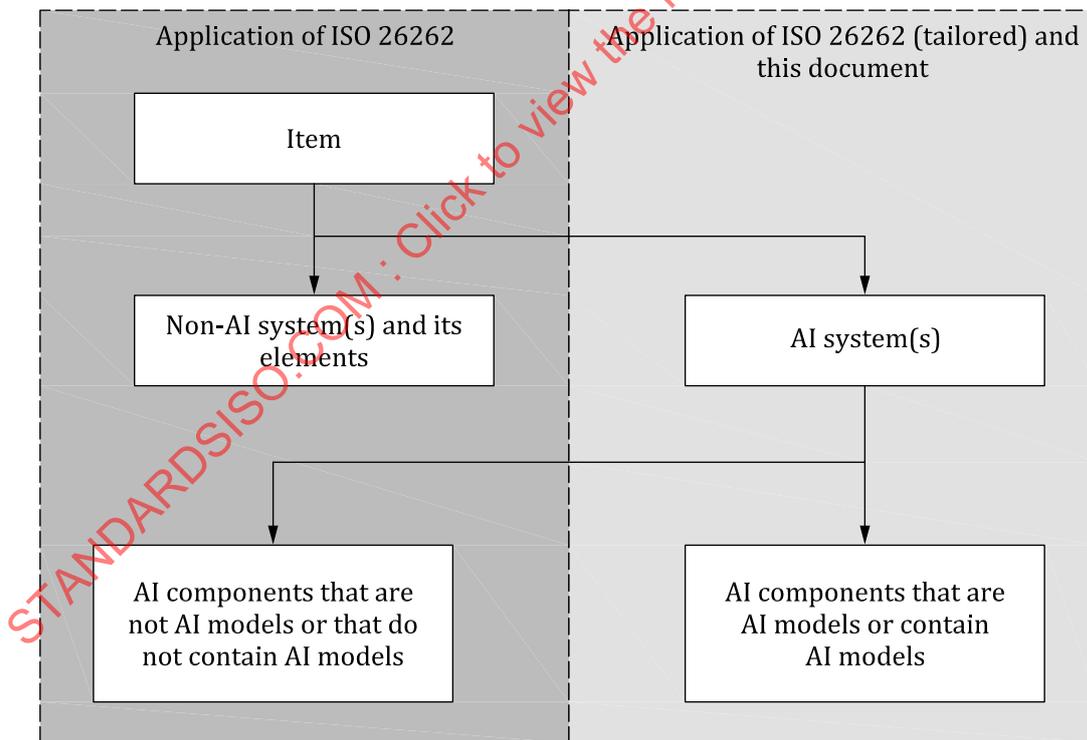


Figure 6-1 — Visualization of the applicability of the ISO 26262 series and this document to the item and its elements

NOTE See [Annex C](#) for a possible tailoring of ISO 26262-4 and ISO 26262-6 for ML.

## 6.2 Interactions with encompassing system-level safety activities

An example interaction of AI system development with the encompassing system development based on the ISO 26262 series and ISO 21448 can be found in [Table 6-1](#) and in [Table 6-2](#). These tables also include remarks regarding the interaction and the applicability of the corresponding standards with this document and the development of the AI elements.

NOTE 1 ISO 21448 is applicable to intended functionalities where proper situational awareness is essential to safety and where such situational awareness is derived from complex sensors and processing algorithms. This applies in particular to functionalities of emergency intervention systems and systems that have SAE levels of driving automation from 1 to 5. However, this document applies to all AI systems whose errors can impact safety independent of the vehicle-level functionality. Systems can utilize AI technologies that do not fall within the scope of ISO 21448 but that do fall within the scope of this document.

NOTE 2 Although this document and ISO 21448 both focus on functional insufficiencies, conformity to one does not automatically imply conformity to the other.

During the encompassing system architecture design phase, encompassing system requirements are decomposed and allocated to the AI systems as well as other elements of the encompassing system.

**Table 6-1 — Example interaction of the AI element development in conformity to this document and with the ISO 26262 series**

ISO 26262:2018	Interaction of the AI system with the encompassing system (the AI system is not an item)	AI system: System component consisting of hardware and software components	AI component: Conventional <sup>a</sup> hardware component	AI component: AI model implemented by software component(s)
Part 2, Clause 5: Overall safety management	-	Adapted to also address the management of AI safety	Directly applicable	Adapted to also address the management of AI safety
Part 2, Clause 6: Project dependent safety management	The management of the AI safety is part of the safety management of the encompassing system	Adapted to also address the management of AI safety	Directly applicable	Adapted to also address the management of AI safety
Part 2, Clause 7: Safety management regarding production, operation, service and decommissioning	The management of the AI safety is part of the safety management of the encompassing system	Adapted to also address the management of AI safety	Directly applicable	Adapted to also address the management of AI safety
Part 3, Clause 5: Item definition	Potential source of input space specification	-	-	-
Part 3, Clause 6: Hazard analysis and risk assessment	-	-	-	-
Part 3, Clause 7: Functional safety concept	Potential source of safety requirements allocated to the AI system	-	-	-
Part 4, Clause 6: Technical safety concept	Potential source of safety requirements allocated to the AI system	Applicable (tailoring can be necessary)	Hardware safety requirements are derived from technical safety requirements allocated to the AI system	Software safety requirements are derived from technical safety requirements allocated to the AI system

<sup>a</sup> Conventional hardware is hardware that is not specifically designed to implement an AI model, e.g. CPUs, GPUs or FPGAs.

Table 6-1 (continued)

ISO 26262:2018	Interaction of the AI system with the encompassing system (the AI system is not an item)	AI system: System component consisting of hardware and software components	AI component: Conventional <sup>a</sup> hardware component	AI component: AI model implemented by software component(s)
Part 4, Clause 7: System and item integration and testing	AI system as a system component to be integrated into the encompassing system	Integration of the hardware and software components of the AI system (tailoring can be necessary)	-	-
Part 4, Clause 8: Safety validation	Potential source of additional validation strategies and requirements	-	-	-
Part 5: Product development at the hardware level	Potential source of hardware safety requirements allocated to the AI elements	Applicable (tailoring can be necessary)	Applicable	Refinement of hardware-software interface
Part 6: Product development at the software level	Potential source of software safety requirements allocated to the AI elements	Applicable (tailoring can be necessary)	Refinement of hardware-software interface	Applicable (tailoring can be necessary)
Part 7: Production, operation, service and decommissioning	AI elements can be part of the production process	Potential source of requirements and work products relevant for production, operation, service and decommissioning	Potential source of requirements and work products relevant for production, operation, service and decommissioning	Potential source of requirements and work products relevant for production, operation, service and decommissioning

<sup>a</sup> Conventional hardware is hardware that is not specifically designed to implement an AI model, e.g. CPUs, GPUs or FPGAs.

Table 6-2 — Example interactions with ISO 21448

ISO 21448:2022, Clause	Interaction of the AI system development with encompassing system activities, motivated by ISO 21448:2022
Clause 5: Specification and design	<p>Clause 5 activities:</p> <ul style="list-style-type: none"> <li>— provide the interfaces of the AI system with the encompassing system;</li> <li>— determine the semantic input space;</li> <li>— provide the functionality required by the AI system;</li> <li>— provide safety requirements allocated to the AI system, including, but not limited to, safety-related KPIs.</li> </ul> <p>Activities of this document:</p> <ul style="list-style-type: none"> <li>— provide triggering conditions and functional insufficiencies of the AI system;</li> <li>— provide achieved safety-related KPIs;</li> <li>— provide a description of deployment measures required to support the AI and data lifecycles.</li> </ul>
Clause 6: Identification and evaluation of hazards	Clause 6 is a potential source of safety requirements, including, but not limited to, safety-related KPIs, allocated to the AI system.

Table 6-2 (continued)

ISO 21448:2022, Clause	Interaction of the AI system development with encompassing system activities, motivated by ISO 21448:2022
Clause 7: Identification and evaluation of potential functional insufficiencies and potential triggering conditions	Clause 7 activities: <ul style="list-style-type: none"> <li>— a potential source of safety requirements allocated to the AI system;</li> <li>— a potential source of AI triggering conditions.</li> </ul> Activities of this document: <ul style="list-style-type: none"> <li>— provide potential triggering conditions of the AI system.</li> </ul>
Clause 8: Functional modifications addressing SOTIF-related risks	Activities of this document: <ul style="list-style-type: none"> <li>— modification request to the encompassing system in case safety requirements allocated to the AI system cannot be fulfilled.</li> </ul>
Clause 9: Definition of the verification and validation strategy	Clause 9 is a potential source of safety-related KPIs, allocated to the AI system.
Clause 10: Evaluation of known scenarios	Activities of this document: <ul style="list-style-type: none"> <li>— provide triggering conditions of the AI system and the associated AI error modes, error patterns or a combination of both;</li> <li>— provide achieved safety-related KPIs.</li> </ul>
Clause 11: Evaluation of unknown scenarios	This document provides means to achieve safety-related KPIs.
Clause 12: Evaluation of the achievement of the SOTIF	This document provides a safety assurance argument for the AI system which can be used within the evaluation of the achievement of the SOTIF.
Clause 13: Operation phase activities	This document provides the AI system requirements to the encompassing system regarding the operation phase.

During development and as part of continuous assurance activities during operation, it can become necessary to adjust the safety requirements allocated to the AI system leading to an iterative feedback cycle to the encompassing system safety concept and safety requirements. Iterations of the requirements are triggered for example if:

- an AI system that is capable of fulfilling its assigned safety requirements and associated safety-related properties cannot be feasibly developed (e.g. due to inherent performance limitations in the ML algorithm used);
- suitable training data and test data cannot be found;
- evidence to demonstrate that the safety requirements and associated safety-related properties are fulfilled cannot be collected with sufficient confidence.

In each of these cases, changes to the encompassing system safety concept can be defined, leading to a set of updated, realisable requirements on the AI system.

NOTE 3 Measures on the encompassing system safety concept to reduce the safety load of an AI system towards better feasibility can be (see also [10.5](#)):

- restrictions in the ODD;
- implementation of diversity such as different processing algorithms or sensing modalities;
- implementation of redundancy such as multiple hardware components in parallel;

- a combination of the aforementioned measures.

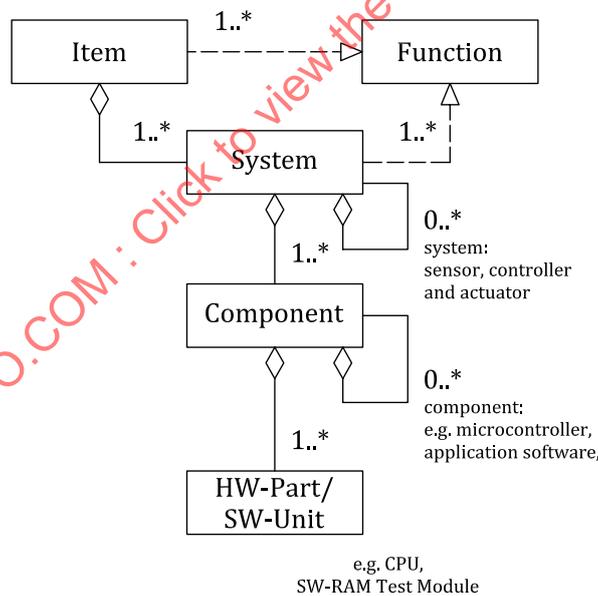
Once the AI system has reached an adequate level of performance in relation to the safety requirements, the AI system can be integrated into the encompassing system including evidence to support the achievement of the safety requirements. This can result in the need for further iterations of the AI safety life cycle, for example, due to the following conditions:

- Integration tests of the encompassing system reveal previously undiscovered faults or functional insufficiencies in the AI system that require additional development cycles.
- The encompassing system assurance case requires additional evidence to support safety claims related to the AI system that require additional effort to collect the required evidence.

Collected field data and observations made during operation related to the performance of the AI system (e.g. increased number of false positive errors under certain traffic conditions or an increased rate of out-of-distribution inputs) can indicate changes in the input space. It may not be possible to address these changes by refining the safety requirements on the AI system or through additional development activities. Changes may need to be made at the encompassing system level. This can lead to changes in the safety requirements assigned to the AI system and a repetition of the safety life cycle.

### 6.3 Mapping of abstraction layers between the ISO 26262 series, ISO/IEC 22989 and this document

The ISO 26262 series uses the following levels of abstraction: item, system, component, software unit and hardware part. The relationship between these is visualized in [Figure 6-2](#) (ISO 26262-10:2018, Figure 3). It is typically used when a given requirement can be applied on different levels of abstraction, e.g. on hardware components as well as on hardware parts. An example item composition is shown in [Figure 6-3](#) (ISO 26262-10:2018, Figure 4).



**Key**

- ↑ realization: one instance is realized by another instance (e.g. a function or part of a function is realized by an item)
- ◇ aggregation: one instance has a set of other instances (e.g. a system has a set of components)

NOTE 1 Depending on the context, the term “element” can apply to the entities “system”, “component”, “hardware part” and “software unit” in this chart, according to ISO 26262-1:2018, 3.41.

NOTE 2 “\*” means N elements are possible, where N is a positive integer number.

**Figure 6-2 — Relationship of item, system, component, hardware part and software unit**

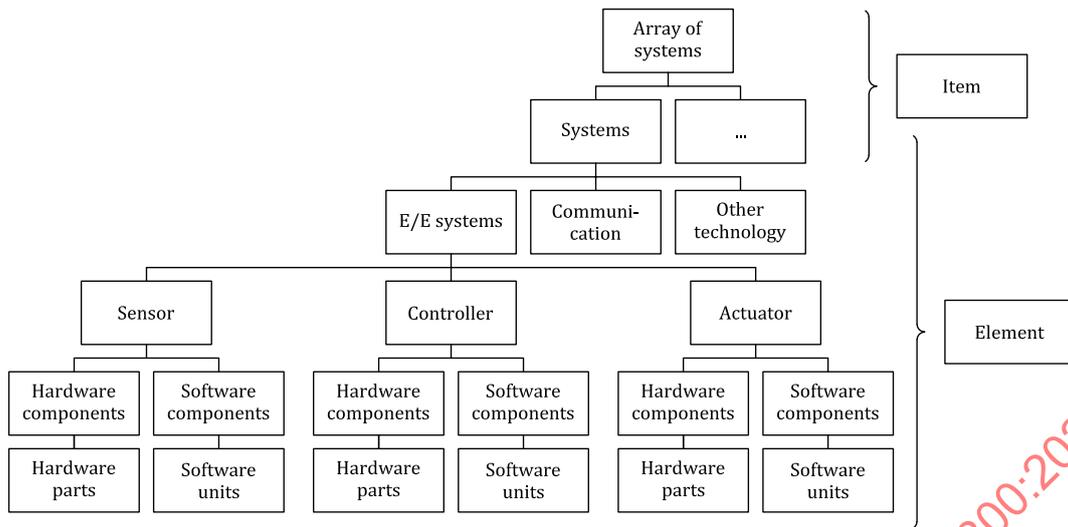
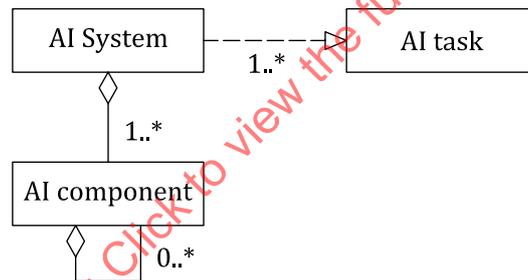


Figure 6-3 — Example item composition

ISO/IEC 22989 uses the following abstraction layers: AI system and AI components, where the AI system consists of AI components. ISO/IEC 22989 does not explicitly state if a given AI component can itself consist of AI components. In this document this is possible. The AI component can be an AI model, a conventional element, i.e. an element not considered to be an AI model, or a combination of both. The AI system contains at least one AI model and realizes the AI task. Figure 6-4 uses the same notation as Figure 6-2 to visualise the relationship between AI system and AI components.



Key

- ↑ realization: one instance is realized by another instance
- ◇ aggregation: one instance has a set of other instances

NOTE 1 Depending on the context, the term “AI element” can apply to “AI system” and “AI component”.

NOTE 2 "\*" means N elements are possible, where N is a positive integer number.

Figure 6-4 — Relationship of AI system and AI component

This document uses terms from both ISO/IEC 22989 and ISO 26262-1. The terms from the different standards do not map one-to-one. So, depending on the context, multiple mappings are possible as shown in Table 6-3.

**Table 6-3 — Possible mappings between ISO/IEC 22989 and ISO 26262-1 terms, depending on the context**

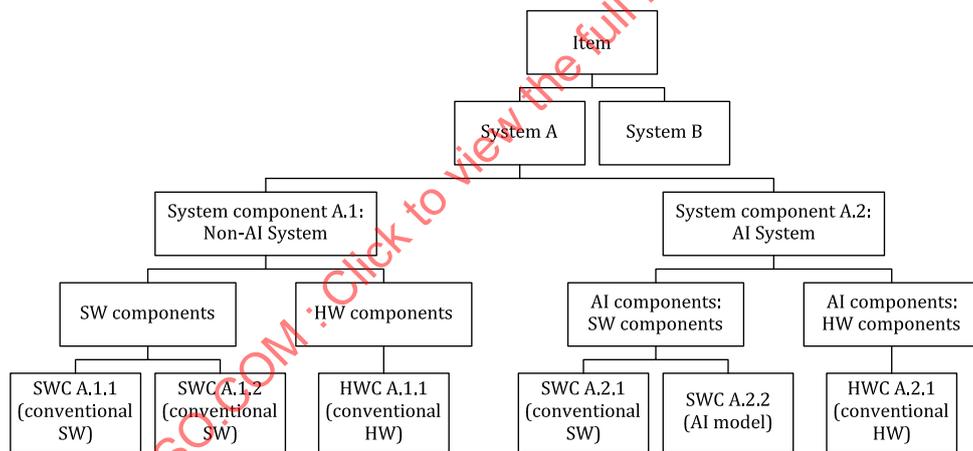
ISO/IEC 22989 terminology	ISO 26262-1 terminology
AI system	Item, system, component, software unit or hardware part Element
AI component	System, component, software unit or hardware part Element

In this document, the term “AI system” is the top level of abstraction of the content to be developed. As such it is possible in a distributed development that what one party considers to be an AI component, the other party considers to be an AI system, as for them it represents the top level of the content they develop.

Figure 6-5 provides an example of an item decomposed into its elements. In this example, the item consists of two systems: system A and system B. For the sake of simplicity, system B is not further decomposed. System A is composed out of the system components A.1 and A.2. System component A.1 does not contain an AI model. As such it cannot be an AI system. System component A.2 contains an AI model and is declared to be the AI system in this example.

**NOTE** It would have been possible to declare system A as the AI system, as it too contains at least one AI model. The decision regarding the scope of the AI system is made as part of negotiations between the development organisations responsible for the encompassing system.

The system components themselves consist of hardware and software components. System component A.2 is considered to be an AI component as it composes the AI system. A further breakdown into software units and hardware parts of the components has been omitted for the sake of simplicity.



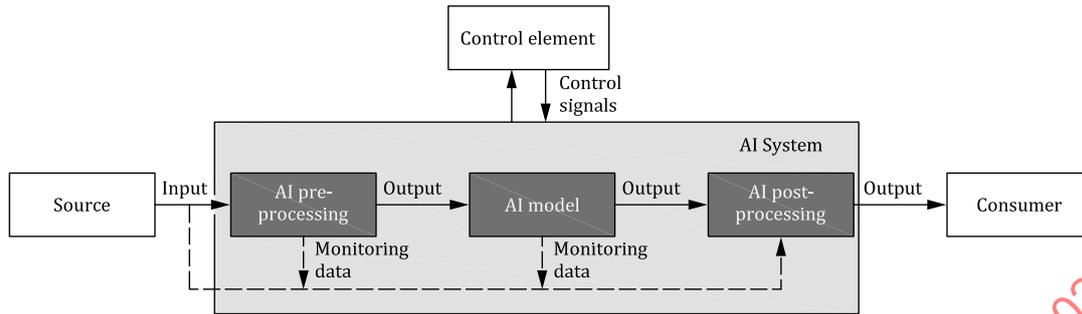
**Figure 6-5 — Example of a hierarchical decomposition of an item into its elements down to the component level - decomposition tree view**

## 6.4 Example architecture for an AI system

This document uses the architecture shown in Figure 6-6 as an example architecture. The AI system receives its input from the source, executes its task based on the input and the control signals and then provides its output to the consumer. The AI system itself consists of the AI components AI pre-processing, AI model and AI post-processing. The AI post-processing uses data provided from the previous process steps (i.e. AI pre-processing and the AI model) in combination with the original input data for monitoring purposes.

**NOTE 1** This architecture is just an example and has no claim of representing all possible architectures of AI systems.

NOTE 2 If the AI system is implemented for a real-time task, the execution of the AI system can be triggered synchronously or asynchronously depending on the behaviour of the control element, sources of the input streams and the consumers of the outputs. These considerations can be relevant to the definition of the safety requirements on the AI system.



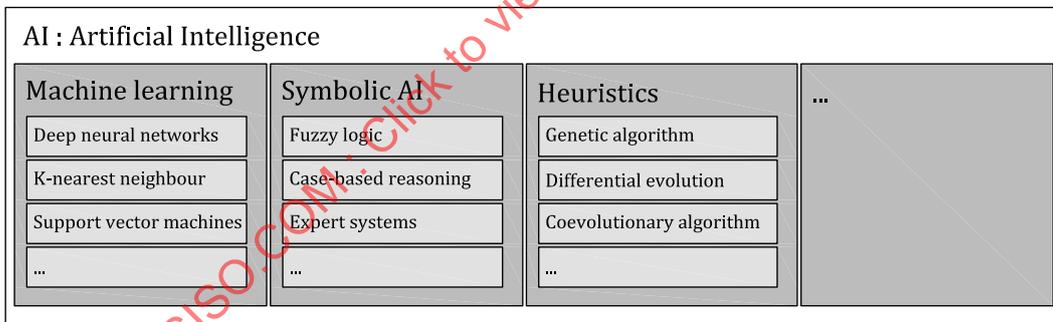
**Key**

- AI system
- AI component
- element not belonging to this AI system

**Figure 6-6 — Example architecture of an AI system**

**6.5 Types of AI models**

Examples for types of AI models include, but are not limited to, deep neural networks, k-nearest neighbours, support vector machines, decision trees, symbolic AI and fuzzy logic. These can be clustered in different categories. For example, deep neural networks, k-nearest neighbours, support vector machines and decision trees can be categorized as ML models as shown in [Figure 6-7](#).



**Figure 6-7 — Example of different types of AI models**

**6.6 AI technologies of a ML model**

[Figure 6-8](#) and [Figure 6-9](#) show an example of possible AI technologies utilized for an ML model that is implemented in hardware and software. In addition to the AI method itself, the AI technology also contains the tools and procedures to generate the AI model.

## ISO/PAS 8800:2024(en)

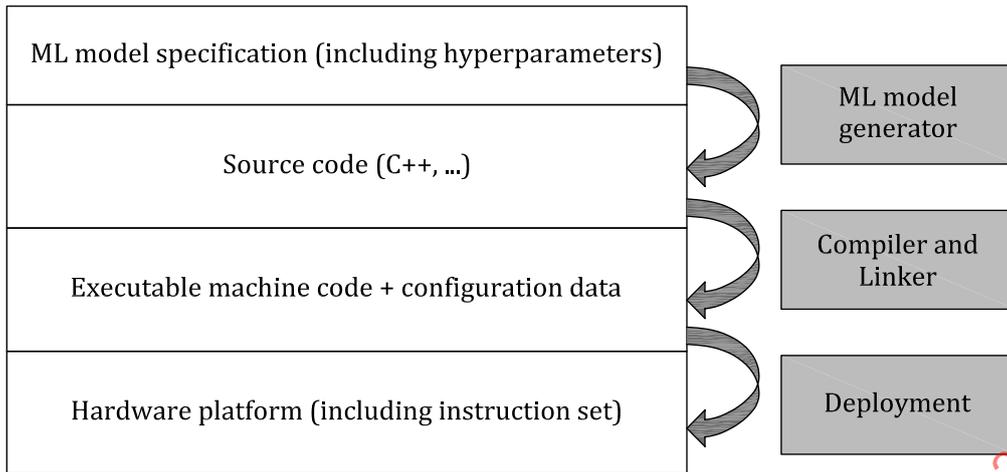


Figure 6–8 — AI technologies to create an executable ML model (application of ISO 26262)

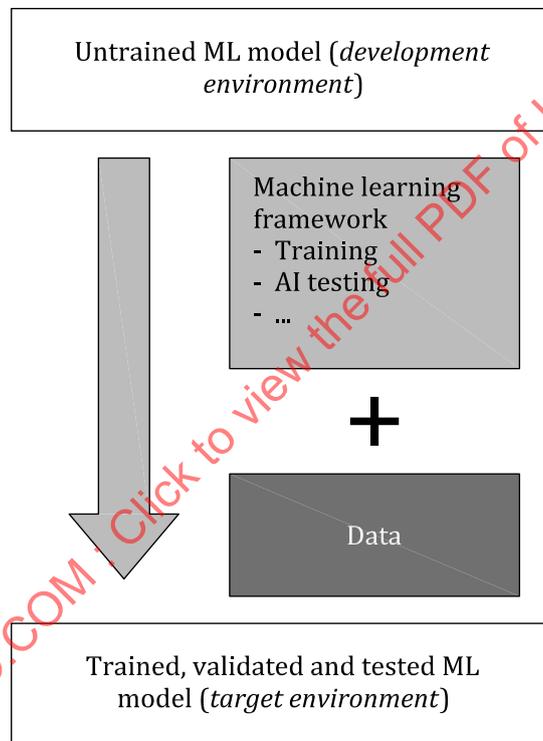


Figure 6–9 — AI technologies to create a trained ML model (application of this document)

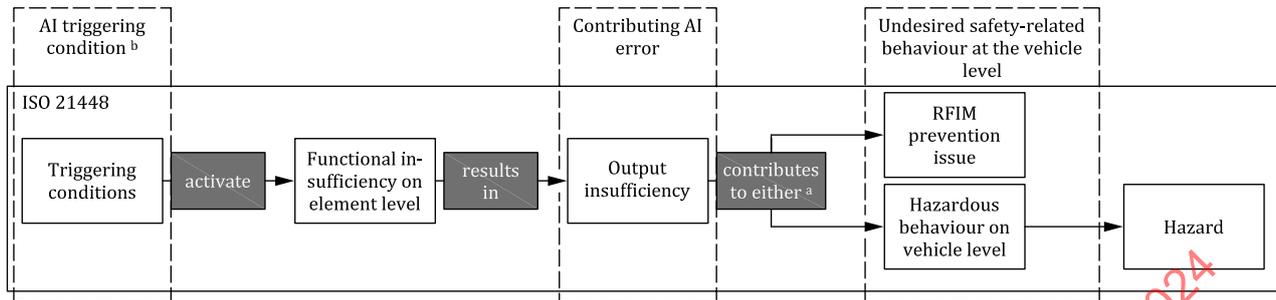
The AI technologies listed in [Figure 6–8](#) are not considered to be relevant sources for functional insufficiencies, e.g. compiler and linker are well known technologies already utilized by non-AI systems. For these technologies, it is enough to apply the ISO 26262 series to achieve AI safety. The AI technologies listed in [Figure 6–9](#) are considered relevant sources of functional insufficiencies for which the application of the ISO 26262 series alone is not considered to be sufficient to achieve AI safety. For these technologies, the remaining clauses of this document are applied.

## 6.7 Error concepts, fault models and causal models

### 6.7.1 Cause-and-effect chain

This document utilizes the concept of AI triggering conditions, faults, functional insufficiencies, AI errors and the undesired safety-related behaviour at the vehicle level. The mapping of the terms of this document

to the cause-and-effect chain used by ISO 21448 can be found in [Figure 6-10](#). The mapping of the terms of this document to the cause-and-effect chain used by ISO 26262 series can be found in [Figure 6-11](#). The AI triggering condition activates a fault or a functional insufficiency, resulting in an AI error in case of a fault and a contributing AI error in case of a functional insufficiency.

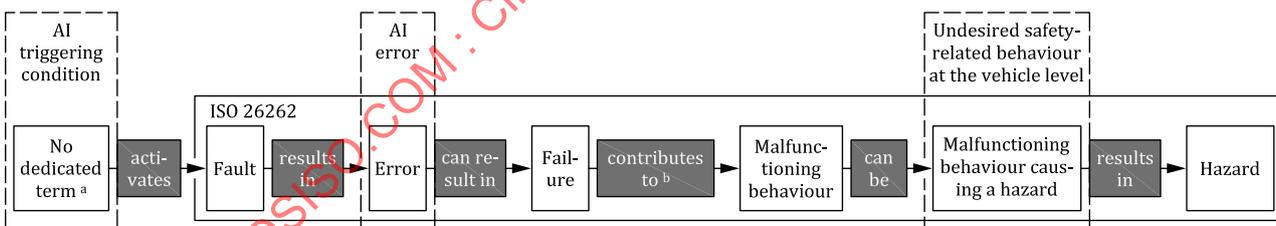


- a An output insufficiency, either by itself or in combination with one or more output insufficiencies of other elements, contributes to either a hazardous behaviour at the vehicle level or an inability to prevent or detect and mitigate a reasonably foreseeable indirect misuse.
- b Since a triggering condition of ISO 21448 results in a contributing AI error in the context of this document, they represent a subset of all AI triggering conditions.

**Figure 6-10 — Mapping of the cause-and-effect chain of ISO 21448 to the terms of this document**

**EXAMPLE** An insufficiency of specification can be a missing object in the AI training, AI validation and AI test dataset of an AI system utilizing an ML model for object classification. In this case, encountering this object during operation in the field is the AI triggering condition that activates this insufficiency of specification, resulting in the occurrence of a contributing error. The contributing AI error would be the incorrect classification of the object by the ML model and consequently by the AI system.

An AI error of an AI component can propagate through the AI system and can result in an AI error of the AI system. The AI error of the AI system can propagate through the encompassing system and can contribute either by itself or in combination with one or more other errors or output insufficiencies of the elements of the encompassing system to an undesired safety-related behaviour at the vehicle level.



- a In the ISO 26262 series, there is neither a dedicated term for the condition that activates a fault nor is this concept explicitly utilized.
- b A failure, either by itself or in combination with one or more failures of other elements, contributes to a malfunctioning behaviour.

**Figure 6-11 — Mapping of the cause-and-effect chain of the ISO 26262 series to the terms of this document**

The undesired safety-related behaviour at the vehicle level is used as an umbrella term for the corresponding terms of the ISO 26262 series (i.e. the malfunctioning behaviour at the vehicle level which can cause hazards) and ISO 21448 (i.e. the hazardous behaviour at the vehicle level and RFIM prevention issue).

6.7.2 Root cause classes

Different error classes can be distinguished depending on the root cause. The correlation between contributing AI errors and their different root causes is shown in Figure 6-12.

NOTE 1 A contributing AI error of the AI element can lead to the undesired safety-related behaviour at the vehicle level by itself or in combination with one or more other AI errors.

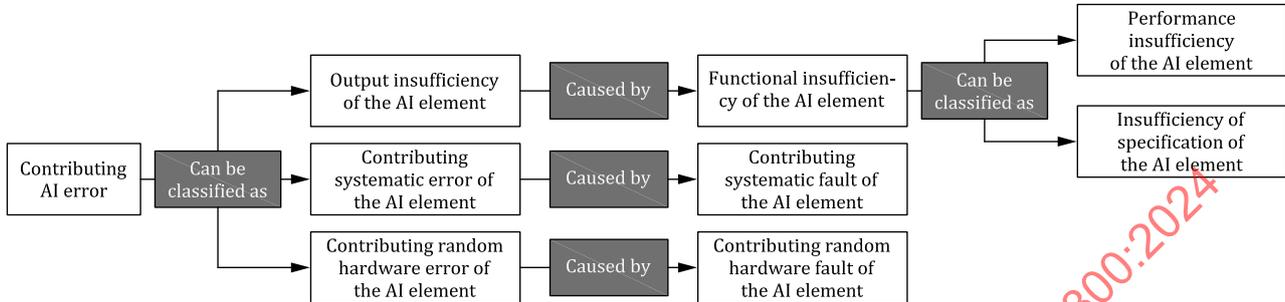


Figure 6-12 — Correlation of safety-related errors with their different classes of root causes

The root causes for the different kinds of AI errors are:

- insufficiency of the specification

EXAMPLE 1 An insufficiency of the specification can be missing datasets in the AI training, AI validation or AI test dataset. The resulting output insufficiency can be a misclassification when exposed to the missing datasets.

EXAMPLE 2 Specification of a neural network model with insufficient complexity.

EXAMPLE 3 An inadequate training loss function.

EXAMPLE 4 Inadequate labelling specification.

- performance insufficiency

EXAMPLE 5 A performance insufficiency can be an insufficient range of a sensor in case of certain environmental conditions. The resulting output insufficiency can be a false negative detection of an obstacle in the trajectory.

NOTE 2 In the case of ML models, performance insufficiencies can be specifically caused by training and test dataset related issues, e.g. insufficiencies in the coverage of the respective input space. These data-related issues are in turn considered to be insufficiencies of specification, or more precisely as insufficiency of specification of the data.

- contributing systematic fault

EXAMPLE 6 A contributing systematic fault can be to divide by zero in the software or to use incorrect variable names.

EXAMPLE 7 Overfitting the DNN resulting in wrong high-confidence classification outputs of corner cases can be regarded to be a contributing systematic fault in the training procedure.

NOTE 3 Sometimes the classification of a given issue in either a systematic fault or a functional insufficiency can be ambiguous. Independent of the classification, a safety assurance argument is provided to argue that this issue does not represent an unreasonable risk. As long as this safety assurance argument is available, the exact classification is not relevant.

- contributing random hardware fault

EXAMPLE 8 Physical defect causing a short to ground.

When evaluating the effectiveness of safety mechanisms, it can be necessary to distinguish the different classes of errors.

EXAMPLE 9 Homogenous redundancy can be effective in detecting random hardware errors in one of the redundant elements, but it is not effective in detecting systematic errors.

The different fault classes can also be used within the safety assurance argument by addressing each root cause class with a dedicated set of safety measures.

EXAMPLE 10 Coding guidelines are a measure to avoid systematic faults in SW. In combination with other fault avoidance measures, e.g. as specified in ISO 26262-6, the absence of unreasonable risk due to systematic faults can be argued.

It can also be useful to classify the AI triggering conditions in different categories.

EXAMPLE 11 For ML based AI models, the AI triggering conditions can be distinguished as:

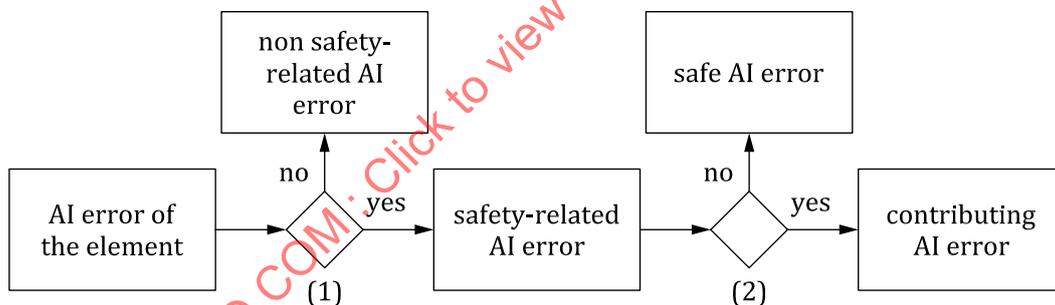
- cases that are similar to known training data samples ("in-distribution" cases), such as: cases at decision boundaries ("hard" cases), and undecided cases (high aleatoric uncertainty, e.g. due to label noise);
- unseen cases ("out-of-distribution" cases), such as: novel objects (semantic shift), or novel image styles (covariate shift);
- addition of small changes to a sample which causes no error (e.g. addition of adversarial crafted perturbation);
- applying a maliciously inserted trigger pattern (e.g. inserted via data or model poisoning).

Statistically obtained ML models like deep neural networks exhibit for any output an inherent uncertainty about their correctness. Therefore, besides classification of AI errors, root causes, and AI triggering conditions, a classification of uncertainty types and their associated sources can also be helpful in practice.

EXAMPLE 12 The two major types used to model the uncertainty of the ML model are epistemic and aleatoric uncertainty. While epistemic uncertainty stems from uncertainty of having the right model for the given sample and can often be fixed with more data, aleatoric uncertainty stems from intrinsic noise in the training data.

### 6.7.3 Error classification based on the safety impact

Compared to the criticality classification of random hardware faults as described in ISO 26262-5, this document uses a simplified scheme as shown in Figure 6-13.



**Key**

- (1) Is the element a safety-related element?
- (2) Can the error significantly contribute to the occurrence of a safety-related undesired behaviour at the vehicle level?

**Figure 6-13 — Error classification scheme based on the potential to lead to an undesired safety-related behaviour at the vehicle level**

## 7 AI safety management

### 7.1 Objectives

The objectives of this clause are:

- a) to define an AI safety lifecycle and its activities to ensure that contributing errors of the AI system do not lead to unreasonable risk of undesired safety-related behaviour at the vehicle-level. The AI safety lifecycle includes:
  - 1) a definition of activities necessary to develop the AI system, to provide the assurance and the evidence that the AI system is safe and to ensure the AI safety during operation;
  - 2) in case of the utilization of ML based AI technologies: a data-driven, iterative approach for the development, evaluation and continuous assurance of AI system within the context of an encompassing system's safety lifecycle;
- b) to ensure that overall and project specific safety management processes and activities are appropriate to ensure the safety of the AI system;

NOTE ISO 26262-2 provides suitable guidance on overall safety management and project-specific safety management. This guidance can require extensions based on recommendations in this document.

- c) to plan, initiate and conduct the AI safety activities.

### 7.2 Prerequisites and supporting information

The following information shall be available (from external sources, e.g. the encompassing system development):

- a) the AI system definition, including:
  - 1) the AI system functionality;
  - 2) the interfaces of the AI system with the encompassing system, including if applicable, the ASIL capability of the inputs to the AI system;
  - 3) the safety requirements allocated to the AI system, including if applicable:
    - i) the ASIL rating of the safety requirements;
    - ii) the acceptance criteria or validation targets derived in conformity to ISO 21448:2022, Clause 6 or 9.

NOTE 1 The safety requirements allocated to the AI system from external sources are typically requirements regarding the avoidance or control of safety-related faults, the allowed maximum error occurrence rate of contributing AI errors, the identification of AI triggering conditions, and the robustness against certain environmental conditions.

NOTE 2 For an elaboration of the fault model, the causal model, and the error concepts used by this document, see [6.7](#).

### 7.3 General requirements

**7.3.1** An AI safety lifecycle shall be defined that specifies the activities necessary to develop the AI system, to provide the assurance and the evidence that the AI system is safe and to ensure the maintenance of AI safety during operation. It can be based on the reference lifecycle ([Figure 7-2](#)) and can be tailored according to project-specific needs. The tailoring is supported by a rationale for why the tailored AI safety lifecycle is appropriate to achieve AI safety.

**7.3.2** At each phase within the AI safety life cycle, work products shall be defined to support the safety assurance claims of the AI system.

**7.3.3** The activities of the AI safety lifecycle shall be coordinated with the safety lifecycle activities of the encompassing system as defined by ISO 26262-2 and, if applicable, ISO 21448.

NOTE The AI system can be developed as a safety element out of context. The concept of a safety element out of context is described in ISO 26262-10:2018, Clause 9.

**7.3.4** The activities described in ISO 26262-2 shall be adapted in order to address the management of AI safety, including:

- a) the integration of the AI safety lifecycle into the ISO 26262 series safety lifecycle (see ISO 26262-2:2018, Figure 2);
- b) the enhancement of “functional safety” to “AI safety”;
- c) measures to ensure that a sufficient level of cross domain competences regarding safety and AI are available, in conformity to ISO 26262-2:2018, 5.4.4.1;
- d) adding this document as a relevant standard of ISO 26262-2;
- e) the use of a safety assurance argument as part of the safety case of ISO 26262-2;
- f) extending the safety plan of ISO 26262-2 to include the safety activities of this document;
- g) tailoring ISO 26262-2:2018, Table 1, to address the work products of this document (see [Table 7-1](#)).

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

Table 7-1 — Required confirmation measures, including the required level of independence

Confirmation measure	Level of independence <sup>a</sup> applies to					Scope
	QM	ASIL A	ASIL B	ASIL C	ASIL D	
Confirmation review of the safety plan Independence with regard to the developers of the item, <sup>b</sup> project management and the authors of the work product	-	I1	I1	I2	I3	Applies to the highest ASIL among the safety requirements
Confirmation review of the AI system validation report Independence with regard to the developers of the item, <sup>b</sup> project management and the authors of the work product	-	I0	I1	I2	I2	Applies to the highest ASIL among the safety requirements
Confirmation review of the AI safety analyses Independence with regard to the developers of the item, <sup>b</sup> project management and the authors of the work product	-	I1	I1	I2	I3	Applies to the highest ASIL among the safety requirements

<sup>a</sup> The notations are defined as follows:

- -: no requirement and no recommendation for or against regarding this confirmation measure;
- I0: the confirmation measure should be performed; if the confirmation measure is performed, it shall be performed by a different person in relation to the person(s) responsible for the creation of the considered work product(s);
- I1: the confirmation measure shall be performed, by a different person in relation to the person(s) responsible for the creation of the considered work product(s);
- I2: the confirmation measure shall be performed, by a person who is independent from the team that is responsible for the creation of the considered work product(s), i.e. by a person not reporting to the same direct superior;
- I3: the confirmation measure shall be performed by a person who is independent, regarding management, resources and release authority, from the department responsible for the creation of the considered work product(s).

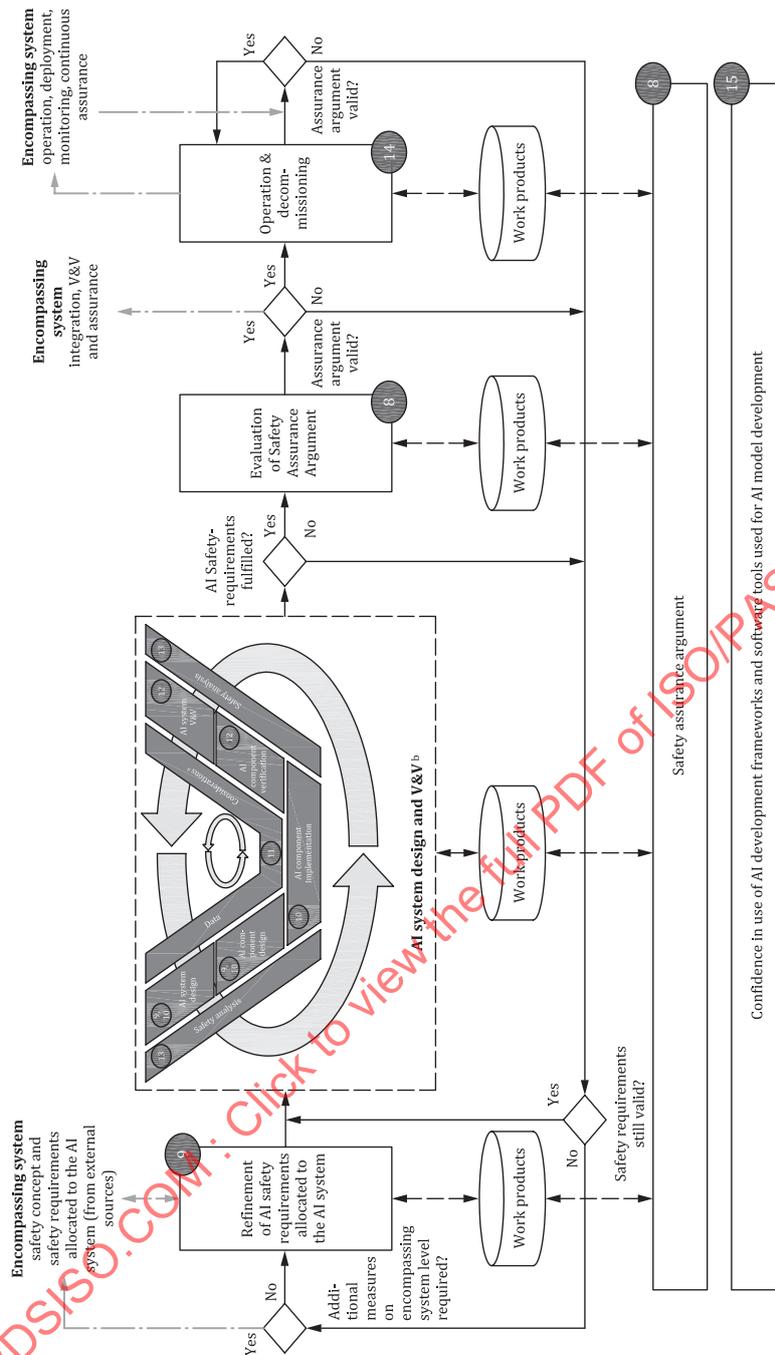
<sup>b</sup> The developers of the item include the developers of the AI system.

Table 7-1 (continued)

Confirmation measure	Level of independence <sup>a</sup> applies to					Scope
	QM	ASIL A	ASIL B	ASIL C	ASIL D	
Confirmation review of the safety assurance argument Independence with regard to the authors of the safety assurance argument	-	I1	I1	I2	I3	Applies to the highest ASIL among the safety requirements
AI safety audit Independence with regard to the developers of the item <sup>b</sup> and project management	-	-	I0	I2	I3	Applies to the highest ASIL among the safety requirements
AI safety assessment Independence with regard to the developers of the item <sup>b</sup> and project management	-	-	I0	I2	I3	Applies to the highest ASIL among the safety requirements
<sup>a</sup> The notations are defined as follows: — -: no requirement and no recommendation for or against regarding this confirmation measure; — I0: the confirmation measure should be performed; if the confirmation measure is performed, it shall be performed by a different person in relation to the person(s) responsible for the creation of the considered work product(s); — I1: the confirmation measure shall be performed, by a different person in relation to the person(s) responsible for the creation of the considered work product(s); — I2: the confirmation measure shall be performed, by a person who is independent from the team that is responsible for the creation of the considered work product(s), i.e. by a person not reporting to the same direct superior; — I3: the confirmation measure shall be performed by a person who is independent, regarding management, resources and release authority, from the department responsible for the creation of the considered work product(s). <sup>b</sup> The developers of the item include the developers of the AI system.						

**7.4 Reference AI safety life cycle**

The reference AI safety life cycle described in this clause covers the activities at the different phases of AI system development, deployment and operation: safety-related requirements derivation, AI system design, verification and validation, deployment and operation. The AI safety life cycle is summarised in [Figure 7-1](#) and is used to structure the remainder of this document. A detailed view of the AI system design and V&V phase is shown in [Figure 7-2](#). [Clause 8](#) through [Clause 15](#) (as indicated by the numbered black circles in [Figure 7-1](#) and [Figure 7-2](#)) are used to describe the activities within the safety life cycle in more detail.



Key



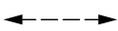
clause number(s)



process flow to/from the development and operation of the encompassing system



process flow



relation

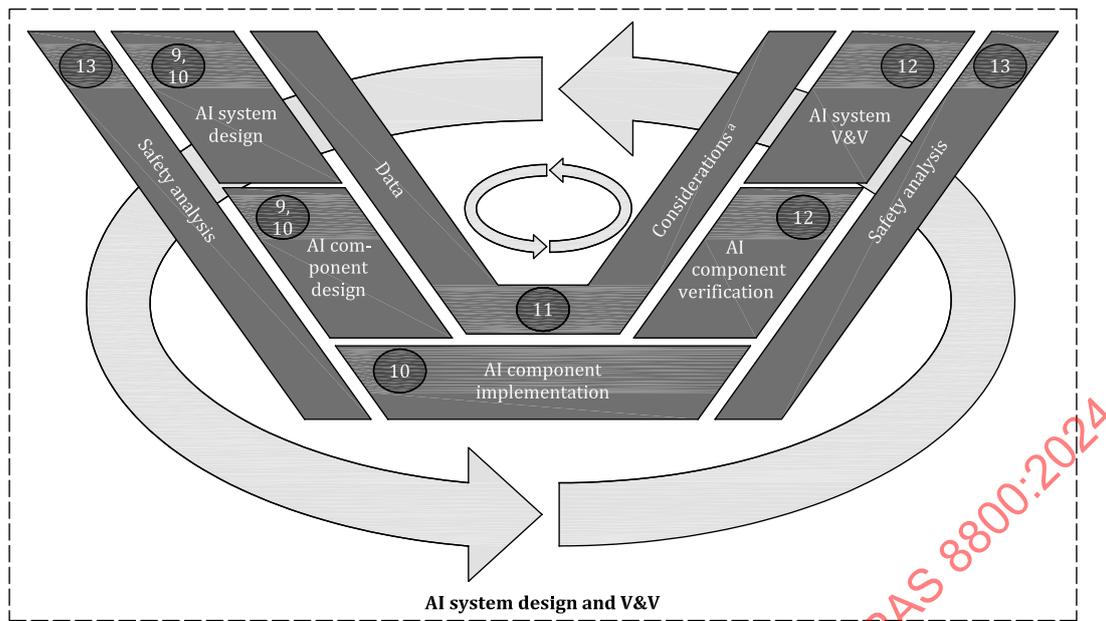


represents the iterative nature, in particular of the AI component design and verification

a Specific to ML-based AI technologies.

b See [Figure 7-2](#).

Figure 7-1 — Reference AI safety lifecycle



**Key**



clause number(s)



represents the iterative nature, in particular of the AI component design and verification

a

Specific to ML-based AI technologies.

**Figure 7-2 — Detailed view of the AI system design and V&V phase of the reference AI safety lifecycle**

## 7.5 Iterative development paradigms for AI systems

The AI system development phase of the AI safety life cycle covers all activities required to design, implement, verify and validate the AI system.

The AI system development activities described in [Clause 8](#) through [Clause 13](#) are iteratively performed until a sufficient level of performance with respect to the safety requirements and the associated safety properties can be demonstrated and associated work products are generated.

An essential characteristic of this development process is the analysis to identify potential functional insufficiencies, their root causes, and their impact on safety. This analysis is used to derive appropriate measures to reduce these functional insufficiencies during design (including the selection of training data in case of ML) as well as to reduce the impact of contributing errors (through architectural measures).

NOTE 1 During the AI system development, the property of “independence” is considered wherever appropriate (e.g. independence between training datasets and test datasets). An analysis similar to the dependent failure analysis (DFA) described in ISO 26262-9:2018, Clause 7 can be used.

The development model described here reflects the iterative development model typically used in the area of ML with a focus on activities to identify, analyse, reduce and mitigate functional insufficiencies in the trained model and the cumulative collection of evidence to support claims regarding safety requirements allocated to the function. Furthermore, for machine learning, the specification and collection of suitable training and test data is one of the most influential factors for the performance of the function. Therefore, the specification, planning, collection, acquisition, preparation and labelling of data related to AI component implementation and verification is treated as a safety-related development activity in this document, with specific objectives and associated safety artefacts (see [Clause 11](#)).

The iterative development of the AI system is guided by a set of performance indicators, including safety-related properties associated with the safety requirements allocated to the AI system.

NOTE 2 “Iteration” in the context of AI system development can be defined as a single repetition of one or more AI system safety lifecycle phases. It is applied until conditions defined by a set of KPIs or other target parameters are either fulfilled or demonstrated to be unachievable.

NOTE 3 Considering that there can be multiple, potentially conflicting target KPIs, the presence of contributing errors in the AI system can be inevitable.

Despite the inherently iterative nature of AI system development, different iteration cycles of the development can focus on different sets of performance indicators, depending on the maturity of the development.

Examples of such distinct iteration cycles can include:

- Proof of concept: In this phase of development, AI models are designed, implemented and tested against a set of initially defined safety requirements. The objective of this phase is to evaluate the potential of the chosen AI technologies to fulfil the safety requirements, as well as defining a set of measures required to minimise the number or the impact of functional insufficiencies in the function (e.g. optimisation of model parameters and defining a data collection strategy).
- Series development of the AI system: In this phase of development, the AI system is iteratively developed against all safety requirements. Errors of the AI system or the AI model are analysed with respect to their potential root causes. Measures are defined to reduce functional insufficiencies through design and data selection or to minimise their impact through architectural measures. This phase of development can use a host platform to implement and test the AI model and continue until an adequate level of performance for all safety-related properties has been met.
- Deployment to the target hardware platform: In this phase of development, the AI system is transferred to the target hardware platform and target software platform. This can include, for example a change in numerical precision used to calculate the results as well as the consideration of constraints such as timing, memory limitations (e.g. resulting in the need to prune the computational graph) and robustness against potential random hardware faults. The focus of this development phase is to ensure that the safety requirements are met, despite any limitations of the target hardware platform and the target software platform.
- Further improvement after deployment: In this phase of development, the AI system is iteratively updated based on observations made during operation in the field and new requirements mandated by the developer. This can include compensating for previously unknown triggering conditions (e.g. concept drift) and distributional shift in the environment (e.g. domain shift). Performance indicators within this phase of development are monitored to ensure a monotonic safety improvement with respect to previous iterations of the AI model. This phase of development can also include development versions of the function running in “shadow” mode within an operational environment in order to collect suitable data and evaluate the potential performance under realistic conditions.

## 7.6 Work products

7.6.1 AI safety lifecycle resulting from [7.3.1](#) to [7.3.4](#).

7.6.2 Work products of ISO 26262-2:2018, 5.5, resulting from [7.3.4](#).

7.6.3 Work products of ISO 26262-2:2018, 6.5, resulting from [7.3.3](#) and [7.3.4](#), in particular the safety plan.

7.6.4 Work products of ISO 26262-2:2018, 7.5, resulting from [7.3.4](#).

## 8 Assurance arguments for AI systems

### 8.1 Objectives

The objectives of this clause are:

- a) to develop an assurance argument demonstrating that the safety requirements allocated to the AI system are fulfilled;

NOTE 1 The assurance argument for the AI system contributes to the safety assurance argument of the encompassing system.

NOTE 2 The assurance argument can be developed independently from the encompassing system as a safety element out of context (SEooC) development activity (see ISO 26262-10:2018, Clause 9). In such cases, the assurance argument documents all necessary assumptions on the encompassing system for arguing that the safety requirements allocated to the AI system are fulfilled.

- b) to evaluate whether the assurance argument reflects the actual residual risk of the AI system violating its safety requirements.

### 8.2 Prerequisites and supporting information

The following information shall be available at the initiation of these activities:

- a) the AI system definition (from external sources), e.g. the encompassing system development, including:
  - 1) a specification of the safety requirements allocated to the AI system;
  - 2) a definition of the technical context within the encompassing system (e.g. definition of interfaces to and from the encompassing system and, if applicable, the environment, conditions under which the AI system functionality is triggered, etc.);

NOTE 1 This includes the ASIL capability and noise to signal ratio of the input signals provided by the source, if applicable.

  - 3) a specification of the input space;
- b) requirements on the assurance argument and work products for the AI system (from external sources). These requirements can be derived from the assurance argument of the encompassing system as well as safety management procedures from [Clause 7](#);

The following information shall be available for the finalization of these activities:

- c) the work products of the AI safety lifecycle;

NOTE 2 This body of evidence can be cumulatively collected as part of the iterative development process phases or produced within dedicated development process cycles.

The following information can be considered for the finalization of this phase:

- d) required properties of the encompassing system to achieve AI safety;
- e) relevant properties of the input space (e.g. distribution of critical events, physical constraints on changes of input values over time);
- f) evidence of organization-specific rules and processes for AI safety, evidence of competence management and evidence of a quality management system, from 7.6.2;
- g) evidence that the organization-specific rules and processes have been followed and that the work products have the required maturity and quality (see 7.6.3 and 7.6.4).

### 8.3 General requirements

**8.3.1** An assurance argument for the fulfilment of the safety requirements allocated to the AI system shall be provided.

NOTE 1 The assurance argument can be part of the encompassing system's safety case in accordance with ISO 26262-2:2018, 6.4.8.

NOTE 2 The assurance argument can be constructed at the level of the encompassing system. In this case, the development of the AI system-specific contributions and supporting evidence are considered part of the AI safety life cycle and therefore within the scope of this document.

**8.3.2** The assurance argument shall use the relevant work products generated during the AI safety lifecycle to support the assurance claims.

NOTE Changes to the work products and their impact on the assurance argument are considered as part of change management throughout the AI safety life cycle.

**8.3.3** Confirmation measures of [7.3.4](#) shall be applied to the assurance argument.

NOTE The evaluation of the validity of the assurance argument can be performed as part of confirmation measures of the encompassing system (see ISO 26262-2:2018, 6.4.9 and ISO 21448:2022, 12.3), or as an independent activity, for example in the case of a SEooC (see ISO 26262-10:2018, Clause 9) or as part of a distributed development.

### 8.4 AI system-specific considerations in assurance arguments

A number of AI system-specific considerations impact the creation of the assurance argument.

- a) The formulation of the AI safety requirements (see [Clause 9](#)):
  - These include quantitative properties expressed in the form of probabilities (e.g. proportion of false positive classifications).
  - Arguments are expressed that demonstrate that these properties have been achieved with a level of statistical confidence appropriate to the quantitative targets associated with the requirement's acceptance criteria.
  - This can lead to additional requirements on the nature of evidence to support the claim and how the validity of this evidence is evaluated.
- b) Statistical arguments related to aggregated performance metrics:
  - These might not be sufficient to argue a suitable level of safety in rare but critical situations (e.g. edge cases, sensor defects or adversarial perturbations).
  - Arguments can be required to demonstrate that such input conditions nevertheless lead to a suitable level of AI safety.
  - The probability of unknown triggering conditions leading to a violation of safety requirements can depend on:
    - features of the input space not directly related to the function (e.g. due to spurious correlations in the training data);
    - predictions based on past inputs (e.g. the accuracy of previous detections of dynamic objects can impact the future behaviour of a planning task).
- c) Verification of the AI system:
  - Direct introspective approaches of AI models might not be effective.

- Alternative means of arguing the correct behaviour of the AI model or an increased reliance on indirect verification (e.g. test) can be required.

EXAMPLE Due to the lack of transparency with respect to the individual contributions of the large number of parameters used in some ML models, introspective approaches to verification can have limited applicability.

d) Reliance on training and test data:

- In machine learning-based AI methods, the behaviour of the AI system as well as its verification and validation are predominantly reliant on the selection of suitable training and test datasets as well as the training procedures themselves.
- Dedicated assurance arguments for demonstrating how the data selection process supports the achievement of AI safety can be required (see [Clause 11](#)).
- These arguments consider the training process and associated tools (see [Clause 15](#)).

e) Conditions during operation:

- Conditions can occur during operation that invalidate the assurance argument due to the complex nature of the environment in which vehicles containing AI systems are deployed.
- These conditions might include distributional shift of the input space (e.g. new types of road vehicles, changes in road infrastructure), changes to the technical system (e.g. replacement or upgrade of sensors) or previously undiscovered unknown triggering conditions.
- A continual, periodic re-evaluation and adaptation of the assurance argument is therefore performed, including an impact analysis of which parts of the assurance argument and associated evidence are to be re-evaluated (see [Clause 14](#)).

NOTE The degree to which continual, periodic re-evaluation is required depends on the properties of the input space and operating environment. If the input space, its distribution and change of distribution over time (e.g. due to ageing) are well known, re-evaluation can be performed within regular software update activities.

## 8.5 Structuring assurance arguments for AI systems

### 8.5.1 Context of the assurance argument

Within the scope of this document, assurance relates to the claim that the AI system achieves AI safety. An assurance argument communicates the relationship between evidence and the AI safety requirements.

The level of confidence in the assurance argument should be appropriate to the required level of integrity (functional safety) and acceptance criteria assigned to the AI system within the context of the encompassing system.

NOTE A model-based graphical representation of the assurance argument can aid the communication and evaluation of the assurance argument. Examples of graphical notations for assurance arguments include the goal structuring notation (GSN)<sup>[22]</sup> and claims argument evidence (CAE)<sup>[23]</sup> based on the structured assurance case metamodel (SACM)<sup>[24]</sup>.

The structure of the assurance argument can appeal to one or a combination of the following perspectives:

- features of the implemented item (product argument);
- features of the development measures and assessment process (process argument);
- factors impacting the residual risk associated with the AI system (e.g. potential causes of insufficiencies and failure modes).

EXAMPLE 1 Process-focused aspects of the assurance argument for the AI system can include an argument for the appropriate tailoring of the AI safety life cycle and the effectiveness with which each activity has been performed, based on an evaluation of the work products developed in each phase.

EXAMPLE 2 A risk-oriented assurance structure can include an argument that all possible causes of contributing AI errors are identified (e.g. via safety analyses), and suitable countermeasures for each cause have been identified and implemented, either through specific development measures or dedicated architectural measures.

The assurance argument for the AI system begins with a claim that the safety requirements allocated to the AI system are achieved. This can include statements related to a reasonable level of residual AI errors with respect to the AI safety requirements and the target functionality of the AI system.

An explicit definition of the context, as well as relevant assumptions, increases the transparency of the assurance argument and limits the scope of the argument to the specific AI system, its technical context within the encompassing system and its operating conditions.

EXAMPLE 3 Examples of context information that can be referenced by the assurance argument include:

- definition of the technical system context of the encompassing system;
- definition of the set of environmental conditions and operating context for which the assurance argument is valid;
- potential causes of contributing AI errors considered as part of the assurance argument.

EXAMPLE 4 Examples of assumptions that might be discharged as separate arguments can include:

- assumptions on the usage and operational profile of the AI system;
- assumptions on the reliability of inputs to the AI system;
- assumptions on the fundamental performance potential of the chosen AI technology.

An example of an assurance argument for an AI system structured according to a strategy that addresses possible sources of insufficiencies can be found in [Annex B](#).

### 8.5.2 Categories of evidence

The following categories of evidence in the form of work products created during the AI safety life cycle can be considered for use within the assurance argument.

a) Addressing insufficiencies in the specification of the AI safety requirements:

- Evaluation of the completeness of the definition of the environmental conditions and operating context (input space). This is used to confirm completeness requirements on training and test datasets (see [Clause 9](#)).
- Evaluation of the validity of the AI safety requirements derived from the safety requirements allocated to the AI system (from external sources). This includes traceability to safety requirements allocated to the AI system and a review of the completeness and consistency of the safety-related properties used to define the AI safety requirements (see [Clause 9](#)).

b) Addressing performance insufficiencies in the design of the AI system:

- Justification for the selection of the chosen AI methods, AI technologies and AI system architecture. This can include references to performance benchmarks and analysis indicative of the fundamental potential of the chosen technology and AI system architecture to meet the safety requirements (see [Clause 10](#)).
- Evaluation of the effectiveness of architectural and development measures. This can include an evaluation of the ability of architectural and development measures (see [Clause 10](#)) to limit the impact of contributing AI errors in the AI model.

NOTE 1 These measures can include hyperparameter optimization (a development measure) as well as monitoring components (an architectural measure) that detect inconsistencies in the outputs and trigger a dedicated AI error reaction to ensure AI safety. This can include components that ensure a continuous availability of the functionality through redundancy and voting, and dynamic adaptation of vehicle behaviour based on the evaluated performance of the AI system.

- Evaluation of robustness against hardware and software faults during execution. This can include an evaluation of the impact of random hardware faults and systematic design faults (including software) on AI safety (see ISO 26262-5:2018 and ISO 26262-6:2018).
- Evaluation of the impact of differences between the development and the target execution environment.

NOTE 2 This supports the argument that AI safety is achieved (see [Clause 10](#)) under the condition that some parts of the evaluation were performed within a development environment, e.g. software-in-the-loop tests, or using synthetic data (see [Clause 12](#)).

- c) Suitability of AI training and AI test datasets. This can include an evaluation of the suitability of AI training and AI test datasets to achieve and demonstrate that the safety requirements have been fulfilled (see [Clause 11](#)).

NOTE 3 This includes evidence of the independence of the AI test datasets from training datasets.

- d) Evaluation of the fulfilment of the safety requirements. This demonstrates the extent to which the safety requirements are fulfilled and, where necessary, provides rationale (e.g. risk analysis) on the requirements that are not fulfilled, or where such fulfilment cannot be demonstrated. This can include a quantitative evaluation of functional insufficiencies in the AI system with respect to target metrics and safety-related properties used to define the requirements (see [Clause 9](#)). Approaches to collect this category of evidence can make use of real, synthetic or hybrid datasets (see [Clause 12](#)).
- e) Evaluation of the impact of AI errors. This can include an evaluation of specific properties of the AI system that can lead to AI errors and consequently hazardous behaviour of the system. This can be based on targeted testing and analysis approaches to evaluate the presence and magnitude of known causes of AI errors such as insufficient generalization capability and insufficient robustness. This evaluation is made based on an analysis of potential causes of insufficiencies and AI errors in the AI system (see [Clause 13](#)) and includes a definition of a set of suitable measures to address the AI errors.
- f) Addressing AI errors during operation:
- Identification and analysis of previously undiscovered AI errors. This includes the continual evaluation of the behaviour of the AI system during operation (see [Clause 14](#)).
  - AI errors discovered during operation are analysed to understand their criticality, and a set of mitigation measures are identified, including a repetition of relevant phases of the AI safety life cycle.
  - Re-evaluation of robustness against changes in the operating conditions over time (distributional shift). This supports the argument that the AI system maintains its safety-related properties despite reasonably expected changes in its deployment environment.

EXAMPLE An analysis of the resilience of the AI system to shifts in the distribution of its inputs or the effectiveness of architectural measures to detect out of training/test distribution conditions (see [Clause 14](#)).

## 8.6 The role of quantitative targets and qualitative arguments

Safety requirements allocated to the AI system (from external sources) can include quantitative risk acceptance criteria and validation targets (see ISO 21448:2022, Clause 6).

These quantitative targets are considered during the derivation of AI safety requirements and are used to define target metrics for the safety-related properties (see [Clause 9](#)).

A direct mapping between quantitative targets (e.g. accident rates) of the safety requirements allocated to the AI system and the safety-related properties of the AI system (e.g. robustness to small changes in inputs) might not be possible.

Safety analyses (see [Clause 13](#)) that evaluate the impact and potential causes of AI errors in the AI system can provide a qualitative argument that the residual risk of violation of quantitative targets defined in safety requirements allocated to the AI system is acceptably low.

A demonstration of the correlation between causes of AI errors, the safety-related properties and the fulfilment of the AI safety requirements increases confidence in the effectiveness of the safety analysis and thereby the associated assurance arguments and evidence.

EXAMPLE 1 The safety analysis hypothesises that an inability to generalise on inputs outside of the training distribution leads to an unacceptably high rate of AI errors under certain conditions. An out-of-distribution detection as a post-processing function is therefore proposed as an architectural measure. To argue the effectiveness of this measure, the assurance argument demonstrates both the achieved coverage of out-of-distribution inputs as well as the actual contribution of out-of-distribution inputs to the overall contributing AI error rate of the AI system. Thus, both the effectiveness and the appropriateness of the out-of-distribution detection as a safety measure are argued.

To ensure that evidence referenced by the assurance argument provides sufficient confidence in the fulfilment of safety requirements, the following assumptions can be supported with dedicated assurance arguments and additional evidence:

- a) The measurement targets are an adequate proxy for measuring the achievement of the safety requirements. There is a demonstrable correlation between the collected evidence, measurement targets of safety-related properties and risk acceptance criteria associated with the safety requirements allocated to the AI system.
- b) The approach to measuring the achievement of the target values of the AI safety requirements is appropriate. This includes assurance arguments for the applicability of methods used and how representative and indicative the datasets are that are used to collect evidence. In particular:
  - 1) The datasets (e.g. test inputs) are representative of the input space.
  - 2) The datasets used to collect evidence are sufficient to detect critical classes of AI errors in the AI system, e.g. by covering known edge cases and triggering conditions.
  - 3) The datasets used to collect evidence are representative of the actual AI error rate for all inputs satisfying the system assumption, including in the presence of unknown triggering conditions. This includes an assessment of the statistical confidence of performance evaluations and overall coverage of the input space.

EXAMPLE 2 A method for obtaining a reliable target measurement for computer vision classification tasks based on a single image might not apply to object detection tasks involving the processing of real-time video streams.

## 8.7 Evaluation of the assurance argument

Confirmation measures according to 7.3.4 as well as methods and criteria for evaluating SOTIF according to ISO 21448:2022, 12.3 can be used to evaluate the achievement of AI safety on the basis of the assurance argument and associated evidence. The confirmation of the assurance argument for the AI system can be used as a precondition for the recommendation for SOTIF release at the level of the encompassing system (see ISO 21448:2022, 12.4).

In case of "conditional acceptance", the conditions required for a final release of the system can be documented in the assurance argument (see ISO 21448:2022, 12.4).

EXAMPLE 1 Restricted usage within the operational environment can be required to confirm assumptions regarding the distribution of triggering conditions. Once sufficient evidence has been gathered to support these assumptions, a final release can be accepted.

Sources of potential uncertainty in the assurance argument can be used to structure the evaluation procedure and to identify potential for strengthening the argument. This includes the identification of defeaters which might contradict assertions within the argument<sup>[25]</sup>. The following types of assertions can be identified for scrutiny within the assurance argument:

- Asserted context: These assertions are associated with the contextual information and assumptions that are used to scope the claims within the argument. If these assertions cannot be demonstrated to be valid, then the conditions under which the assurance argument is valid will be restricted.

EXAMPLE 2 Changes in the operational environment in which the AI system is deployed can undermine the assumptions made on the input space of the AI system, thus undermining the validity of the assurance argument.

EXAMPLE 3 An incomplete documentation or understanding of safety requirements allocated to the AI system will undermine the validity of the statements made within the assurance argument.

- Asserted evidence: These assertions relate to the evidence used to support claims in the assurance argument being both appropriate and trustworthy. This relates to individual pieces of evidence as well as combinations of evidence that are used to support a specific claim. If these assertions related to evidence cannot be argued with sufficient confidence, then the veracity of the claim can no longer be asserted.

EXAMPLE 4 Tests based on samples taken from within a test dataset are used to provide evidence for the robustness of the AI system against rare events (edge cases). However, there are not enough data points representing such events, which results in test results that are insufficient to demonstrate that the robustness claim has been fulfilled within a given statistical confidence interval. Therefore, additional or alternative forms of evidence are required.

EXAMPLE 5 Synthetic test data can be used to generate a sufficiently large and diverse number of edge cases that are not commonly found in samples taken directly from the operating environment. Additional assurance arguments can ensure the appropriateness of the testing approach to support the claim based on a validation of the fidelity of the synthetically generated data in comparison to the target environment and the inclusion of real-world samples in the test set.

EXAMPLE 6 When tools are used to produce evidence, the level of confidence in the evidence is directly linked to the confidence in the usage of the software tools. This is achieved by applying the requirements outlined in [Clause 15](#). Work products from [Clause 15](#) can be used in the assertion of the validity and integrity of evidence in the assurance argument.

- Asserted inference: These assertions relate to the reasoning behind the structuring of the assurance argument itself. In particular, how top-level claims are iteratively refined into detailed sub-claims that can be directly supported by evidence.

EXAMPLE 7 The set of causes of AI errors in the AI system used to structure a risk-based assurance argument overlook critical exacerbating factors (e.g. variation in sensor positioning and calibration) resulting in an assurance argument that demonstrates a set of properties that are not sufficient to ensure all safety requirements allocated to the AI system are met.

NOTE The use of assurance claim points [\[22\]](#), [\[25\]](#) can be used to elaborate those assertions within a GSN argument that require additional confidence arguments.

## 8.8 Work products

8.8.1 **Safety assurance argument**, resulting from [8.3.1](#) and [8.3.2](#).

8.8.2 **Confirmation measure reports**, resulting from [8.3.3](#).

## 9 Derivation of AI safety requirements

### 9.1 Objectives

The objectives of this clause are:

- a) to specify a complete and consistent set of AI safety requirements that are sufficient to ensure AI safety;
- b) to refine AI safety requirements based on learnings from development, verification and validation;
- c) to specify the limitations of an AI system over its input space to be escalated to its encompassing system development process.

## 9.2 Prerequisites and supporting information

The following information shall be available at the initiation of this activity:

- a) AI system definition (from external sources, e.g. the encompassing system development), including:
  - 1) safety requirements allocated to the AI system;
  - 2) input space definition;
  - 3) functional requirements;
  - 4) impacted stakeholders;
  - 5) the interfaces of the AI system with the encompassing system, including if applicable, the ASIL capability of the inputs to the AI system;
  - 6) interfaces to the environment, if applicable.

NOTE Safety requirements allocated to the AI system are allocated from the encompassing system development process. They can be motivated by different aspects, e.g. functional safety, SOTIF, or indirectly due to security (e.g. robustness against adversarial attacks and data poisoning). In this document, safety requirements are not explicitly distinguished by these aspects.

The following information can be considered during further iterations of this activity:

- b) safety analysis report, from [Clause 13](#);
- c) evaluation report of functional insufficiencies detected during operation, from [Clause 14](#);
- d) specification of the necessary off-board and on-board measures ensuring AI safety during operation, from [Clause 14](#).

## 9.3 General requirements

**9.3.1** The input space definition of the AI system shall be refined to the degree suitable for initiating the AI safety lifecycle.

**9.3.2** To provide a connection between each AI safety requirement and the addressed problem, the refined AI safety requirements shall either:

- a) trace to the safety requirements allocated to the AI system (from external sources), assumptions or critical scenarios; or
- b) address and trace to the potential influencing factors or root causes of functional insufficiencies and triggering conditions.

**9.3.3** A justification shall be provided that the refined AI safety requirements are reasonable to either ensure the achievement of the safety requirements allocated to the AI system (from external sources) or prevent or mitigate the functional insufficiencies at the AI system level.

NOTE Refined AI safety requirements to address functional insufficiencies at the AI system level are identified by safety analysis as necessary to fulfil the safety requirements allocated to the AI system (from external sources).

**9.3.4** To argue for the absence of unreasonable risk due to random hardware faults and systematic faults, the requirements of the ISO 26262 series shall be fulfilled.

NOTE 1 Combining [9.3.1](#) to [9.3.4](#), all of the causes defined in [Figure 6-12](#) can be addressed, i.e. functional insufficiency ([9.3.1](#) to [9.3.3](#)) and contributing systematic faults and contributing random hardware faults ([9.3.4](#)).

NOTE 2 Some adaptations can be necessary, in particular for safety requirements motivated by ISO 21448 activities with no ASIL rating and since the target is to achieve AI safety and not only functional safety.

**9.3.5** The following cases shall be identified and reported to the encompassing system development process:

- a) the AI system does not fully conform to the AI safety requirements;
- b) the AI safety requirements are only fulfilled for a limited part of the input space.

**9.3.6** AI safety requirements shall be identified to support the measures ensuring AI safety during operation.

EXAMPLE Different goals of field observation can be addressed with requirement [9.3.6](#):

- a) monitoring of the uncertainty in the current situation to indicate this to the encompassing systems allowing for measures to prevent hazardous behaviour at the vehicle level;
- b) monitoring of performance and detection of failure to support mitigation measures internal to the AI system (also referring to [10.5](#));
- c) supporting the continuous improvement of the AI system to ensure AI safety (e.g. recording devices to identify inconsistency among different sensor modalities and collect data for training and updating AI models);

## 9.4 General workflow for deriving safety requirements

[Figure 9-1](#) explains the general requirements in [9.3](#) for deriving AI safety requirements and establishes their connections to the objectives.

- The safety requirements allocated to the AI system (from external sources) are part of the AI system definition provided by the encompassing system development process. These safety requirements have SOTIF and functional safety aspects and are refined into an initial set of AI safety requirements. SOTIF requirements are typically identified as allowed maximum error rates while being exposed to the input space.

NOTE 1 Safety requirements allocated to the AI system (from external sources) are not work products in this document. "AI safety requirements" refers to all requirements derived within the scope of this document.

- Refined AI safety requirements (quantitative or qualitative) will be derived either by referencing requirements from past products, or by utilising the safety-related properties of AI systems that can be relevant to the application. These requirements will control uncertainty in the development process of the AI system in order to achieve the development quality of AI models and in addressing the safety-related issues in the AI system in order to achieve SOTIF. As the work product [9.6.2](#) AI safety requirements, these requirements are distributed to further development tasks in different phases of the AI safety lifecycle as described in other clauses (see [9.5.1](#) and [9.5.2](#)).

NOTE 2 Organizations define criteria to evaluate the uncertainty in the AI development process, i.e. rigour in training, evaluation, data collection, labelling, etc., to achieve the integrity of AI models, the safety-related errors in the AI system, and their overall impact to the encompassing system's safety. These criteria capture the organizational acceptance of uncertainty in the AI development process and safety-related errors in the AI system. Thus, it guides requirements elicitation, development tasks and decisions.

NOTE 3 Derived AI safety requirements include SOTIF and functional safety aspects and are allocated to AI components (see [Clause 10](#)), AI systems, encompassing systems, systems, vehicles, mobility services, etc. SOTIF requirements do not have ASIL values. Only functional safety requirements have ASIL values, including the level "QM", (quality management) and conform to the applicable parts of the ISO 26262 series.

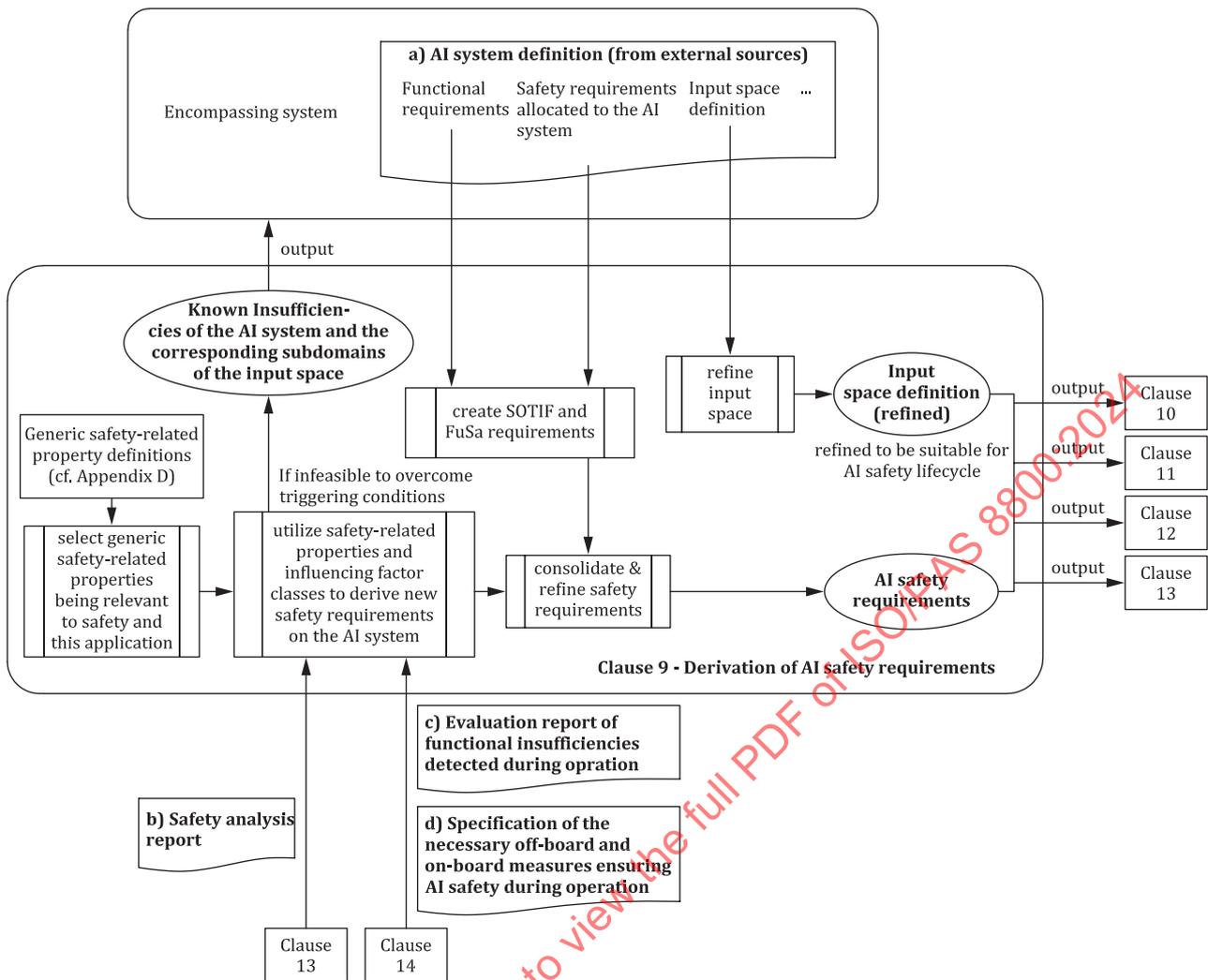
NOTE 4 Derived AI safety requirements can be allocated to the AI system or the AI system and further to AI components and tested at the respective level. In particular, the requirements allocated to the AI components, with appropriate safety metrics and test targets, are derived from the requirements allocated to the AI system. The derivation of the requirements allocated to AI components considers the complexity of the involved AI models, their task, and the environment they operate in.

EXAMPLE 1 The AI safety requirements allocated to an AI system for automatic braking, including test targets, can be derived from vehicle behaviour related to speed, distance to an obstacle, etc. If the AI system is composed of multiple AI components, the AI system decomposition needs to be accounted for in the derivation of the requirements allocated to the AI components. Examples of such decompositions are parallel AI model ensembles, parallel heads in multi-task architectures and sequential components in multi-stage object detectors.

- The input space definition, e.g. ODD in the automated driving context, in the prerequisite [9.2](#), list item a) AI system definition (from the encompassing system), is provided by the encompassing system development process. The definition can require further refinement to be used in AI safety lifecycle, as described in [9.5.3](#) and distributed as the work product [9.6.1](#) input space definition (refined).
- For functional safety requirements with ASILs (including QM), conformity to the ISO 26262 series can be demonstrated. For SOTIF requirements, the occurrence of particular error (sequence) patterns might be too high with respect to the criteria used to evaluate the uncertainty of the AI system during the engineering of the AI system or during operation. Refined AI safety requirements can be derived to inhibit the error (sequence) patterns produced by the AI system by understanding the triggering conditions or the safety-related errors. When further performance improvements are unlikely, performance limitations and relevant triggering conditions can be reported to the encompassing system development process. This is described in [9.5.4](#).
- Safety analysis ([Clause 13](#)) or observations from the field ([Clause 14](#)) can be used to determine the required thresholds for particular error (sequence) patterns.
- The activities described in [Clause 13](#) either evaluate the residual risk of the AI system with respect to the AI safety requirements or identify the safety-related errors in the AI system which can lead to the violation of the safety requirements allocated to the AI system (from external sources). The work products resulting from [Clause 13](#) activities are used as an input for a subsequent iteration of [Clause 9](#) activities.

EXAMPLE 2 A lane detection function produces AI errors at night and during heavy snow. Continued operation under these conditions can lead to the violation of AI safety requirements. This information can be communicated to the encompassing system development, so that measures can be taken at the system or vehicle level to restrict the conditions of operation or otherwise mitigate against AI errors under these conditions.

- During the operation phase, field monitoring can be used to detect unknown triggering conditions and violations of assumptions.



**Key**

- prerequisite work products
- ▭ activities
- ▭ information from/ to external or other clauses

NOTE For ease of understanding, the following activities are omitted in the diagram:

- a) reviewing the selection of safety-related properties of AI systems based on safety analysis in an iterative manner;
- b) refinement of input space definition based upon e.g. safety-related errors or component-level triggering conditions;
- c) reporting to the encompassing system development process for the technical limitations of the AI system to ensure the AI safety requirements and the safety requirements on the encompassing systems and the vehicle reach consistency.

**Figure 9-1 — Conceptual diagram reflecting major activities in derivation of AI safety requirements**

## 9.5 Deriving AI safety requirements on supervised machine learning

### 9.5.1 The need for refined AI safety requirements

Based on the ISO 21448 (SOTIF) activities at the vehicle level to identify and evaluate risks, potential functional insufficiencies and potential triggering conditions, safety requirements allocated to the AI system (from external sources) are refined into AI safety requirements.

For AI systems operating within a specific input space, AI safety requirements can contain acceptance criteria similar to the concept of probability of failure on demand (PFD), e.g. “Probability (occurrence of an error pattern)  $< \alpha$ ” (see IEC 61508-4:2010).

**EXAMPLE 1** An AI error (sequence) pattern of a DNN “Consecutive misdetection (False Negatives, FNs) of a nearby pedestrian for more than 0,1 seconds” can be refined into “The occurrence of more than three consecutive FNs of a pedestrian @30FPS” if the DNN is triggered at a rate of 30 FPS. Such an AI error (sequence) pattern with FNs can be verifiable only if the ground truth is available.

Methods such as systems-theoretic process analysis (STPA) can be used to refine AI safety requirements where a detailed design is available (see [Annex E](#)), i.e. low-level control structure in STPA terms. However, specific types of AI models (e.g. those implemented by supervised learning) cannot be decomposed further. This clause focuses on the case that the AI model cannot be decomposed.

**EXAMPLE 2** An example detailed design in an AI system can be:

- redundant DNNs for perception task;
- a majority voting on the individual results;
- an out-of-distribution (OOD) detector to abort the validity of the majority voting upon detection of an OOD input.

Based on a detailed design, refined AI safety requirements can be derived to address unsafe conditions, e.g. OOD detector provides an incorrect output in rainy situations.

For ML applications, the probability of an error in an AI system may not be computable due to the complexity of its input space, e.g. pedestrian detection for a wide variety of pedestrians in automated driving, and the inherent nature of ML, i.e. in-distribution and OOD error gap. The relative frequency of an error in an AI system often can be estimated only if enough samples are observed in its input space. The inaccuracy of relative frequency, e.g. PFD, is the uncertainty of an AI system that refined AI safety requirements can control in the AI development process. Requirements on a sufficiently low level of probability of AI errors are refined into quantitative requirements and qualitative requirements. [Figure 9-2](#) illustrates the underlying approach.

**EXAMPLE 3** A safety requirement allocated to an AI system (from external sources) “Probability (occurrence of two consecutive FNs of a nearby pedestrian per hour of operation)  $< 10^{-6}$ ” can be refined into an AI safety requirement “The DNN does not produce two consecutive FNs of a nearby pedestrian in  $10^8$  hours of driving data with sufficient diversity demonstrated” under the condition the validity of the refinement is provided. It is assumed that more driving hours can be used to prove the hypothesis that the requirement holds.

**NOTE 1** Relative frequency, empirical probability or experimental probability of an event is the ratio of the number of outcomes in which a specified event occurs in the total number of trials<sup>[26]</sup>.

**NOTE 2** For estimating very low probability using relative frequency, the denominator, i.e. the number of samples to be tested, can be large. Sample size determination for estimating the probability of occurrence of a particular error pattern can be based on estimation principles and the statistical confidence of experimental results.

Requirements on probability of error occurrence being low



Requirements on relative frequency of error occurrence being low (by sampling and counting)

and

Requirements on additional engineering means to control or reduce the uncertainty in evaluating quantitative requirements

Quantitative requirements

Qualitative requirements

**Figure 9-2 — Understanding refinement of AI safety requirements for supervised learning**

9.5.2 Derivation of refined AI safety requirements to manage uncertainty

For engineering machine learning components, risks can be reduced by managing sources of uncertainty in the AI development process. Uncertainties exist for each phase of the AI development process, which are referred to here as potential influencing factors. Influencing factors fall into the following categories: observation, label, model, and operation. [Table 9-1](#) describes typical approaches to managing influencing factors in the AI development process, which constructs an abstracted framework to reduce the uncertainty regarding limited knowledge about the true state of the world.

NOTE 1 Such limited knowledge is expected in any data-driven AI development. Perceptual uncertainty refers to uncertainty associated with the performance of perceptual components.<sup>[27]</sup>

**Table 9-1 — Managing influencing factors in the AI development process**

Influencing factor class	Uncertainty management approach
Observation	(1) Enumerate patterns influencing AI outputs both in deductive (top-down) and in inductive (bottom-up) ways (2) Define data coverage policy for comprehensive AI performance and safety (harm avoidance) trade-off
Label	Reduce variation of labelling policy and enhance/ensure labelling work quality
Model	Use appropriate model selection and de facto standard approaches for tuning and evaluation of hyperparameters
Operation	Detect deviation in model performance between development and deployment

Addressing the above influencing factors, [Table 9-2](#) describes a non-exhaustive set of example generic AI safety requirements that may be considered for specific applications.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

Table 9-2 — Example generic AI safety requirements to manage influencing factors

Influencing factor class	Example generic AI safety requirements
Observation	<ul style="list-style-type: none"> <li>— Derive scenes relevant to attributes of inference targets, environmental conditions, and system configurations (relevance)</li> <li>— Include scenes which would potentially lead to errors (e.g. FN, FP, substantially incorrect regression)</li> <li>— Consider the effects of system configuration change during operation (e.g. sensor mounting position shift, sensor hardware degradation caused by ageing)</li> <li>— Include specific scenes which empirically have led to errors when developing the current product and when developing and operating legacy products (safety)</li> <li>— Specify data quantity for each scene or combination of scenes (diversity)</li> <li>— Explain the rationale for the scenes for which one cannot define target data quantity</li> </ul>
Label	<ul style="list-style-type: none"> <li>— Specify inference targets</li> <li>— Describe the necessary quality of labels, e.g. FP rate, bounding box size precision, etc.</li> <li>— Specify labelling procedure                             <ul style="list-style-type: none"> <li>— Define the appropriate type of label, e.g. class, numerical.</li> <li>— Publish and circulate a labelling policy guideline, e.g. treatment of occluded objects, number of annotators annotating the same data etc., to all relevant stakeholders</li> <li>— Provide clear criteria and lines of accountability about the labelling of data involving protected characteristics or special category data (or both)</li> </ul> </li> <li>— Specify label evaluation policy, i.e. evaluation timing, procedure and measures, and acceptance criteria                             <ul style="list-style-type: none"> <li>— Make sure appropriate processes are in place if using crowdsourced labellers</li> <li>— Assess the risk of incorrect labelling through sample checks of submitted labels</li> </ul> </li> <li>— Specify label incident reporting and management process to handle labelling errors and ambiguity in labelling policy</li> </ul>
Model	<ul style="list-style-type: none"> <li>— Specify metrics and acceptance criteria for model performance</li> <li>— Specify model selection policy and ensure an appropriate model is selected according to the policy</li> <li>— Specify hyperparameter search policy and conditions, e.g. bounds, and document the determined hyperparameter values</li> <li>— Perform safety analysis of the AI model</li> <li>— Introduce the technical capability of checking the plausibility of the output</li> <li>— Introduce the technical capability of de-noising/quantizing the input</li> </ul>

Table 9-2 (continued)

Influencing factor class	Example generic AI safety requirements
	<ul style="list-style-type: none"> <li>— Introduce the technical capability of robust training techniques (e.g. training against noise, quantized neural networks)</li> <li>— Introduce the technical capability of providing interpretation (explainability) over the decision being made                             <ul style="list-style-type: none"> <li>— For DNN, pre-hoc global explainability (i.e. converting decision-making mechanisms into human-comprehensible rules) can be difficult. Post-hoc local explainability (i.e. the reason why a particular decision is made for a specific input) is likely, via techniques such as local linearization or heatmaps [28]</li> </ul> </li> <li>— Introduce the technical capability to measure epistemic uncertainty for the AI model outputs</li> </ul>
Operation	<ul style="list-style-type: none"> <li>— Introduce the technical capability of identifying situations matching a known triggering condition or measuring model uncertainty</li> <li>— Introduce the technical capability of identifying concept/domain drift</li> <li>— Introduce the technical capability of identifying implausible or otherwise non-trustworthy AI system behaviour</li> </ul>

While AI safety requirements derived using influencing factors are largely qualitative, how safety-related properties of AI systems can be used to assist in deriving refined quantitative AI safety requirements is considered in the following.

In this document, we establish a differentiation between safety-related properties and AI safety requirements.

- Safety-related properties are inherent characteristics of engineering AI systems which may either lead to the insufficiency of the generalization or merely make it difficult to argue the safety of the AI system. Safety-related properties are a subset of generic AI properties; they are application-independent and may include aspects such as robustness or domain shift, incomplete specification as suggested by ISO/IEC 22989 and related guidance on safety in ISO/IEC TR 5469. See [Annex D](#) for a list of AI properties that may be safety-related.
- AI safety requirements are application-dependent and can be mapped to one or more safety-related properties.

NOTE 2 The term "safety-related property" in this document is stated differently in other contexts with similar meanings. In the French national project DEEL, it is referred to as high-level properties[29], while in the German national project KI-Absicherung (EN: AI-assurance), it is referred to as safety concerns[30].

The list of safety-related properties of AI systems from [Annex D](#) can be used to refine AI safety requirements. The safety-related properties of AI systems considered relevant for the AI system under consideration are identified and instantiated into refined AI safety requirements. These requirements can be associated with the safety-related properties of AI systems to support the selection of AI technologies and architectural and development measures, as detailed in [Clause 10](#).

Deductive safety analysis approaches such as FTA might not directly lead to the identification of safety-related properties and associated AI safety requirements. This is because a causal analysis of AI errors might have limited effectiveness due to the complexity of the model and interaction of numerous, potentially unknown causes. Instead, inductive approaches such as hypothesis testing are used to derive the safety requirements. First, an initial set of AI safety requirements is identified, then AI safety requirements are updated and validated through safety analysis.

NOTE 3 These safety-related properties of AI systems are currently described conceptually rather than using formal definitions in mathematics. The list of safety-related properties on AI systems, e.g. [Annex D](#), is not exhaustive. The list is based on discussions in AI safety research, etc.

Refined AI safety requirements can be either qualitative or quantitative, where for quantitative AI safety requirements, the design of the acceptance thresholds is application dependent. The validity of these thresholds requires justification with respect to the evaluation criteria and the overall residual risk acceptable for the AI system. For the example of robustness, the validation means “Can setting such a threshold positively increase the robustness?”, which is an activity to be evaluated via experiments in [12.5.7](#).

**EXAMPLE** Appropriate safety-related properties are hypothesized at the early stage of development, e.g. based on past product experiences and identified by safety analysis in an iterative manner to ensure their contribution to system safety. Safety analysis can include the impact of changing safety-related properties, specific noise robustness for the equipped hardware characteristics, etc., on the system’s safety. For example, a safety-related property of AI systems AI robustness is identified by such safety analysis and further refined to different kinds of robustness, then the following quantitative AI safety requirement can be derived along with justification for the validity of these thresholds: “For all clear images in the input space, if noise perturbations characterized by  $L_1$  norm  $< 0,001$  are added on the image, the AI system shall at most introduce 0,01 % of new errors”. This requirement, derived from “AI robustness,” will lead to additional engineering efforts in [Clause 10](#), such as using robust training techniques.

Using safety-related properties of AI systems to derive refined AI requirements is an inductive approach, and exhaustiveness is not ensured. Therefore, it can be only a complementary part of the verification and validation strategy elaborated in ISO 21448:2022, Clause 9.

### 9.5.3 Refinement of the input space definition for AI safety lifecycle

The input space defines what workspace, what conditions and around what dynamic elements a system will operate to ensure its safe design and operation. If a system depends on AI systems using data-driven methods to build ML models, the input space leads to the definition of the dataset requirements, which form the basis of the data used throughout the lifecycle of the system.

**EXAMPLE** For autonomous driving, the workspace can be the road or area where the system operates. The conditions can be the weather, visibility, illumination and connectivity. The dynamic elements can be traffic and moving debris.

The refinement of the input space is an iterative process. The initial definition provided as a prerequisite may not be suitable to be used for ensuring the intended functionality of the AI system. This is also reflected in the ISO 21448 document, emphasizing that environmental factors are essential issues, while systems and their elements have different concerns depending on the hierarchical layers.

The refinement of the input space definition aims at conforming to the capabilities of the available systems (e.g. sensors). Its refinement may also consider newly discovered or known triggering conditions that may lead to performance insufficiencies. Regarding where and what kind of refinement of the definition of the input space is needed, the information can be derived from the results of the safety analyses as described in [Clause 13](#).

The refined input space definition forms the basis for the verification and testing of the target system performance. That is, the refined input space bounds the performance expectations of the AI system, reducing the uncertainty associated with its operation and consequently ensuring the bounds of its safe operation.

### 9.5.4 Restricting the occurrence of AI output insufficiencies

Overall, when engineering an AI system, a particular scenario (or a set of scenarios) where the error rate of the AI system is too high may be observed, violating the initially derived AI safety requirements. The predefined list of safety-related properties of AI systems can be used again as a checklist to examine the potential root causes and subsequently, derive a new set of AI safety requirements with associated metrics. In the iterations beyond the initial iteration, the safety analysis described in [Clause 13](#) provides additional input to this activity. Based on observations from the field or from V&V activities, the safety analyses identify the safety-related properties which maximally correlate with the safety-related issue. The insights into the

correlation (or if applicable causality) between output insufficiencies and safety-related issues can be used as a basis for the refinement of the AI safety requirements after the initial iteration.

NOTE 1 Currently these safety-related properties are described conceptually rather than using formal definitions. The consequence is that the derivation of AI safety requirements is not viewed as a (logical) causal derivation but rather a hypothesis to be validated.

NOTE 2 The vast majority of ML-type statistical models estimate the likelihood of events, which is inferred from correlations in the data. In such cases, clear causal relationships between the inputs and outputs of these models (and ultimately the actions taken by the system) cannot be established. Thus, the classic causal fault analysis tools are not necessarily applicable in this context.

In addition to the analysis results from [Clause 13](#), which are driven by V&V issues or issues observed in the field, analytic methods can also support the creation and refinement of AI safety requirements, which may include plausibility checks on the inputs to the AI system and bounding on the outputs from the AI system.

As illustrated in [Figure 9-3](#), an issue associated with performance insufficiencies may be identified from a too-high error rate during operation. By conducting appropriately designed experiments (e.g. counterfactual analysis or hypothesis testing), it may be discovered that the issue is strongly correlated with one or more intrinsic AI properties that are safety-related (e.g. robustness), with evidence of correlation empirically manifested in some metrics. The activity can be continued to further refine the influencing factors as described in [Table 9-3](#), thereby deriving reasonable AI safety requirements with the aim to prevent or mitigate performance insufficiencies. The effectiveness of the refined AI safety requirements as countermeasures should still be validated (see [12.5.7](#)).

NOTE 3 The use of correlation-driven analysis techniques is driven by the practical need as demonstrated in the field of supervised ML. It does not inhibit the establishment of an assurance argument for AI, provided that additional supporting evidence such as experiments can be provided. This document also does not inhibit logical causal analysis for AI as used in establishing the safety case, provided that causality can be rigorously demonstrated.

NOTE 4 Regarding systematic errors, technical safety requirements allocated to the AI system are achieved by safety mechanisms as the countermeasures for systematic errors.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

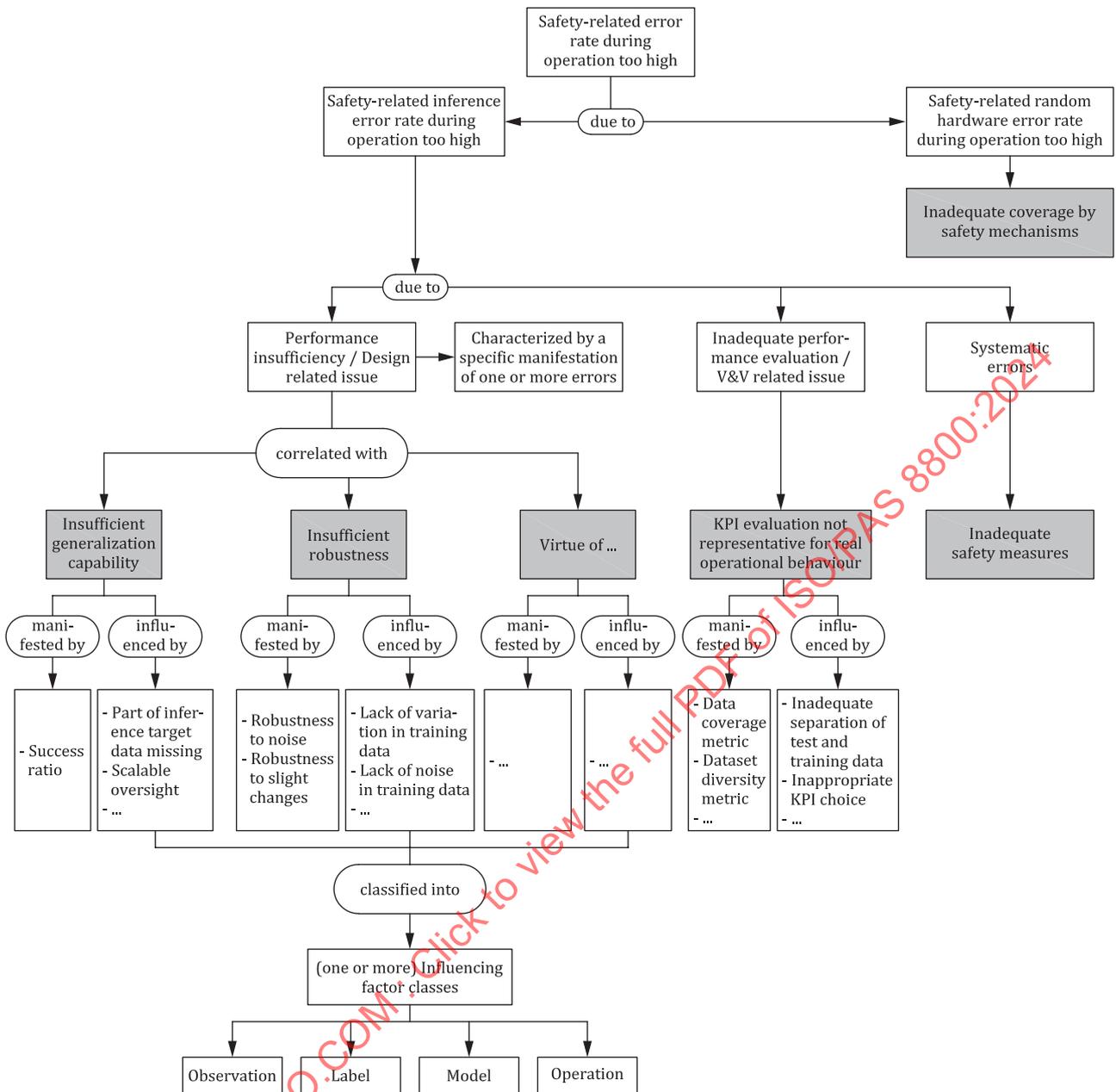


Figure 9-3 — Conceptual diagram illustrating safety-related properties

Table 9-3 — Influencing factor classes

Influencing factor class	Description	Examples
Observation	Influencing factors related to data representing the input space	<ul style="list-style-type: none"> <li>— Diversity of training and test data</li> <li>— Coverage of the inference target data domain</li> <li>— Distribution according to features of the inference target domain</li> </ul>
Label	Influencing factors related to the data labels	<ul style="list-style-type: none"> <li>— Incorrect labels</li> <li>— Inconsistent labels</li> </ul>
Model	Influencing factors related to the ML model itself	<ul style="list-style-type: none"> <li>— Choice of ML models</li> <li>— Choice of hyperparameters</li> <li>— Training procedure</li> <li>— Limited computing power/timing/memory, or precision change/constraint during deploying AI systems into target environment, e.g. hardware</li> </ul> <p>NOTE The impact of differences due to deployment constraints between the ML model developed in the off-board environment and its on-board implementation cannot be analytically determined due to the complexity of ML models and is left as uncertainty in the AI system.</p>
Operation	Influencing factors related to changes within the input space during operation	<ul style="list-style-type: none"> <li>— For ML models classifying traffic signs: Introduction of new traffic signs during its operational phase</li> </ul>

Table 9-4 illustrates an example of connecting insufficiencies, safety-related properties and concrete AI safety requirements.

Table 9-4 — Examples for utilizing safety-related properties to derive new safety requirements

Insufficiencies observed from the field	Associated safety-related properties	Concrete AI safety requirement derived from insufficiencies	Metrics as KPIs	Original error pattern being considered
Under slight input variations, the object being detected is missing (FN)	Robustness (inherent in model training, decision boundary)	For any training data, if a Gaussian noise bounded by L-infinity norm $< 0,01$ is added, the predicted output class shall remain the same with a high probability (e.g. 99,5 %)	For each image in the training/test data, evaluate the robustness by creating 1 000 variations of each image with additional Gaussian noise, and derive the number of incorrect predictions KPI returns “violation / false” if there is an image where the number of incorrect predictions subject to noise is larger than 5	The safety requirements allocated to the AI system (from external sources) related to FN
When hardware faults occur where a small area of pixels has turned black, the prediction can be constantly wrong (FN)	Robustness (against random HW faults)	For any training data, if 1 % of the pixels is arbitrarily changed to completely black, the predicted output class shall remain the same with a high probability (e.g. 99,5 %)	For each image in the training/test data, evaluate the robustness by randomly setting 1 % of the pixels to black, repeat 1 000 times, and derive the number of incorrect predictions KPI returns “violation / false” if there is an image where the number of incorrect prediction subject to noise is larger than 5	The safety requirements allocated to the AI system (from external sources) related to FN

NOTE L-infinity norm quantifies noise by taking the maximum absolute value of the noise on every input dimension. For example, if there are three inputs to the DNN and the noise on each input is 0,1, -0,2, and 0,15 respectively, then L-infinity norm =  $\max \{ |0,1|, |-0,2|, |0,15| \} = 0,2$ . The use of L-infinity norm is only an option; there are other norms (e.g. L1 or L2) for quantifying the noise.

9.5.5 Metrics, measurements and threshold design

Quantitative evidence requires the definition of a set of metrics and target thresholds. The metrics are used to define, characterise and theoretically analyse specific properties of AI systems. The selection and design of these metrics and associated measurements and target thresholds involves the following considerations:

- metrics and measurement methods for specific properties of AI technologies referenced in the safety requirements allocated to the AI system (from external sources);
- measurements and evaluation of the AI systems with respect to these properties;
- analysis of how the AI methods used to develop the AI system impact the metrics/measurement methods chosen for the AI safety requirements;
- definition and creation of suitable datasets used to measure each property (see [Clause 11](#)). The datasets are designed to provide a statistically relevant analysis of the property and be representative of the input space.

Selection of suitable target thresholds for each metric requires justification. The justification can include:

- past product experiences
- commonly agreed industry consensus
- system analysis
- expert judgements

- experiments

The decision may also involve multiple factors where trade-offs (e.g. bias-variance trade-off and robustness-accuracy trade-off phenomena in machine learning) are taken into consideration.

**EXAMPLE** Robustness against noise can be demonstrated by the minimum perturbation needed to change the output to an incorrect prediction. Different minimum perturbation values can be selected as the acceptance criterion such as  $\varepsilon_1$  (e.g. 0,5) rather than  $\varepsilon_2$  (e.g. 0,25). There can be different ways of justifying the threshold by providing convincing arguments, e.g. the robustness shall be larger than  $\varepsilon_1$ , since:

- $\varepsilon_1$  is the maximum possible noise communicated from the complete input pipeline, or
- $\varepsilon_1$  is derived from experiments where experts agree that it is hard for humans to make reasonable predictions when noise exceeds  $\varepsilon_1$ , or
- $\varepsilon_1$  is the value that is used in similar standards or products.

### 9.5.6 Considerations for deriving safety requirements

While previous clauses provide general principles in deriving refined AI safety requirements via the assistance of safety-related properties of AI systems, the following is a summary of some practical example considerations for the derivation of AI safety requirements.

**NOTE** This list is not exhaustive. Nevertheless, it addresses complementary topics that can be relevant to safety.

- Given the lack of guarantees on the ability of ML models to generalize, it is important to specify the input space that defines the workspace, conditions and dynamic elements, thereby limiting the performance requirements and also the training data needed for that space.
- Results from statistical learning theories are commonly based on idealistic assumptions. These theories (e.g. Vapnik–Chervonenkis theory)<sup>[31]</sup> can nevertheless be used to derive a lower bound on the number of samples needed to ensure a tight generalization error for a given type of ML model (e.g. support vector machine). The derived number of samples for model training commonly reflects the best-case scenario where the data used in training (in-sample) has the same distribution as the data in operation (out-of-sample).
- The possibility of relevant foreseeable adversarial attacks, i.e. foreseeable attacks that are judged to be realistic, and their impact on the overall AI system can be considered. However, cybersecurity has not been considered in this document.
- The side effects of region-specific privacy considerations (e.g. the General Data Protection Regulation in the European Union) can be considered, since they can indirectly influence the quality of the data and impact the model's performance.
- A practical purpose of improving AI explainability is to ease the engineering of AI systems. For validation and performance improvement purposes, interfaces to AI components can be specific to improve the understanding of the AI component response. Such interfaces can allow for introspective approaches for validation and performance improvement, particularly isolating errors of the black-box ML based AI components.
- To reduce gaps or non-conformities to a particular AI safety requirement (e.g. by improving the related performance), different performance aspects, which are inherent to the selected AI methods, can become entangled. If overall efforts cannot be increased to meet either of the conflicting safety requirements, a trade-off (e.g. between robustness and accuracy) is found. As shown in [Figure 9-1](#), such limitations can be forwarded to the encompassing system development process so that safety requirements on the AI system and the encompassing system design are updated, and the entire safety requirements are met.
- Results from the Neyman-Pearson approach (hypothesis testing) can be used to minimize the missed detection (number of FNs) under a fixed number of false alarms (number of FPs) by considering inputs like sample size and signal-to-noise ratio (SNR).<sup>[32]</sup>

## 9.6 Work products

9.6.1 **Input space definition (refined)**, resulting from [9.3.1](#).

9.6.2 **AI safety requirements**, resulting from [9.3.2](#), [9.3.3](#), [9.3.4](#), and [9.3.6](#).

9.6.3 **Known insufficiencies of the AI system and the corresponding subdomains of the input space**, resulting from [9.3.5](#).

## 10 Selection of AI technologies, architectural and development measures

### 10.1 Objectives

The objectives of this clause are:

- a) to select and justify appropriate AI technologies for use in the AI system;
- b) to identify appropriate architectural and development measures to fulfil the safety requirements prior to deployment;
- c) to identify appropriate architectural measures to mitigate residual functional insufficiencies of the AI system revealed after deployment;
- d) to identify measures for ensuring the safety requirements of the AI system are fulfilled within its target execution environment.

### 10.2 Prerequisites

The following information shall be available:

- a) safety requirements on the AI system, [Clause 9](#);
- b) training and validation datasets, [Clause 11](#);
- c) AI component or AI system architecture, if available;
- d) AI component or AI system development process, if available.

### 10.3 General requirements

**10.3.1** A justification shall be provided that the selected AI technologies and AI methods are capable of fulfilling the AI safety requirements.

**NOTE** AI technology, its application to road vehicle functionality, as well as methods for ensuring that AI safety requirements are met are rapidly evolving. However, the most advanced technologies are not necessarily the most suitable for safety-related applications. This is due to the lack of an appropriate set of methods for ensuring the safety of such technologies (see Technology class III, which is outlined in ISO/IEC TR 5469).

**EXAMPLE** After analysing the benefits and limitations of alternate technologies, an argument is made to justify why DNNs are selected and used in combination with a set of redundancy and monitoring measures for a potentially safety-related functionality despite the challenges of demonstrating safety requirements for these approaches.

**10.3.2** AI safety requirements shall be allocated to AI components.

**NOTE 1** In some exceptional cases, it is possible that not all AI safety requirements are allocated to an AI component, e.g. diagnostics components.

**NOTE 2** The AI safety requirements allocated to the AI component depend on the functionality of the AI component and on the AI safety requirements of the system.

**10.3.3** Sufficient measures, such as architectural, development or a combination of both, shall be defined to ensure the AI safety requirements are fulfilled by the AI components.

NOTE Architectural and development measures help prevent, by design, AI errors of the AI components.

**10.3.4** Sufficient measures, such as architectural, development or a combination of both, shall be defined to reduce the risk resulting from contributing AI errors of the AI components.

EXAMPLE For out-of-distribution error detection, implementation of reject classes implies development measures as well as architectural measures for ML systems.

**10.3.5** The effectiveness of the chosen combination of architectural and development measures resulting from [10.3.3](#) and [10.3.4](#) shall be supported by an argument.

**10.3.6** Safety analysis of the AI system outputs and, where reasonably practicable, of its architectural elements shall be performed to determine whether the safety requirements allocated to the AI system can be met.

NOTE For DNNs, safety analysis of the architectural entities of the AI model are not always reasonably practicable.

EXAMPLE Safety analysis can include the analysis of the computational graph of AI components to identify if the intermediate (latent space) or final outputs meet the relevant AI safety requirement allocated to the AI components.

**10.3.7** The differences between the development environment and the target execution environment shall be identified and evaluated regarding their potential impact on the safety requirements. If necessary, appropriate AI architectural and development measures shall be defined.

**10.3.8** AI components that are AI models or contain AI models shall be trained using the training dataset and evaluated using the validation dataset.

## 10.4 Architecture and development process design or refinement

The architecture of an AI component or AI system is designed or updated, to meet its AI safety requirements. Similarly, an AI system or AI component development process is tailored or designed according to the AI safety requirements.

AI safety requirements are satisfied by two types of measures, architectural and development. There are two categories of AI safety requirements, some derived from safety requirements allocated to the AI system (from external sources) and some derived from safety related properties of AI system. For the latter category, [Table 10-1](#) provides some examples of measures that can help meet them.

It is possible that, during the design process, there is no architecture and/or development process that can meet all AI safety requirements. The challenges can then be discussed with the requirement stakeholders and the AI safety requirements updated (see [Clause 9](#) for guidance on updating AI safety requirements).

Similarly, there can be more than one architecture and/or development process that meets all AI safety requirements during the design process. The benefits and the cost of each option will then be discussed with the system integrator to choose the most appropriate option for the AI component or AI system application.

Once a candidate architecture and development process are identified, the AI component or AI system is trained using the training dataset and its potential to achieve its allocated safety requirements is evaluated using the validation dataset.

AI model training is a critical step in the development process where the model learns from data in an iterative manner by using a preferred optimization algorithm. In “supervised learning”, for example, the learning process involves learning the weights and biases that minimize the error between the model's prediction and the ground truth. The training process relies on several tuneable parameters known as hyperparameters (e.g. learning rate, regularization strength). Hyperparameter tuning helps optimize the model's performance. The training process can also involve one or more of the steps such as feature engineering, regularization, dropout, error analysis, etc.

## 10.5 Examples of architectural and development measures for AI systems

[Table 10-1](#) provides guidance on which measures can support the achievement of the AI-property-specific KPIs and targets associated with AI safety requirements and safety-related properties of AI systems.

NOTE 1 [Table D-1 \(Annex D\)](#) provides a definition of the safety-related properties of AI systems listed in [Table 10-1](#).

NOTE 2 [Annex G](#) provides a short description of each architectural and development measure listed in [Table 10-1](#).

NOTE 3 The applicability of a safety-related property of an AI system depends on the use case, encompassing systems AI models, etc. For example, while a self-driving vehicle's actions, such as acceleration and steering, can be controllable, the property of AI controllability might not apply to the outputs of a DNN model for object detection in the perception pipeline.

Acknowledging that there is an overlap between the safety-related properties of AI systems, rationales used for the allocation of measures to the appropriate AI property(ies) can include:

- a) AI resilience supports AI robustness. The following rationales were considered in the selection of recommended measures to distinguish AI robustness and AI resilience:
  - Measures that can guarantee that the system maintains its nominal performance under bounded input perturbations are classified under AI robustness;
  - Measures that mitigate a failure for which the system impact is unclear are classified under AI resilience;
  - Some measures fall under both categories, AI robustness and AI resilience.
- b) AI controllability can support AI resilience, e.g. a supervisory system detects an error and switches between redundant systems. The same supervisory system can also detect an error and trigger a risk mitigation measure that stops the AI system operation. The following rationales were considered in the selection of recommended measures to distinguish AI controllability and AI resilience:
  - The aim of AI resilience is to keep the system running.
  - The aim of AI controllability is to keep the system safe.
- c) AI alignment and AI predictability have the same objective of establishing confidence in the correctness of the AI system's prediction but achieve it with different means:
  - Measures to demonstrate that the AI system's behaviour is aligned with the user's expectations and values are classified under AI alignment.
  - Measures that provide supporting evidence that the system behaves as expected, i.e. the system's behaviour can be reasonably predicted based on its inputs, are classified under AI predictability.

Table 10-1 — Example of measures for AI robustness

Type of measure	Measure	Remarks
Architectural	Architectural measures that foster AI robustness against OOD inputs (G.3.2.1)	Architectural artefacts such as reject classes can be required to detect OOD inputs.
	Diverse redundant models (G.1.2.1) Model ensembles (G.1.2.2) N-version diverse programming (G.1.2.3) Selection techniques for architectural redundancy (voting and switching) (G.1.2.5)	Architectural redundancy combined with a voting system provides confidence in the AI system generating a correct output despite some ML elements providing wrong predictions.
	Fault-aware training (G.4.2)	Training the AI system to recognize faults helps to handle those faults whilst maintaining a nominal or degraded mode.
	Adversarial training (G.3.2.2, G.4.2)	Adversarial training reduces the AI system's sensitivity to external malevolent perturbations and increases its reliability.
Development	Transfer learning (G.4.3)	Transfer learning relies on the robustness a given AI model has shown within its source input domain to leverage this robustness within the target input domain.
	Augmentation of data (G.4.4)	Using data augmentation to increase the diversity of data the model is exposed to helps it generalize better.

Table 10-2 — Example of measures for AI generalization capability

Type of measure	Measure	Remarks
Architectural and development	Transfer learning (G.4.3)	Transfer learning can leverage the generalization capability of the foundation model to the target application.
Development	Regularization (G.4.2)	Regularization methods help the model adapt better to small shifts in the input domain and reject OOD inputs.
	Hyperparameter tuning (G.4.1)	Hyperparameter tuning can help reduce the underfitting or overfitting, thereby improving generalization.

Table 10-3 — Example of measures for AI reliability

Type of measure	Measure	Remarks
Architectural and development	All measures supporting AI robustness, AI resilience and AI generalization support reliability	None

Table 10-4 — Example of measures for AI resilience

Type of measure	Measure	Remarks
Architectural and development	OOD data and its mitigation (G.3.1) Distributional shifts and its mitigation (G.3.2)	Mechanisms that detect OOD samples and distributional shifts trigger the need to update and that contain actions to ensure the system's safety until the OOD or distributional shift is addressed.
	Qualitative and quantitative analysis of AI architectures (Clause G.2)	Safety analyses on the architecture help identify the need for redundant systems that maintain the AI system to an initial or degraded performance in case of AI errors.
	Usage of AI-model based and conventional software (G.1.2.6)	Non-AI components can be used to detect errors and switch to redundant or fallback systems. Redundant systems provide a means to continue the AI system operation with the initial performance.
Development	Targeted and controlled model update (G.5.2.2)	Partial and targeted updates can allow simplified degradation of safety performance testing and speed up an issue resolution.

Table 10-5 — Example of measures for AI controllability

Type of measure	Measure	Remarks
Architectural	Usage of AI components and non-AI components (G.1.2.6)	Non-AI components are used to detect errors and switch to redundant or fallback systems. Fallback systems guarantee the AI system's controllability in all situations.
	Supervisory, limiting logic and non-AI backup system (G.1.2.4)	The implementation of safety monitors provides a means to detect errors and take control over one element of the AI system to ensure the AI system's safety.
Architectural and development	Qualitative and quantitative analysis of AI architectures (Clause G.2)	Safety analyses on the architecture help identify the need for a monitoring system that can detect and control AI errors.

Table 10-6 — Example of measures for AI explainability

Type of measure	Measure	Remarks
Development	Attention or saliency maps (G.4.7)	Attention/saliency maps provide information that explains the input characteristics that strongly influence the prediction of the ML system.
	Structural coverage of AI component (G.4.9.1)	This is a white (or open) box method. It can build confidence by identifying which features of inputs are important for the decision/prediction of AI models.
	Identification of SW units (G.2.1)	Breaking down the architecture into SW units aims to understand the function and performance of each unit, fostering some explainability of the AI component output outcome.

Table 10-7 — Example of measures for AI predictability

Type of measure	Measure	Remarks
Architectural	Model ensembles (G.1.2.2) Techniques for selection of architectural redundancy (G.1.2.5)	Ensembles typically increase the accuracy of the prediction, which fosters predictability.
Architectural and development	Criteria for retraining (G.5.2.1)	Monitoring criteria for retraining, e.g. distributional shift, help identify a decrease in the performance of the AI model to predict the correct output. Consequently, this monitoring also helps to detect a reduction in robustness, resilience and generalisation capability.
Development	Confidence calibration and uncertainty quantification of AI models (G.4.4)	Calibrating the AI model's uncertainty fosters confidence in the correctness of the model's predictions.
	Structural coverage of AI component (G.4.9.1)	Structural coverage provides confidence in the comprehensiveness of the testing strategy.
	Monitoring multiple scores (G.4.6)	Monitoring model performance metrics such as precision, recall, F1-score during training provides insight into the model's ability to predict correct outputs consistently.

Table 10-8 — Example of measures for AI alignment

Type of measure	Measure	Remarks
Development	Alignment of intention (Clause G.6)	None

Table 10-9 — Example of measures for bias and fairness

Type of measure	Measure	Remarks
Development	Data coverage techniques for test data augmentation (G.4.9.2)	This method can help identify and reduce bias within the datasets by measuring data distribution across equivalence classes.

## 10.6 Work products

**10.6.1 AI component or AI system architecture** (refined), resulting from [10.3.1](#) to [10.3.7](#).

**10.6.2 AI component or AI system development process** (refined), resulting from [10.3.1](#) to [10.3.7](#).

**10.6.3 Implemented AI component**, resulting from [10.3.8](#).

## 11 Data-related considerations

### 11.1 Objectives

The objectives of this clause are:

- a) to define the dataset lifecycle of activities related to the gathering, creation, analysis, verification and validation, management, and maintenance of the datasets used in the development of the AI system;
- b) to identify the dataset insufficiencies that may impact the safety of the AI system;
- c) to identify the data-related safety properties that have a bearing on the safety of the AI system and that support dataset safety analysis;
- d) to define the countermeasures to prevent or mitigate dataset insufficiencies using dataset safety analysis methods at different steps in the dataset lifecycle;
- e) to define the data-related work products that support providing evidence of the safety of the AI system.

This clause applies to AI systems whose development and/or testing relies on data, including AI systems based on supervised, semi-supervised, and unsupervised learning techniques.

### 11.2 Prerequisites and supporting information

The following information shall be used at the initiation of this phase of activities:

- a) AI system definition, including:
  - 1) AI safety requirements, [Clause 9](#);
  - 2) input space definition (refined), [Clause 9](#);
- b) field data and functional insufficiencies detected during operation, [Clause 14](#);
- c) safety analysis report, [Clause 13](#).

### 11.3 General requirements

**11.3.1** A dataset lifecycle shall be defined for the datasets used in the development of the AI system.

**11.3.2** The dataset lifecycle shall be defined such that it supports iterative development of the dataset taking into account changes in the AI safety requirements and any insufficiencies observed during the AI system deployment phase.

**11.3.3** The dataset lifecycle shall include activities that relate to the gathering, creation, safety analysis, verification, validation, management and maintenance of the datasets used to develop the AI system.

NOTE An example dataset lifecycle covers requirements development, design, implementation, verification, validation, safety analysis and maintenance of the dataset. The dataset verification activity can ensure traceability from the dataset requirements to the dataset design and implementation. The dataset validation can involve integration of the AI system and be performed as part of the AI system verification.

11.3.4 Data-related safety properties of the dataset shall be identified and be used as inputs at different phases of the dataset lifecycle.

11.3.5 The dataset lifecycle activities shall include safety analyses to identify potential dataset insufficiencies, their root causes and their potential to cause a violation of AI safety requirements.

11.3.6 Dataset requirements shall:

- address the dataset insufficiencies that can lead to violation of the AI safety requirements;
- specify countermeasures to prevent the dataset insufficiencies, to mitigate them, or both.

11.3.7 Traceability shall be ensured between the dataset requirements and the AI safety requirements.

## 11.4 Dataset life cycle

### 11.4.1 Datasets and the AI safety lifecycle

Datasets play a crucial role in AI system development and testing. ML, in particular, typically involves an off-line training process to determine values for the parameters of an AI model. Three types of datasets enable this training and its usage afterwards: training datasets, AI validation datasets and AI test datasets.

Example flows for dataset creation and supervised learning are shown in [Figure 11-1](#) and [Figure 11-2](#).

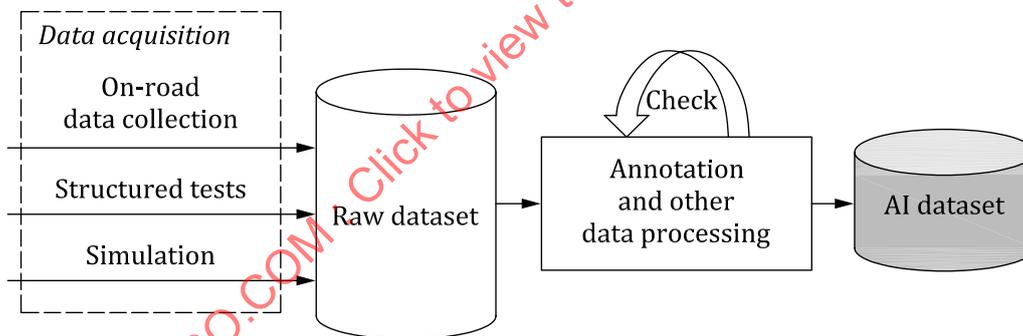


Figure 11-1 — An example dataset creation flow

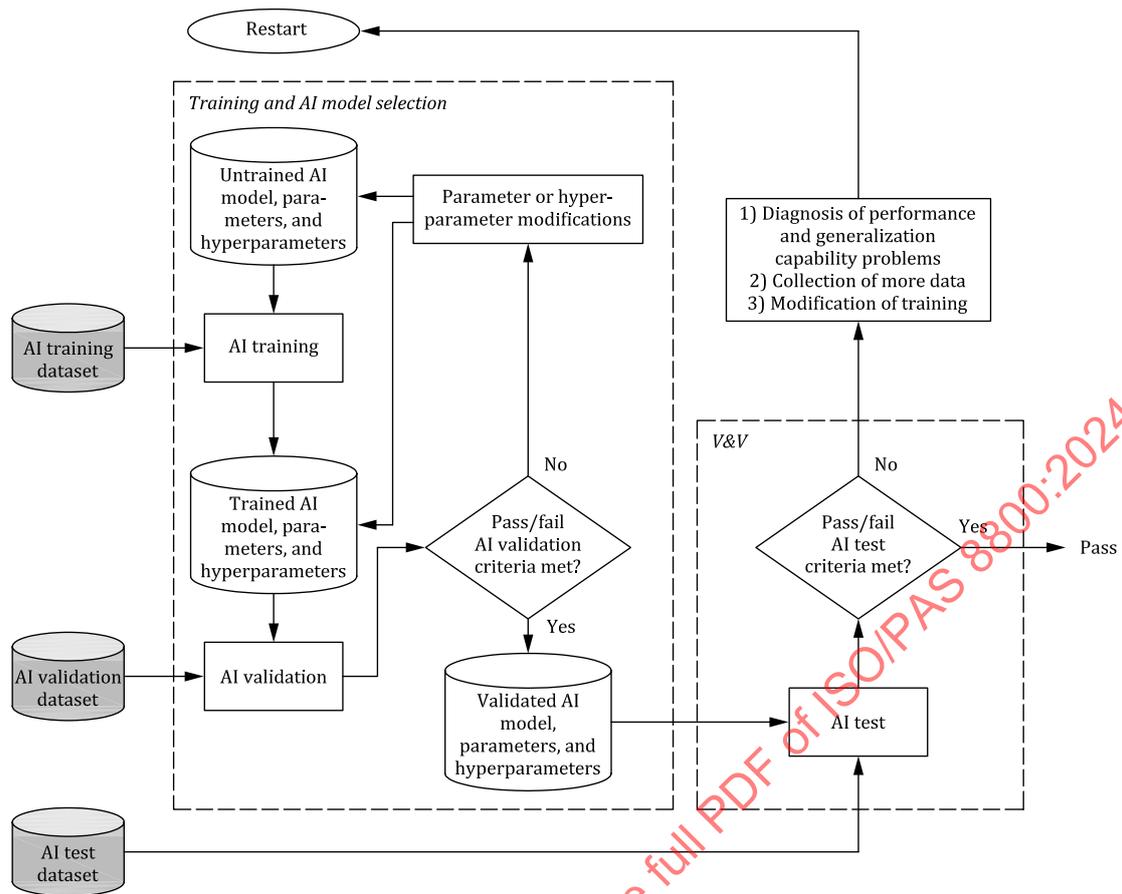


Figure 11-2 — An example supervised learning flow

The training dataset and AI validation dataset are used in the iterative AI training process of an AI model. The training dataset is input into the AI model while its parameters and hyperparameters are optimized based on the AI model’s performance. This proceeds until the predetermined AI training pass/fail criteria are reached. The AI validation dataset is then input into the AI model and the AI model is evaluated against the AI validation pass/fail criteria. If the results are not satisfactory, hyperparameters of the AI model are refined and the AI training process is repeated.

Once the AI training is complete, the AI model is evaluated with the AI test dataset using the AI test pass/fail criteria as part of the verification and validation activities. If the verification or the validation fail, the process is continued after more data are collected and/or training is modified.

#### 11.4.2 Reference dataset lifecycle

A typical dataset lifecycle describes the set of data-related activities carried out during the entire lifecycle of AI system development, including after deployment. The lifecycle serves as a means to manage the datasets and supports the realization of the AI safety requirements (and ultimately the safety requirements of the encompassing system).

A dataset lifecycle can consist of the following phases:

- dataset safety analysis;
- dataset requirements development;
- dataset design;
- dataset implementation;
- dataset verification;

- dataset validation;
- dataset maintenance.

A dataset lifecycle is created for the AI training, AI validation, and AI test datasets used in the AI system development (an individual dataset lifecycle can be created for each dataset role, if appropriate).

A dataset lifecycle can be aligned with or defined as part of the dataset creation and management activities at the level of the encompassing system, since system-level validation typically also relies on datasets.

Figure 11-3 provides an example dataset lifecycle based on the traditional V-model of development. Some of the salient features of the lifecycle are traceability of AI safety requirements to the dataset requirements (which impact the dataset’s design and implementation) and an iterative workflow that extends into operation, where new data can influence dataset revision.

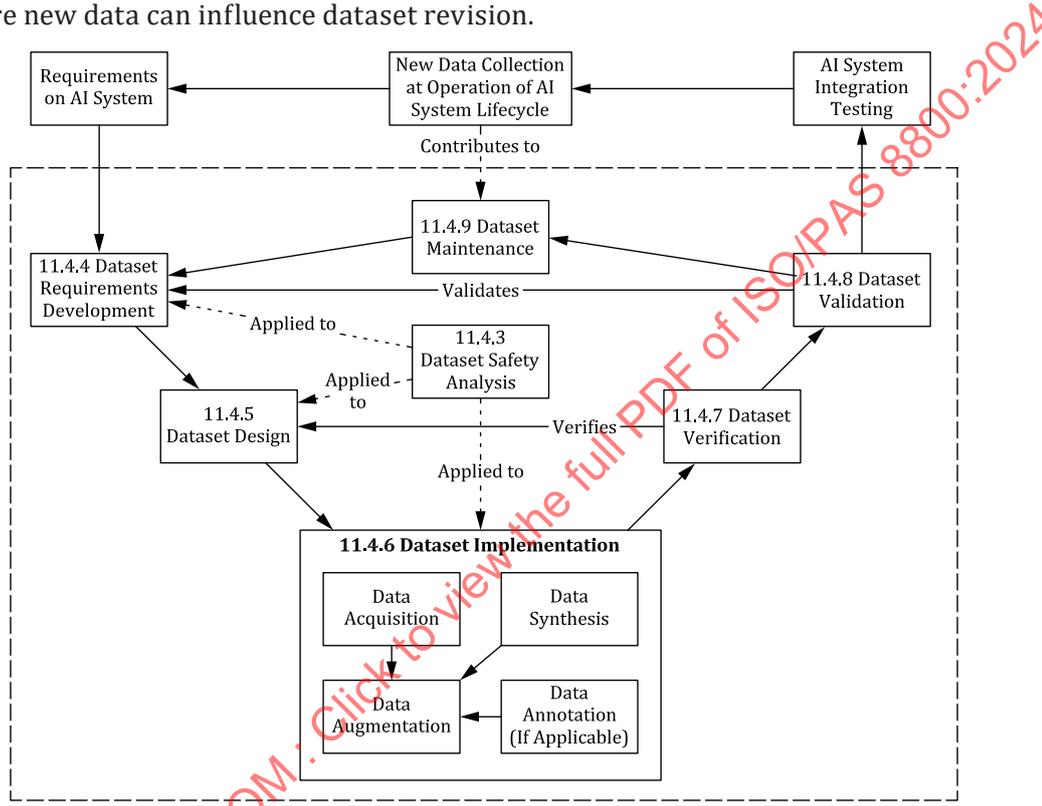


Figure 11-3 — Dataset lifecycle model

Subclauses 11.4.3 to 11.4.9 discuss each dataset lifecycle phase in more detail.

### 11.4.3 Dataset safety analysis

#### 11.4.3.1 General considerations

Dataset safety analyses focus on identifying safety-relevant dataset insufficiencies. When these dataset insufficiencies have been examined and the causes and consequences have been identified (including risks of the AI system and encompassing system), that information is fed as inputs to the dataset requirements development, dataset design, and dataset implementation to realize:

- countermeasures to prevent or mitigate dataset insufficiencies;
- metrics to judge effectiveness of measures to avoid dataset insufficiencies.

Depending on the dataset lifecycle phase, different approaches to dataset safety analyses can be used as outlined below.

Dataset requirements development phase:

- A guideword-based approach (such as that applied in HAZOPs) can be used to identify how dataset insufficiencies impact the safety of the AI system. Using this approach, one can determine what characteristics of the dataset(s) lead to the AI system performing a function incorrectly, seldom, too often, too little, too early, or too late.
- A qualitative risk analysis approach can be used to define the rigour applied in the dataset lifecycle to avoid dataset insufficiencies. This approach is similar to the application of HARAs to determine ASIL in ISO 26262. The approach considers:
  - the severity of the outcome associated with dataset insufficiencies (including the risks of the AI system and encompassing system);
  - the likelihood of the outcome associated with dataset insufficiencies;
  - existing countermeasures that can prevent or mitigate dataset insufficiencies;
  - additional countermeasures that can be applied to prevent or mitigate the dataset insufficiencies and the degree of rigour with which such countermeasures can be applied.

Generally, more significant risks at the AI system and encompassing system can be associated with a qualitatively higher assurance level and hence a greater rigour in the robustness of the countermeasures that can be applied.

Last, a residual risk analysis approach can be applied after the impact of the countermeasures in risk reduction has been considered.

Dataset design phase:

Various deductive and inductive analysis approaches can be conducted considering the proposed dataset design and can generate additional countermeasures that were not originally introduced at the dataset requirements development phase.

Dataset implementation phase:

A process failure mode and effects analysis (PFMEA) approach can be employed to identify potential issues in the processes, methods and tools of the data preparation and labelling and link these issues with the dataset insufficiencies and/or violations of the AI safety requirements (see [Clause 9](#) of this document and ISO 26262-9).

**EXAMPLE** An AI system classifies in-path objects for alert and trajectory planning purposes. The AI system uses a colour image from a camera mounted on the windshield. The AI training data, in particular, is hand-labelled. A PFMEA raises issues that can generate safety-relevant dataset insufficiencies such as:

- using images that were collected from the same camera on a different vehicle with different mounting and calibrations (data reuse impact);
- using images that were collected 10 years ago, even though there has been a considerable change in the types of vehicles that are on the road in the domain of interest (data ageing impact);
- using images that were collected from one region even though the system is targeted to operate in multiple regions with varying types of in-path object behaviour (data bias impact);
- using images that have not been labelled in a standard and correct manner (labelling inconsistency impact).

A summary of outputs of the dataset analyses serves as evidence to demonstrate that safety-relevant dataset insufficiencies are prevented or sufficiently mitigated. This summary can motivate or reference further artefacts (e.g. dataset tool qualification plan).

11.4.3.2 Dataset-related safety properties

Dataset insufficiencies are insufficiencies of the dataset regarding data-related safety properties under consideration. Examples of these data-related safety properties are given in [Table 11-1](#). They encompass both general properties and properties specific to AI applications.

**Table 11-1 — Examples of data-related safety properties**

Property	Definition
Accuracy	The data correspond to their source with respect to semantical representation and interpretation.
Completeness	The data elements (including metadata) are populated and the data have defined coverage of the input space, safety-relevant cases and plausible data perturbations.
Correctness (or fidelity)	The data correspond to the phenomenon they intend to capture and include features and metadata which help to characterize the phenomenon.
Independence of datasets	The datasets sufficiently avoid leakage of information amongst themselves with respect to data sources and the methods used to capture, gather, generate and process the data.
Integrity	The data are not altered by natural phenomenon (e.g. noise) or intentional action (e.g. usage of lossy data compression without consideration of impact to model, poisoning).
Representativeness	The distribution of data corresponds to the information in the environment of the phenomenon to be captured; it is free of biases.
Temporality	The data gives sufficient consideration to time-based characteristics (e.g. timeliness, ageing, lifetime, time contributing to distribution shift).
Traceability	The derivation of the data from their origin (including information on how they were captured, gathered, generated and processed) is demonstrated.
Verifiability	The data include sufficient features to be amenable for verification as prescribed by their requirements and properties.

NOTE 1 ISO/IEC 5259-1 and Reference [33] detail additional data-related safety properties (e.g. portability, understandability, auditability).

NOTE 2 Properties might not necessarily be mutually exclusive (e.g. a well-known property that is not listed is independence and identical distribution, or IID, which is covered by the correctness, completeness and independence properties).

Regarding a dataset insufficiency due to lack of independence of datasets, independence between the AI training and AI validation datasets supports detecting overfitting. Though not required, the K-fold cross validation technique can support this as the AI training and AI validation datasets are independent in each of the folds. Independence between the AI training and AI test datasets, on the other hand, supports providing a reliable statistical estimation of the residual risk of the trained AI component.

Regarding a dataset insufficiency due to lack of representativeness, biases can manifest in different forms. Human cognitive biases impact how engineering decisions are made and how datasets are sampled. Non-human cognitive biases (such as sensors failing during data collection) result in systematic dataset insufficiencies. More information on these forms can be found in ISO/IEC TR 24027:2021, Clause 6.

Generally, a bias in the training dataset impacts the performance of the AI system and is unwanted. However, intended biases can be used as a design measure to put the AI training focus on some important but rare features critical to the safety of the AI system (e.g. a training dataset with a higher occurrence rate of corner

cases can be used to complement an AI test dataset based on the real-world distribution). This design measure can still be insufficient to capture the true variability of rare safety-relevant cases, however.

Specific examples of dataset insufficiencies and the potential actions to avoid such dataset insufficiencies can be found in [Table 11-2](#).

**Table 11-2 — Examples of dataset insufficiencies**

Property	Dataset insufficiency example(s)	Potential actions to avoid insufficiency
Accuracy	<p>The resolution of camera images is not sufficient according to expected AI model inputs for object detection.</p> <p>An AI system operates with sensors detecting certain types of obstacles at a given distance range and high speed, but the camera used is not adapted for the range and speed, yielding blurry images.</p> <p>The mesh used in LiDAR imaging of objects is not fine enough (number of points, spacing) to properly detect target obstacles.</p>	<ul style="list-style-type: none"> <li>— Selection of source sensors appropriate for the input space and application;</li> <li>— Inspection of manually labelled data.</li> </ul>
Completeness	<p>Few images have obstacles close to the camera in a dataset for obstacle detection.</p> <p>No night images are in the dataset even though the input space includes night-time.</p> <p>Perturbations like noise, brightening, darkening, vibration, rotation, turbulence, blurring, blooming, smear and interference are not reflected in the dataset.</p> <p>An AI system for traffic signal identification is trained with a dataset that does not contain data elements that have all of the possible variations of traffic signal shape, height, positions, etc. outputted by the AI system.</p> <p>Missing information on the location of captured data does not allow one to analyse the geographical distribution of data and can cause undetected bias.</p>	<ul style="list-style-type: none"> <li>— Investigation of general use cases;</li> <li>— Calculation of distribution of the data and verification that the data cover the input space;</li> <li>— Collection of data from different geographical experts' perturbations on data that represent realities within the input space;</li> <li>— Addition of data through selection, generation, augmentation or synthesis if there are gaps identified (e.g. by analysis such as examination of saliency maps, training and testing);</li> <li>— Monitoring and collection of new or changing items within input space;</li> <li>— Corner and edge case collection.</li> </ul>

Table 11-2 (continued)

Property	Dataset insufficiency example(s)	Potential actions to avoid insufficiency
Correctness (or fidelity)	<p>Annotators manually create bounding boxes around objects inconsistently, which leads to object size per scenario to be calculated differently.</p> <p>No distinction has been made between a motorcycle and its rider in an image label, though this is relevant for the driving task that the AI system performs.</p> <p>A scene is marked as rain by an annotator although it is in snow.</p> <p>LiDAR provides ground truth for a camera outputting distance to an object. The vehicle collecting data drives through rain, which causes the ground truth to be noisier than in nominal conditions.</p> <p>Adhesive body markers provide ground truth for a driver monitoring system outputting head position, and the markers shift during the data collections.</p>	<ul style="list-style-type: none"> <li>— Characterization of the essential features of the target phenomenon;</li> <li>— Determination of the adequacy of the sensors to detect, observe and capture the phenomenon;</li> <li>— Redundancy of sensors.</li> </ul>
Independence of datasets	<p>One frame in a sequence is in the AI test dataset and the next frame is in the training dataset. Due to frame rate, the frames are nearly identical.</p> <p>One frame of a certain geo-position is in the AI test dataset and another taken later at the same geo-position is in the training dataset. For object detection, this can compromise the independence of the datasets.</p> <p>All datasets used to develop an AI model come from the same exact environment (e.g. same city street, time intervals, weather conditions and traffic load).</p> <p>All datasets are collected relying upon the same means (e.g. only one sensor or database).</p> <p>The task for dataset creation is always based upon the same technique, algorithm or parameters and is conducted by the same person.</p>	<ul style="list-style-type: none"> <li>— Use of data management system;</li> <li>— Use of different sources of data;</li> <li>— Separation of the teams preparing the different datasets;</li> <li>— Use of different technical means for data capturing, e.g. different sensors, vendors and brands;</li> <li>— Deployment of different processes/methods for dataset creation, e.g. applying two algorithms for data sample generation.</li> </ul>

Table 11-2 (continued)

Property	Dataset insufficiency example(s)	Potential actions to avoid insufficiency
Integrity	<p>Corruption of hardware storage introduces error(s) in the dataset.</p> <p>Failure of database memory introduces error(s) in the dataset.</p> <p>Untrained/careless user inadvertently introduces inconsistent data element (e.g. altered label).</p> <p>Error is introduced during dataset processing/manipulation due to transfer over lossy channel.</p>	<ul style="list-style-type: none"> <li>— General inspection;</li> <li>— Analysis of robustness to adversarial attacks on the dataset (e.g. random erasing, corruption);</li> <li>— Standard access controls (e.g. authorized users with passwords, denial of service protection techniques);</li> <li>— Built-in features for integrity checks in databases and other data storage;</li> <li>— Inclusion of integrity check codes (e.g. CRC, checksum, hash) for storage and transfer over lossy channels.</li> </ul>
Representativeness	<p>An AI system for driver monitoring is going to be used in a region where drivers have an even distribution from 20 years old to 100 years old, but the dataset does not contain drivers above 80 years old.</p> <p>AI datasets collected by a heavy truck fleet can have geospatial bias for an AI system intended to perceive aspects of roadways with weight limits.</p> <p>Synthetic data do not capture certain real-world features, such as shadows in images, to which the AI system is sensitive.</p> <p>Real-world data have been captured with wrong sensor parameters, resulting in variances to which the AI system is sensitive.</p> <p>Sensor data have disturbances due to a bug on the camera.</p>	<ul style="list-style-type: none"> <li>— Analysis and comparison of the theoretical and experimental distributions of the phenomenon;</li> <li>— Distributional drift analysis.</li> </ul>
Temporality	<p>COVID-19 induced change in distribution of people wearing face masks, but the dataset for a driver monitoring system does not consider this.</p>	<ul style="list-style-type: none"> <li>— Inclusion of metadata containing details like time of creation and validity;</li> <li>— Usage of version control for the dataset.</li> </ul>

Table 11-2 (continued)

Property	Dataset insufficiency example(s)	Potential actions to avoid insufficiency
Traceability	<p>An image lacking information on its source appears to be complete under simple visual inspection and is integrated into a dataset.</p> <p>Two datasets containing the same category of data are integrated into a single dataset without their metadata and attributes. The original datasets and their metadata are deleted.</p> <p>Data samples that were randomly selected for a training dataset decrease performance of the AI model due to those samples being corner cases. The samples are subsequently removed from the training dataset, but their metadata do not have an attribute to label them as corner cases.</p> <p>A new optimization algorithm is applied to datasets without properly tracing its application in the data management process. The optimized datasets are still used to replace the older ones.</p>	<ul style="list-style-type: none"> <li>— Use of data management system;</li> <li>— Updating of data management process to account for all tasks impacting datasets;</li> <li>— Creation and inclusion of appropriate and sufficient metadata to the data element collected and synthesized.</li> </ul>
Verifiability	<p>A framework that relies on a random algorithm generates data samples that violate a safety indicator in a simulation run. Without sufficient mechanisms for reproducing the same run, the violation is likely unverifiable.</p> <p>No checksum, CRC or hash mechanisms were activated at any time in the data management process and the datasets cannot be integrity-checked.</p> <p>Dataset images cover a certain period of the year (e.g. autumn and winter), but the camera was not correctly configured with the actual date and the respective metadata is unreliable.</p> <p>Different camera vendors and brands were used to ensure the independence of the data. However, all images were mixed and processed together and the technical means were not recorded.</p>	<ul style="list-style-type: none"> <li>— Ensuring reproducibility of data generation;</li> <li>— Mechanisms to allow verification of data properties, e.g. built-in features in databases for integrity checks;</li> <li>— Discarding of datasets with unreliable or insufficient metadata;</li> <li>— Manual analysis;</li> <li>— Use of statistical sampling methods.</li> </ul>

**11.4.4 Dataset requirements development**

The dataset requirements development follows the activity flow below, assuming that the method specified in ISO 26262-3 regarding item definition and the method specified in ISO 21448:2022, Clause 7 regarding triggering conditions have been followed:

- a) comprehension of the AI system;
- b) dataset safety analysis;
- c) dataset requirements formulation;
- d) dataset requirements quality assurance.

## ISO/PAS 8800:2024(en)

The comprehension of the AI system activity of the dataset requirements development focuses on understanding the intended functionality of the AI system, including:

- the AI safety requirements, [Clause 9](#);
- the input space definition, [Clause 9](#).

The dataset safety analysis activity is performed in line with [11.4.3](#). The outputs are fed into the dataset requirements formulation.

The dataset requirements formulation activity focuses on formulating the dataset requirements that mitigate the risks associated with the output of the AI system. It specifies:

- the logistical aspects, addressing at least the following items:
  - where the dataset is stored;
  - who has access to the dataset, what type of access they have and when they have access, including consideration given to ensure that this dataset is safe from unintended editing;
  - how the dataset is version controlled and how changes are tracked;
  - requirements on the verification and validation processes to be employed to ensure that the data within the dataset is correct and appropriate for usage;
  - how stakeholders can report known vulnerabilities, risks or biases in the data and/or dataset during any of the dataset life cycle phases.
- the technical aspects, addressing at least the following items:
  - size of the dataset;
  - format of the data within the dataset, including what syntactic and semantic parameters describe the data and what the format for labelling is;
  - boundaries of the data within the dataset (driven by both ground truth and design decisions);
  - dataset's role (AI training, AI validation or AI test) and what ensures that it is sufficient for its given role, including limitations on how many times it can be used (to avoid overfitting);
  - constraints affecting creation of the dataset (e.g. region-specific data privacy regulations);
  - mitigations for the different manifestations of dataset insufficiencies detailed in [11.4.3](#);
  - methods to prevent undetected data failures.

Finally, the dataset requirements quality assurance activity focuses on ensuring that the dataset requirements follow the criteria given in ISO 26262-8:2018, Clause 6. These requirements are the following:

- traceable to the AI safety requirements;
- updatable and maintainable upon a change to the encompassing system, the AI system or input space;
- updatable upon exposure of an insufficiency in the AI system due to the discovery of new safety-relevant scenarios or other triggers.

Table 11-3 — Further considerations for requirements development

Requirements topic	Considerations
Input space	<p>For an ADS, “input space” is equivalent to “operational design domain”.</p> <p>An input space can also include the user demographic and user-driven parameters (e.g. a driver monitoring system can have an input space that includes drivers with face paint).</p>
Role of the dataset	<p>Regarding the AI test dataset role, the AI test dataset is used specifically in the AI test verification part of the process, which can be performed as part of vehicle-level verification and validation. Ideally, a large AI test dataset is used to uncover overfitting.</p> <p>In case of repeated performance measurements on the same AI test dataset, the statistical validity of the test results can be threatened by implicitly optimizing towards the AI test dataset. The dataset requirements can include countermeasures to address this. Example countermeasures include:</p> <ul style="list-style-type: none"> <li>— employing multiple independent AI test datasets for use across different iterations of AI system development;</li> <li>— restricting access to the AI test dataset such that KPIs measured on a random subset of the dataset are returned instead of detailed results.</li> </ul> <p>Datasets can support another role – providing a means to monitor the input space while the AI system is in operation. <a href="#">11.4.9</a> covers this use case in more detail.</p>
Boundaries of data	<p>An example of setting boundaries is a situation in which a data element for a particular dataset has the following format: &lt;Time&gt;, &lt;Day&gt;, &lt;Traffic heaviness&gt;. Time is bounded between 00:00 and 23:59 inclusive, day is bounded within the (Sunday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday) enumerations and traffic heaviness is bounded within 0 and 100 % inclusive.</p>
Constraints affecting creation	<p>An example of a constraint is a situation where a region in which a vehicle is going to operate has restrictions disallowing capturing images of individuals. This can introduce an unwanted bias if the AI system’s function is to classify pedestrians based on the images.</p>
Traceability	<p>Traceability from dataset requirements to AI safety requirements (including those which have been generated from consideration of the input space (refined)) can be evidence that the AI safety requirements have been sufficiently considered. An example of this is as follows:</p> <ul style="list-style-type: none"> <li>— An AI system has an AI safety requirement that specifies the acceptable FP rates for the function that the AI system performs in two different weather conditions: sunny and rainy.</li> <li>— The training dataset has dataset requirements that specify how much of its data are in sunny conditions and how much of its data are in rainy conditions. Similar dataset requirements are created for the other datasets.</li> <li>— These dataset requirements link to the AI safety requirement.</li> </ul>

### 11.4.5 Dataset design

The dataset design outlines details on:

- data elements that are collected physically, created synthetically and/or created through augmentation (and how augmentation is applied in dataset generation and on-the-fly in the AI training pipeline, if applicable);
- which aspects of the data elements comprise the core data;
- the metadata, including any ground truth data associated with objects of interest within the data elements (also known as the labelling specification for supervised learning);
- the operations to be performed on the dataset (e.g. filtering of irrelevant or invalid data, dimensionality reduction, de-identification of data for data privacy purposes, normalization of the data with respect to appropriate metadata parameters);
- any mechanisms to be realized for monitoring the distribution shift in the input data during operation and collecting additional data items for subsequent revision of the dataset.

The details outlined in the dataset design are to be documented and subjected to analysis to ensure that they preserve the AI safety requirements and the dataset requirements.

#### 11.4.5.1 Creation of data elements and identification of core data

There are three main approaches used for creating data elements in a dataset:

- 1) Physical collection of data elements: Data elements in the dataset are directly obtained using either the sensors used in the encompassing system or their surrogates. In the case of surrogate sensors, a gap analysis is done to assess the effects of the difference with the native sensors and countermeasures are taken to contain the effects.
- 2) Synthetic creation of data elements: Certain aspects of the input space might not have been captured during data collection and various simulation tools. Specialized machine learning methods like generative adversarial networks can be employed to create additional data elements that capture those aspects. The adequacy of the synthetic data elements is considered.
- 3) Augmentation of physically-created data elements: A physically-created set of data elements are augmented to create a new set of data elements that have parameters altered to include perturbations such as noise, brightening, darkening, vibration, rotation, turbulence, blurring, blooming, smear and interference. If augmented data are used for testing and if the data even partially replace real-world tests, the adequacy of the augmentation is considered.

**EXAMPLE** For computer-vision based data, the means of altering can include colour space transformation, background modification, superposition of multiple images, flipping, scaling, translation and random cropping.

The aspects of these data elements that serve as the inputs to the AI system during runtime operation of the AI system (e.g. the raw image data from a camera sensor, a three-dimensional 128 x 128 x 3 array of RGB pixel values, for an AI system built on computer vision) are the core data.

#### 11.4.5.2 Design of metadata for data elements and datasets

The metadata associated with a data element provides valuable information about the data element that can be used during training, analysis and verification. For a supervisory ML system, the metadata includes the ground truth or the label used during the training process. The data type, structure and range of values that the labels assume are also identified during the design phase. They conform to the functionality of the AI system and meet the dataset requirements.

**EXAMPLE 1** An AI system used to classify objects for an automatic emergency braking system depends on a training dataset for which the data elements have metadata containing the class of the object (e.g. car, truck, person), the height and width of the object (the bounding box) and the distance of the object from the subject vehicle.

Metadata can also be associated with a dataset as a whole. The metadata associated with a dataset serves to support analysis, verification and validation of the dataset. The metadata can contain the following:

- details on how the dataset was created, e.g. physical collection, details of sensor devices, synthetic methods and tools that were used for generation, augmentation, etc.;
- statistics of syntactic and semantic parameters of data elements in the dataset;
- information that relates the data elements in the dataset to the AI system and the encompassing system, input space, object/event detection and response, dynamic driving tasks, edge cases, etc.

**EXAMPLE 2** An image dataset can be defined to be such that 10 % of its data elements contain traffic signals in the top left corner. In the assessment of completeness, the range of values of various parameters and their combinations can be useful.

### 11.4.6 Dataset implementation

The activities in the dataset implementation phase realize a concrete dataset based on the dataset requirements and dataset design. The activities include:

- defining the processes, methods and tools to prepare a given dataset (e.g. physical, synthetic and/or augmented data generation, cleaning as covered in [11.4.5](#));
- preparing the dataset;
- defining the processes, methods and tools for labelling the dataset;
- labelling the dataset.

**NOTE** ISO/IEC 23053 identifies the data preparation step as a tool in the development of an AI system. The ISO/IEC 5259 series defines a data quality framework that can be used as guidance during this phase, and ISO/IEC TR 24368 provides guidance on having processes in place for stakeholders to disclose/report known vulnerabilities, risks or biases associated with the dataset preparation and labelling, which can then be fed into a dataset safety analysis.

Regarding labelling in particular, it is often a human-labour intensive process involving a label supplier, and label quality tests and audits are applied. The nature and extent of these tests and audits is commensurate with the complexity of the inputs and outputs being labelled.

**EXAMPLE 1** If the data to be used for an AI system development is LiDAR data involving point clouds, the labelling to identify objects in the inputs can be complex and error-prone. To address this, the labelling process can employ multiple levels of human involvement, e.g. labellers, reviewers and auditors, and possibly also semi-automated and automated label quality tests and plausibility checks.

Additionally, for labelling, consideration is given to how any ground truth is obtained to ensure that any instrumentation used to obtain the ground truth does not interfere with how the data are represented.

**EXAMPLE 2** Body markers used to provide the ground truth positioning of a person for a camera-enabled driver monitoring system can cause interference since the body markers will likely not be a part of the real-world operation distribution set.

As part of the dataset implementation activities, the details of records created during preparation and labelling are documented as inputs for dataset safety analysis and dataset verification activities. The coverage of the input and output spaces and the statistical distribution of the datasets are also recorded.

### 11.4.7 Dataset verification

Dataset verification applies to the dataset under evaluation. Its purpose is to confirm that the dataset has been developed correctly. This process comprises product verification complemented by process verification:

- product verification:
  - determining the consistency and correctness of information in a data element;

- determining the consistency and correctness of information at the dataset and metadata levels, e.g. lack of outliers, missing data elements, duplicates, wrong data types;
- verifying the conformance of the dataset against dataset requirements, e.g. metrics on distribution of parameters of the dataset, extreme values and edge cases of parameters and their combinations, noise characteristics, independence between datasets;
- process verification:
  - checking that the design and implementation phases are performed correctly;
  - checking the correctness of the processes, methods and tools used to create the dataset and its metadata (including any other AI systems involved in ground truth labelling).

**NOTE** Product verification can be done either manually or using automated tools, depending on the type of checking that is involved. For instance, verification of the information about the sensing device used for data collection can require manual inspection, while that of the ground truth label can employ running automated software. Checking the correctness of all data elements in a dataset can be impractical, so this can be done using statistical sampling approaches.

**EXAMPLE** For an AI system that is expected to learn certain high-level concepts (e.g. pedestrians, vulnerable road users, traffic signals), the dataset is required to have a sufficient number of data elements containing these concepts so that they can be learned. The AI system metric is that the AI system perceives vulnerable road users under certain lighting conditions with an accuracy of X%. This AI system metric relates to dataset requirements that Y% of the training dataset and Z% of the AI test dataset contain vulnerable road users under those lighting conditions. If these dataset metrics are not met, the dataset is enhanced with additional data elements.

Dataset verification is repeated every time dataset requirements are added or refined. The details of the verification carried out are documented as evidence for the verification of the dataset.

#### 11.4.8 Dataset validation

Dataset validation ensures the correctness of the dataset requirements from the dataset requirements phase, i.e. if the correct dataset and data elements are developed for the AI system with the desired safety properties and if they reflect a correct translation of the AI safety requirements.

There are two approaches to dataset validation activities which can be applied together or individually:

- a) requirement conformance, which involves checking that the derived dataset requirements meet the expected objectives of the dataset;
- b) integration testing, which involves checking that the AI system developed using the dataset(s) (i.e. trained and tested with the dataset(s)) meets the AI safety requirements.

For the first approach, the expected objectives of the dataset are adequately articulated at the AI system development phase and handed over to the dataset development team. Often, the expected objectives are in terms of use cases and edge cases of the AI system, and checking that these are met is done by examining and reviewing the dataset requirements. As part of this evaluation, the consistency and completeness of the dataset requirements can be assessed.

In addition, requirement conformance can be carried out by examining the AI safety requirements. Every safety requirement that has a potential impact on a dataset is covered by one or more dataset requirements (although there can be dataset requirements that do not trace to an AI safety requirement). Requirement conformance involves checking that a correct and desired traceability exists, and can be carried out manually and/or using some automated support of analysis of the requirements.

The second approach to dataset validation is integration testing. In this approach, the AI system is derived or revised using the dataset, and the resulting AI system is verified against its requirements, with the additional objective of the AI system verification being to check that the right dataset was developed. Any failure in AI system verification can be traced to deficiencies in the dataset (and subsequently accounted for in the dataset requirements or design) or to other issues, like an inadequate AI system architecture or an inadequate training process.

The dataset validation phase results in evidence describing the relationship of dataset requirements to the AI safety requirements and summarizing the review results.

#### 11.4.9 Dataset maintenance

Dataset maintenance refers to the set of activities that ensure that a dataset is up-to-date and compliant with the dataset requirements. These activities are carried out across the entire dataset lifecycle.

According to ISO 26262-8, the dataset maintenance activities can include:

- configurations of the dataset, including what they are, what they do and how they are managed;
- management of dataset resources, tools, repositories, access rights and timelines;
- change management of datasets including what triggers changes to the dataset;
- retirement and decommissioning of dataset and their elements.

NOTE Guidance on retirement and decommissioning are available in ISO/IEC 5338 and ISO/IEC 8183.

Dataset maintenance activities include actions taken during operations (see [Clause 14](#)) and involve:

- general field data collection and monitoring (e.g. monitoring the inputs to the AI system for conformation to the AI safety requirements and identifying data elements corresponding to safety-relevant edge cases encountered during operation);
- OOD data identification, collection and processing;
- AI system adaptation.

**Table 11-4 — Examples of dataset maintenance activities**

Sub-category	Example(s)
OOD data identification, collection and processing	<p>When an AI system fails to detect certain objects, causing unwanted behaviours at the encompassing system level, data with those objects are collected and uploaded to fine-tune the AI model.</p> <p>An infrastructure feature that was in the AI system's input space has become obsolete. Data elements containing that feature are removed from the dataset.</p> <p>An AI system for driver monitoring performs sub-optimally when the driver is wearing a certain thickness of glasses. Data elements containing people wearing that thickness of glasses are added to the dataset.</p>
AI system adaptation	<p>When an AI system is deployed in another country, it might be expected to detect different signage and adjust system behaviour according to local laws and regulations. Data elements containing that signage are added to the dataset.</p> <p>An ADS that employs an AI system goes from operating on limited access highways to also operating on highways with at-grade crossings. Data elements containing at-grade crossings are added to the dataset.</p>

### 11.5 Work products

**11.5.1 Dataset lifecycle**, resulting from [11.3.1](#) and [11.3.2](#).

**11.5.2 Evidence for the outputs of the defined phases of the dataset lifecycle**, resulting from [11.3.3](#).

**11.5.3 Evidence for the safety analyses of the dataset**, resulting from [11.3.4](#) and [11.3.5](#).

**11.5.4 Dataset requirements specification**, resulting from [11.3.6](#) and [11.3.7](#).

## 12 Verification and validation of the AI system

### 12.1 Objectives

The objectives of this clause are:

- a) to verify that the AI system fulfils its AI safety requirements;
- b) to validate that the safety requirements allocated to the AI system are achieved when integrating into the encompassing system;

NOTE 1 This clause includes guidance for:

- the stand-alone performance analysis of the AI system itself;
- testing at the AI system and AI component level. Testing at the AI component level can include AI model, and pre- and post-processing elements.

NOTE 2 Within the development of AI systems, the term “validation” is typically used differently to how it is used within safety standards and system safety engineering. In this document, the term “AI system safety validation” ([3.1.21](#)) is used.

### 12.2 Prerequisites and supporting information

The following information shall be available to complete the verification and validation activities associated with the corresponding phases of the AI safety lifecycle:

- a) safety requirements allocated to the AI system (from external sources, e.g. the encompassing system development);
- b) AI safety requirements, [Clause 9](#);
- c) known insufficiencies of the AI system and the corresponding subdomains of the input space, [Clause 9](#);
- d) input space definition (refined), [Clause 9](#);
- e) AI component or AI system architecture, [Clause 10](#);
- f) implemented AI component, [Clause 10](#);
- g) dataset lifecycle, [Clause 11](#);
- h) evidence for the outputs of the defined phases of the dataset lifecycle, [Clause 11](#);
- i) evidence for the safety analyses of the dataset, [Clause 11](#);
- j) dataset requirements specification, [Clause 11](#).

### 12.3 General requirements

**12.3.1** The AI system shall be verified to provide evidence for:

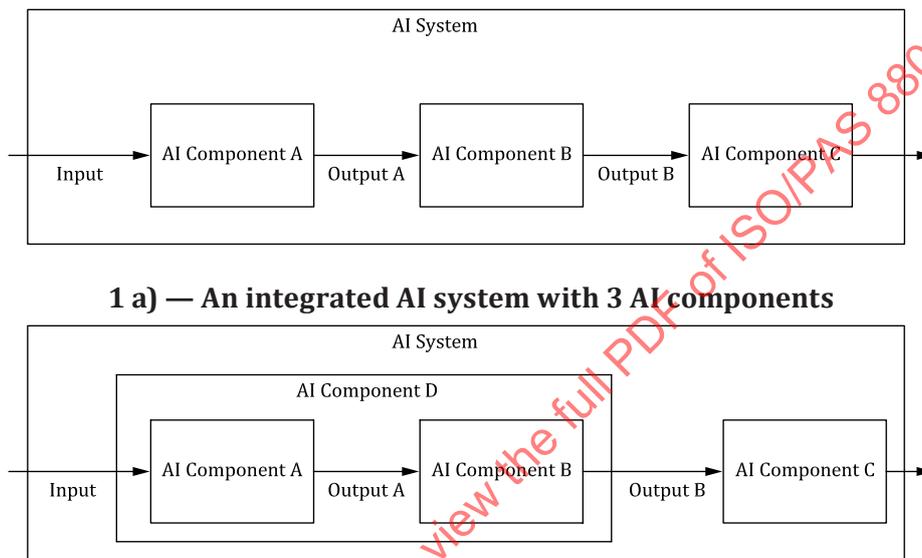
- a) conformity to the AI safety requirements;
- b) confidence in the absence of unintended functionality and properties.

NOTE Confidence in the absence of unintended functionality and properties can be increased, for example, by following safety standards such as ISO 26262, ISO 21448 and this document during development.

**12.3.2** Testing of an AI system shall be performed on the AI components that can be tested stand-alone and on the integrated AI system.

**NOTE** If an AI component cannot be tested stand-alone, the AI component is tested at a higher level of integration. In this case, the AI system integration strategy is devised to accommodate this testing.

**EXAMPLE** [Figure 12-1](#) a) shows an AI system with three AI components. Assume AI component C can be tested stand-alone, but AI components A and B cannot. For example, suppose AI components A and B implement a convolutional neural network to propose bounding boxes and a non-max suppression algorithm to integrate them. In another example, suppose AI components A and B implement a backbone/head (neural network), i.e. feature extraction such as vision transformers, shared with other AI systems and a task-specific neural network for this AI system. Output A is an intermediate value in these cases and cannot be tested stand-alone. However, AI components A and B can be tested if integrated together. In such a case, the integrated AI component is considered and tested stand-alone. [Figure 12-1](#) b) shows the integrated AI component (AI component D) within the AI system. Here, testing can be performed on AI components C and D and on the integrated AI system.



**1 a) — An integrated AI system with 3 AI components**

**1 b) — AI components A and B are integrated so that the integrated AI component D can be tested stand-alone**

**Figure 12-1 — Example of integrating AI components for AI system testing**

**12.3.3** Test cases for the verification of the AI components shall be derived using best practices for test case derivation. This includes using an appropriate combination of the methods listed in ISO 26262-6:2018, Clause 9, i.e. analysis of the requirements, generation and analysis of equivalent classes, analysis of boundary values and error guessing based on knowledge or experience.

**EXAMPLE 1** Analysis of the requirements, including the required safety properties, can be used to select KPIs for the V&V activities.

**EXAMPLE 2** Generation and analysis of equivalent classes can be suited to generate complete test sets for pre- and post-processing algorithms.

**EXAMPLE 3** Error guessing based on knowledge or experience can be suited to identify yet unknown edge cases for testing.

**EXAMPLE 4** Analysis of boundary values can be suited to generate complete test sets for pre- and post-processing algorithms.

**NOTE 1** For ML, analysis of requirements relies on statistical tests to analyse whether the safety relevant performance requirements are met.

**NOTE 2** If relevant, error guessing can include evaluation of known and potential triggering conditions and evaluation of known and potential functional insufficiencies.

NOTE 3 The term “knowledge” is interpreted broadly beyond human knowledge and can include knowledge automatically derived by an algorithm.

NOTE 4 For AI models used in perception modules of an autonomous driving vehicle, verification involves driving scenes in test datasets. See ISO 34502 for more details on creating test scenarios for such models.

**12.3.4** Each test case of an AI component shall include pass/fail criteria.

NOTE Pass/fail criteria can be based on the formulation of the thresholds and parameters provided in the AI safety requirement allocated to the AI component, if applicable.

**12.3.5** Test cases of an AI component shall adequately verify the AI safety requirements allocated to the AI component within the specified input space of the AI system.

NOTE 1 AI test quality and safety-aware AI testing are considered for these test cases. AI test quality refers to the need for rigorous testing of AI models that goes beyond any simple mean performance calculation with a single test dataset. Safety-aware testing refers to testing that uses safety-aware metrics and safety-relevant data points or subsets.

EXAMPLE 1 Test cases are designed to verify the AI model in terms of OOD performance and in-distribution performance on samples underrepresented in the training data.

EXAMPLE 2 Test cases reflect data points contributing to unsafe states of vehicles deductively enumerated by safety analysis such as design FMEA or FTA and inductively known from past products.

NOTE 2 For automated driving applications, the completeness and sufficiency of the test cases can be evaluated considering the acceptance criteria defined in ISO 21448.

NOTE 3 If the AI task is implemented by multiple AI models, the relevant sub-domain of the input space for each AI model is defined, e.g. one AI model can be used to explicitly identify vulnerable road users (VRU) while another can be used to explicitly identify traffic signs, resulting in a relevant input space subdomain VRU and traffic signs. The test cases are then more focused on the relevant input space subdomains and less on the overall input space of the AI system.

**12.3.6** The AI system integration approach shall specify the steps for integrating the individual AI components hierarchically into higher level AI components until the AI system is fully integrated.

**12.3.7** The AI system integration shall be verified to provide evidence that the hierarchically integrated AI components and the integrated AI system achieve:

- a) conformity to the AI system architectural design in accordance with [Clause 10](#);
- b) satisfaction of the AI safety requirements.

**12.3.8** AI system safety validation shall confirm that the safety requirements allocated to the AI system are fulfilled when the AI system is integrated into the encompassing system.

## **12.4 AI/ML specific challenges to verification and validation**

AI systems, especially those developed using data-driven methods, pose unique challenges to verification and validation, such as:

- They lack a precise statement of AI safety requirements. AI systems are often expected to identify and quantify human-interpretable high-level semantic concepts like road objects, traffic signals, attendant driver, etc. Many of these concepts lack precise definitions and are difficult to capture in terms of mathematical descriptions.
- The inputs to an AI system are often versatile/diversified and from different sources, e.g. radar, LiDAR, camera etc. whose representation requires high dimensional objects with a very large range of values and satisfying complex and unknown constraints. Use of traditional input coverage methods would be incomplete or expensive.

- AI systems involving DNNs employ complex architectures, especially for perception-based applications (e.g. long short-term memory networks and encoder-decoder networks) and contain several layers with millions of parameters whose values are tuned during the training process. Verification, i.e. checking that these parameters have the desired values for optimum performance and to satisfy AI safety requirements, leads to scalability issues for many realistic applications.
- Training AI models involves the use of various heuristics for identifying parameter values to optimize appropriate cost functions. These heuristics can lead to locally optimum parameter values without achieving a globally optimum solution (i.e. model generalization) that can lead to AI errors in the AI system. The verification task involves checking that the computed parameters are adequate to satisfy the AI safety requirements.
- AI systems rely on a large dataset for their reliable performance. Demonstrating the validity and completeness of the verification dataset is a non-trivial task given the complexity and scope of the input space.
- The non-predictable erroneous behaviour in data-driven AI systems, for example, based on spurious correlations, limit the ability to predict performance based on a review of the training data or the implemented model.

NOTE This property is closely related to robustness. This challenge can be further exacerbated by the lack of explainability of the trained function when using technologies such as DNNs.

- Limitation of structural coverage: Due to the lack of a detailed specification as well as the dependency on parameter values (e.g. weight in an NN) during execution, both black box and white box coverage metrics have limited use when extrapolating the results of executed tests in evaluating performance over the entire input space.
- Stability of performance due to changes in the environment or the function. Small changes in the input space (e.g. one- or two-pixel values in an image input) can lead to, as yet undiscovered, AI errors in the function. Furthermore, changes to the function (due to re-training) can lead to an unpredictable impact on previously verified properties.
- An AI system while training with an inadequate number of examples can reach a local optimum that results in behaviours not aligned with the desired outcome.

## 12.5 Verification and validation of the AI system

### 12.5.1 Scope of verification and validation of the AI system

[Figure 12-2](#) shows the phases where verification and validation activities of the AI system happen. Verification of the AI system is applicable to the following phases of the AI system safety lifecycle:

- a) When defining the AI safety requirements ([Clause 9](#)), verification ensures that the AI safety requirements are correct, complete and consistent with each other and with respect to the encompassing system technical safety concept and safety requirements. Verification of the AI safety requirements can be performed following ISO 26262-8:2018, Clause 9 and ISO 21448.
- b) During the development phase of the AI system, verification is conducted in different forms, as described below:
  - During the design phase ([Clause 10](#)), verification is the evaluation of the work products, such as architectural design, models or architectural measures, thus ensuring that they meet the AI safety requirements for correctness, completeness and consistency. Evaluation can be performed by methods such as review, simulation or analysis. It is planned, specified, executed, and documented in a systematic manner following ISO 26262-8:2018, Clause 9.
  - In the data lifecycle ([Clause 11](#)), data is verified at each phase for sufficient correctness, consistency and completeness.

## ISO/PAS 8800:2024(en)

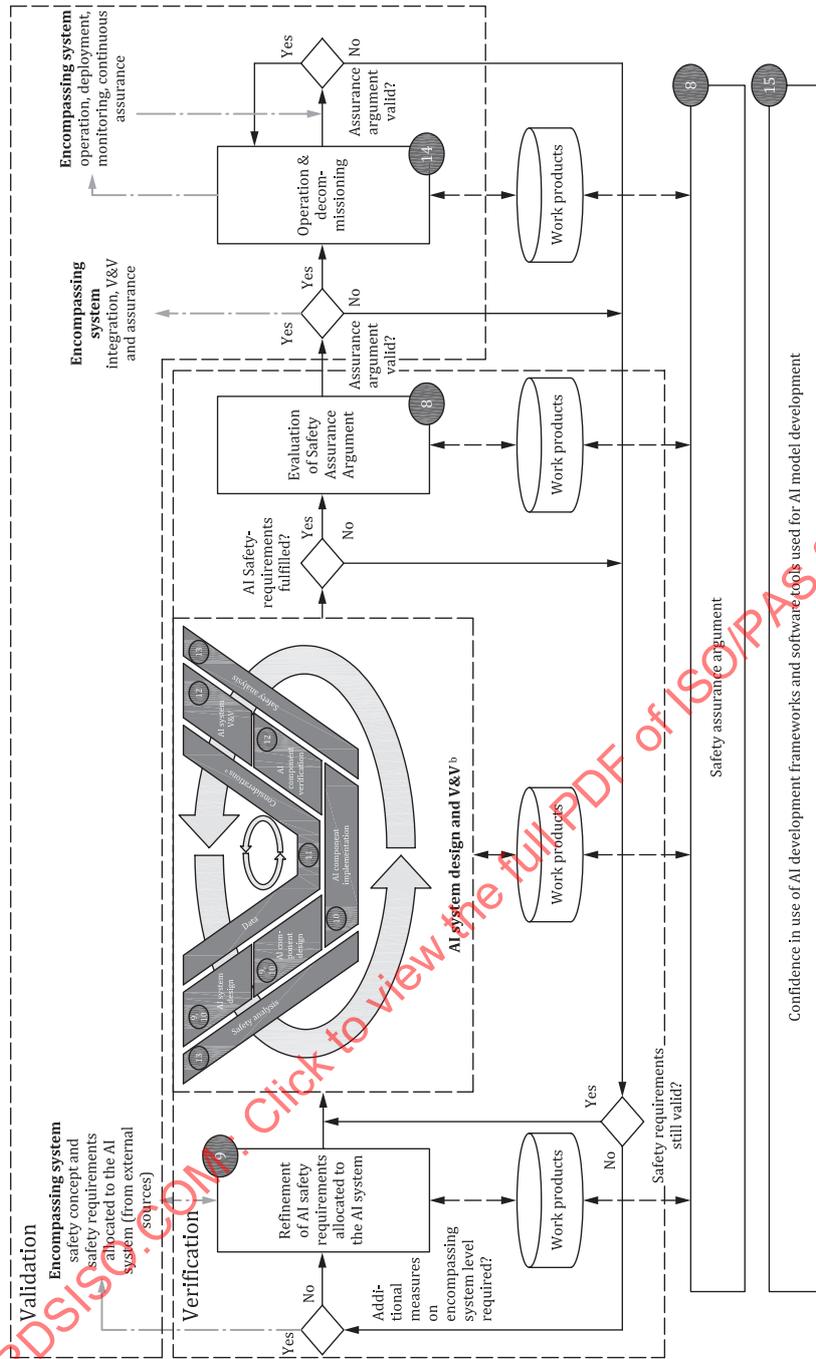
- In the test phase, verification of the AI system is the evaluation of the work products and elements within a test environment to ensure that they conform to the AI safety requirements. The tests are planned, specified, executed, evaluated and documented in a systematic manner.

Testing of an AI system is performed at different levels of system integration.

- AI component testing ([12.5.2](#));
- testing of the integrated AI system ([12.5.4](#));
- integration of the AI system with the encompassing system and AI system safety validation ([12.5.7](#));
- post-deployment validation ([Clause 14](#)).

In [12.5.2](#) to [12.5.7](#), some guidance on the verification and validation of AI systems is provided whilst addressing AI/ML specific challenges.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024



Key



clause number(s)

← process flow to/from the development and operation of the encompassing system

→ process flow

↔ relation



represents the iterative nature, in particular of the AI component design and verification

a Specific to ML-based AI technologies.

b See [Figure 7-2](#).

Figure 12-2 — Phases of verification and validation activities of the AI system

12.5.2 AI component testing

12.5.2.1 Testing workflow of an AI component

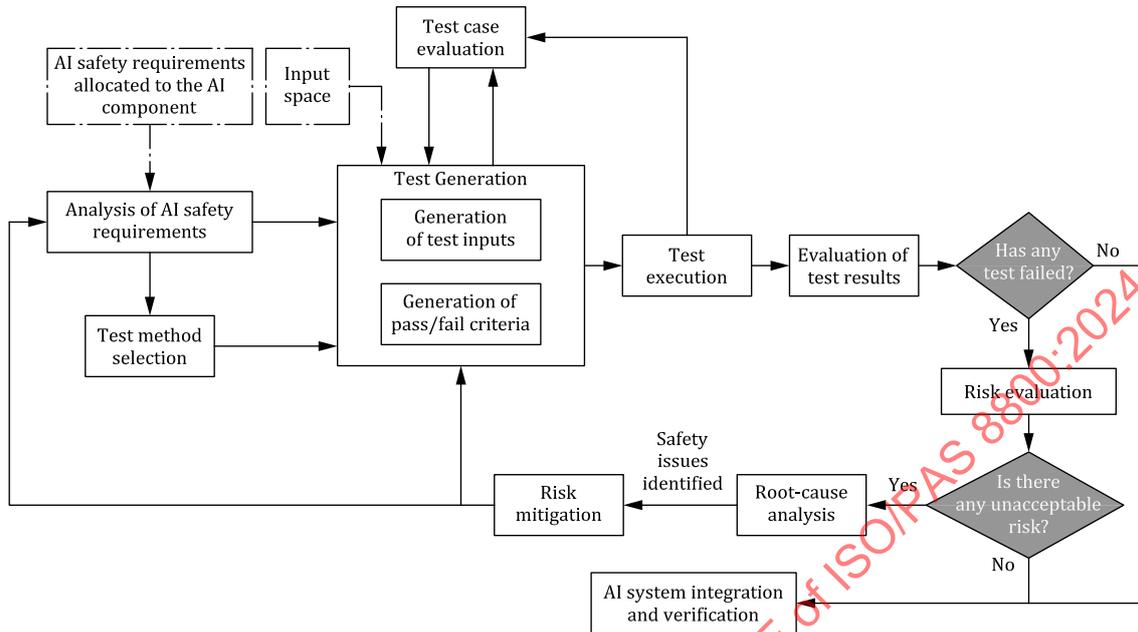


Figure 12-3 — Testing workflow of an AI component

Figure 12-3 shows a typical workflow for the testing of an AI component, based on the process presented in Reference [34] for the testing of ML-based systems. Initially, the AI safety requirements allocated to the AI component are analysed to select the most appropriate combination of test methods. The AI safety requirements can also provide direct input to test case generation approaches which involve (automatically) deriving test cases for the AI component under test. The result of the analysis is a definition of expectations on the inputs, methods, environment and tools for performing the verification.

NOTE 1 Test planning of an AI component includes gathering information about related test methods, test evaluation criteria, such as test coverage or number of tests passed, etc.

NOTE 2 AI safety requirements include safety properties of the AI system and safety performance indicators.

This is followed by the formulation of the pass/fail criteria (test oracles) to be applied to the test results. These criteria can include a combination of coverage criteria (e.g. from the input space or from the AI model) as well as performance targets (e.g. at most x% of inputs lead to an output with a deviation from the ground truth of no more than y%) and safety-related KPIs.

Test inputs are sampled from the collected data or generated using simulation or other methods. See 12.5.2.2 and 11.4 for more details on test input generation. Test cases are evaluated to ensure that they sufficiently cover the scenario space considered and that they completely verify AI safety requirements within the input space of the AI system.

NOTE 3 Adequacy of the test cases for verification of the AI safety requirements can also be evaluated upon the execution of the tests. This is an ongoing area of research and the idea is similar to the structural coverage analysis for conventional software. See 12.5.2.3 for more details.

During the test execution, the AI component is run with the test inputs and test results are generated. The test results are evaluated with respect to the defined test oracles.

NOTE 4 When evaluating the risk due to AI errors, one can consider the required target metrics and KPIs. It is possible that an AI safety requirement is not violated by the occurrence of a single AI error if the corresponding target metrics and KPIs (as defined in Clause 9) associated with the AI safety requirement are nevertheless met.

If a test fails, a safety analysis is conducted to evaluate the impact of the AI error(s) on safety. If the risk due to the failed test is deemed unacceptable, the root cause(s) of the AI error(s) are investigated. Then, depending on the potential root cause(s), appropriate mitigation measures are applied to reduce the risk. If the AI component is modified as a result of a risk mitigation measure, the AI component will be retested. The safety analysis (including risk evaluation, root cause analysis and identifying mitigation measures) will be discussed in more details in [Clause 13](#).

NOTE 5 It is assumed that at the test execution stage, the AI model is ready to be tested. For ML models, this means that the model is already trained and model parameters are set and the model is ready to be tested.

### 12.5.2.2 Test generation

Test generation involves identifying a test dataset that determines whether the AI system meets its AI safety requirements. The test dataset contains data elements which are representative of the inputs that the AI system receives in operation.

EXAMPLE An AI system used in a camera-based perception component of an automated driving system might receive road scenes as inputs which are dependent upon the ODD of the underlying feature. The test input scenes would then be required to have sufficient coverage over the entire ODD.

NOTE Test generation for validation at the vehicle level is described in ISO 34502.

The AI safety requirements can include some requirements related to AI errors. Testing these requirements requires generating edge cases based on the input space definition. Applying methods for the systematic exploration of the input space can support an argument that both nominal and edge case conditions triggering AI errors are covered during testing. In general, there can be an expected distribution of nominal and edge cases over the input space. The test cases are representative of the input space and capable of uncovering critical AI errors.

In general, it is difficult to achieve the desired distribution for the test data when generating the data from actual measurements taken from the environment and using the target sensor set. Simulation environments can be used to generate additional (synthetic) test data to achieve the required coverage and distribution over the input space. The input space can be described by a set of constraints on the input parameters and the desired coverage can be expressed and achieved using techniques such as Design of Experiments. When synthetic test data is used, the test evaluation activities include checking the validity of the generated test data. Due to the black box nature of certain ML methods (e.g. CNN), formal validation of the generated data is not possible. It is also not feasible to determine what features of the data the AI model is sensitive to (e.g. shadow, edge softness, blurs). Proving that insights gained on synthetic data can be transferred to real world data involves training multiple models on various composition of real and artificial data and comparing the performances.

### 12.5.2.3 Test case evaluation

When test cases are generated, they are evaluated to make sure they:

- a) are appropriate and correct to evaluate the AI safety requirements allocated to the AI component, in particular regarding the expected outcome (the ground-truth);
- b) adequately cover the AI safety requirements allocated to the AI component within the specified input space of the AI system;
- c) effectively cover the scenario space considered.

In SOTIF, a triggering condition is always characterized by some range of parameters. For example, for a pedestrian wearing black clothing, there can be different interpretations, reflecting a range of RGB values that can be defined as black. Another example is the rain intensity, where for heavy rain, the intensity of precipitation is not a concrete value but instead included in a pre-regulated range. Therefore, the set of test cases ensures coverage of the range of values within the syntactic space that correspond to triggering conditions defined in terms of the semantic input space.

Once the input space has been characterized, methods such as combinatorial testing (see 12.5.3) can be used to perform such a check on relative completeness, where each relevant input dimension is partitioned into a set of finite valuations (e.g. for “fog”, further partition the intensity into three discrete classes “strong”, “moderate” and “low”), followed by ensuring all combinations being considered in the set of test cases.

NOTE Evaluation of test cases at the vehicle level is described in ISO 34505.

For covering the scenario space, for deep neural networks, a complementary method is to perform interpretability analysis (e.g. saliency maps) to understand the scenario or feature of why a particular neuron gets excited or suppressed. This can be used to define a coverage criterion to ensure the relative completeness of learned-feature combinations (e.g. neuron coverage and its variations or k-way combinatorial testing on neuron activations). However, achieving full coverage using such methods might not be possible, and methods to artificially ensure coverage can lead to generating random images that are never observed in practice.

### 12.5.3 Methods for testing the AI component

A suitable testing method of an AI component needs to consider, under available resources (e.g. time or computing power), both the breadth (to efficiently covering the input space) and depth (to enable effective issue/error detection) of testing. Practically, a portfolio of diversified methods is used for testing AI components. The following is an incomplete list of methods for testing AI models:

- Statistical testing: This method evaluates the achieved values of the metrics defined within the AI safety requirements and associated safety-related properties within a given confidence interval. Effective experimental design plays a crucial role in statistical inference (including statistical hypothesis testing) by safeguarding the validity and reliability of findings, thus preventing the production of spurious or distorted associations.
- Data/scenario replay: This method refers to collecting a set of scenarios (for example, recorded during test drives) and subsequently using the collected data to stimulate the AI model under test and evaluate the responses. Examples of these scenarios include known pre-crash scenarios from the National Highway Traffic Safety Administration for motion planning,<sup>[35]</sup> or WildDash for perception.<sup>[36]</sup>
- Random testing: This method refers to test cases created based on randomly generated parameters from the input domain. For automotive, vision-based noise can include Gaussian noise or artificial occlusions.
- Metamorphic testing:<sup>[37],[38]</sup> This method refers to defining metamorphic relations to transform one test case into another. Metamorphic relations characterize the relationship between the change of input and the change of output.

EXAMPLE 1 The metamorphic relation can describe that when switching from daytime to night-time while keeping the rest of the factors the same, the predicted bounding box size of the pedestrian is the same.

- K-way combinatorial testing<sup>[39]</sup> against pre-specified input space dimensions:<sup>[40],[41]</sup> This class of methods is introduced during the definition of test coverage. Dimensions of the input space are analysed from the input domain to identify equivalence classes in which a uniform behaviour of the system is expected. Then for any K dimensions, test all possible discrete combinations of those parameters on the AI model.

EXAMPLE 2 Consider a simple input domain is characterized using the following dimensions: time-of-day ∈ (daytime, evening, night, dawn), weather ∈ (fine, cloudy, rainy, snow, fog) and road-intersections ∈ (lane-diverging, lane-merging, straight). Two-way combinatorial testing needs to ensure that the set of collected test cases can cover, for any arbitrary two dimensions (e.g. <time-of-day, weather> or <time-of-day, road-intersections>), every pair of elements (e.g. <night, snow> or <dawn, lane-merging>) can be covered with a minimum amount of test cases (e.g. at least 1 test case). Using combinatorial testing, one can argue the relative completeness of testing efforts.

- Boundary value testing or corner case testing: This technique refers to testing boundary values of the input parameters. Typically, the parametrization of the boundaries of the input domain for complex tasks can only be accomplished approximately and with significant effort. Thus, testing AI systems at the boundaries of the input domain requires additional methods as compared to the non-AI case.

NOTE 1 The equivalence classes within the semantic input space do not necessarily correspond to the features within the syntactic input space learned by the AI model.

- Gradient-based search methods or other open-box optimization-based testing methods: This class of methods utilizes the knowledge of internal model parameters of the ML model to guide the generation of test cases. One can design optimization objectives such as creating erroneous prediction and, subsequently, utilize a given input's derivatives (gradient) to move towards the optimization objective. This is an open-box technique as the gradient is made available to the testing tool. Methods in adversarial perturbation (e.g. FGSM or PGD)<sup>[42],[43]</sup> utilize this principle.
- Genetic algorithms or other closed-box optimization-based testing methods: In contrast to open-box testing methods where the gradients are made available, closed-box optimization methods still try to perform the change of the input without utilizing the gradient information. Genetic algorithms use an initial population of test inputs and perform mutation of parameters to generate new test candidates. The next-generation of test cases are based on the pool of current-generation test cases that are best performers, with the definition of best characterized by the degree of violation. Other closed-box methods in the falsification of cyber-physical systems (e.g. simulated annealing or Bayesian optimization)<sup>[44],[45]</sup> use random samples to guess the gradient direction, in order to transform an input into another that can lead to undesired situations.
- Probabilistic sampling-based test methods:<sup>[46]</sup> These methods assume the availability of some prior belief about the distribution of AI errors within the input space. The areas with higher AI error distribution are then sampled more often with the aim of finding triggering conditions of AI errors more efficiently.
- Synthetic test case generation: This method allows generating edge cases that can be dangerous to be reproduced in physical world, as well as creating diversities in the scenarios being collected.
- Testing based on expert knowledge: Knowledge-driven testing refers to applying domain-specific know-how to create test cases, thereby checking if the model under analysis exhibits performance limitations that lead to safety concerns. For example, automotive perception component providers maintain a database of edge cases (e.g. jaywalking on a foggy night, pedestrian walking out of a pile of snow) that are considered challenging scenarios in products of the previous generation. These edge cases correspond to potential triggering conditions to the model under analysis.
- Tests that analyze resource limitations (e.g. runtime): This method covers activities such as ensuring that the model can be operated under the specified frequency (e.g. 10 FPS).
- Robustness testing: This method refers to considering the application's noise patterns or reasonable transformations and checking if the prediction subject to noise or transformation produces results consistent with the input without noise.

EXAMPLE 3 For an AI component working on audio data, robustness with respect to noise is of interest. In robustness testing, noisy audio signals are applied to the component, and it is verified whether the behaviour is as described in the requirements, e.g. that the performance does not fall below a certain threshold for noise up to a certain amplitude.

EXAMPLE 4 Autonomous vehicles rely on their perception capabilities to interact with the surrounding environment, which can be influenced by changes such as weather and lighting conditions. During robustness testing, it is important to consider adding perturbations to the test data across multiple dimensions simultaneously. For instance, in the case of image data captured by cameras, various types and intensities of weather conditions such as dusk and heavy rain can be introduced under different lighting conditions. This enables the evaluation of the system's accuracy in perceiving targets under challenging scenarios and assesses the model's adaptability to different combinations of disturbances. By subjecting the AI system to diverse and realistic environmental variations during testing, the effectiveness and robustness of the AI system in handling real-world conditions can be ensured.

- Tests based on model analysis/review: The performance on subsets of data is analysed to identify weak spots and/or lack of fairness. Subsets of data are identified where an AI component has weak performance. The errors are analysed to determine whether this is due to a systematic problem, e.g. low perception performance for bright frames because there were no bright frames in the training dataset. Furthermore, test for potential assumed weaknesses of the model architecture fall into the class of test methods. For example, if a DNN for object detection is sensitive to rotated objects due to its architecture, this may motivate tests with rotated objects.

Apart from testing, one can also apply methods that make use of formal verification. Methods making use of formal verification can be categorised into those that are sound and complete or just sound.

- Exact methods (sound and complete) via specialized constraint solvers (e.g. Reference [47]) or a reduction to mixed-integer linear programming or convex optimization.
- Sound methods based on methods such as abstract interpretation.[48][49] These types of methods can guarantee safety when the solver returns "safe". However, when the solver returns "unsafe" and provides a counterexample, the counterexample can be spurious due to over-approximating the state-space in the verification process.

NOTE 2 For formal verification of deep neural networks, beyond the issue of scalability, the lack of a precise specification and characterisation of the input space is one of the critical challenges in the application of formal verification approaches to AI-components with high-dimensional input spaces (such as images).[50] Therefore, some state-of-the-art approaches to image models, due to the inability to mathematically characterize the input space, restrict the use of formal verification to the evaluation of robustness against perturbations over selected test samples.

## 12.5.4 AI system integration and verification

Based on ISO 26262-6:2018, 10.2, software integration and verification refers to the activities where suitable integration levels and the interfaces between the software elements are verified according to the software architectural design. Moreover, ISO 26262-4:2018, Clause 7 discusses system and item integration and testing. ISO 21448:2022, Table 10 offers an additional list of methods for integrated-system verification. In principle, activities conducted in the ISO 26262 series and ISO 21448 can also be used to support the AI system integration and verification.

If two components that both contain AI components are integrated and if these AI components are not stochastically independent, test methods that address the statistical nature of the AI components are used to verify the statistical properties of the composition. This can be done using the methods listed in 12.5.3. Stochastic independence, in this context, informally means that the correctness of the output of the first AI component does not influence the probability for a correct output of the second component. Formally, two events are stochastically independent if the probability for the occurrence of both events is equal to the product of the probability of occurrence of the individual events.

EXAMPLE Two object detection components, one based on camera, and one based on LiDAR, are both affected by occlusion of objects. If the component using the camera input does not detect an occluded object, it is more likely that the component using the LiDAR input also does not detect it. Thus, the probability that both components do not detect a certain object is not the product of the probabilities that each component does not detect the object. They are not stochastically independent.

NOTE If statistical properties need to be verified after an integration step, statistical test methods can also be employed. In more detail, statistical properties of a component can be verified on component level. However, typically there are also statistical properties on higher integration levels or AI system level that are relevant. Often these cannot directly be taken over from component level because of the effect of other components. Thus, statistical verification can be deferred to higher integration levels or the system level.

## 12.5.5 Virtual testing vs physical testing

### 12.5.5.1 Virtual test platforms

ISO/IEC TR 5469 provides a detailed discussion on virtual testing and physical testing for functional safety in AI systems. In addition, it provides guidance on how to assess the virtual test platforms.

This subclause briefly describes challenges with the physical testing of AI systems and focuses on some of the advantages of virtual testing for AI systems in road vehicles.

NOTE 1 Requirements for using virtual test platforms in the context of validation of whole vehicle systems are described in ISO 34502:2022, 4.6.4.3.

For AI systems operating within complex environments, there are challenges in physically testing an adequate range of use case conditions (e.g. weather effects, behaviour of other road users, etc.), in particular

to achieve a sufficient coverage of edge cases. In these particular situations, virtual test platforms can be used to simulate all the desired variations.

NOTE 2 Edge cases can be defined as scenarios with very specific and rare conditions, e.g. extreme weather conditions, sun glare, specific environment conditions (erased road marking).

Another capability of virtual test platforms is the synthetic generation of datasets used for AI system development. Synthetic data generation allows for the generation of synthetic ground-truth information automatically from virtual test platforms. Indeed, within a simulation environment, not only the virtual image as seen by the sensor is simulated, but as an omniscient environment, simulation frameworks can also generate additional information (also known as ground-truth data) such as depth, object segmentation, object materials, bounding boxes or optical flow. Those data are important in the evaluation of performances of the AI models as it provides data as seen from a perfect sensor. Ground-truth data are challenging to obtain from real-world observations and require in most cases a complicated, error-prone and time-consuming manual process. Additionally, generating synthetic datasets allows for the training, validation and testing of AI models with independent data, which is key to avoid coincidental correlations, as stated in [11.4.3.2](#).

NOTE 3 Independence between synthetically generated data is a complex subject and can be achieved by coverage analysis of the different environment parameters (daytime, weather conditions, sensors positioning, etc.), by using different sensor types among the available in the catalogues, by using different ground-truth generators, by variation of different simulation parameters (e.g. scenario length), by following independent processes when generating the synthetic data (e.g. relying upon different methods or stakeholders to accomplish process tasks), etc.

As the dataset generated for AI training is produced in a virtual environment, the entire generation workflow is validated, and the correlation with real data is determined. To make sure we can rely on virtual datasets, comparisons between virtual and physical datasets is performed. Real world conditions, for instance, adverse weather, can be reproduced, tuned and measured precisely in the laboratory (for example, see Reference [\[51\]](#)). Once such data are collected, the same scenario and conditions can be set in the virtual environment for advanced correlation. It also helps understanding the gaps and domain of validity.

Moreover, the use of physics-based solvers on different disciplines (optics, electromagnetic, thermal, etc.) can be validated by independent certification bodies. This means being able to handle physics-based data as an input (e.g. materials, emitters, sources), implement laws of physics within the solver but also generate the outputs that imitates the real conditions (e.g. spectral images, point clouds, range-doppler).

#### 12.5.5.2 Using HiL for synthetic data validation

Hardware-in-the-Loop (HiL) testing can be also used to test the accuracy of synthetic datasets; it is one of the most accurate ways to correlate between real and virtual datasets. The sensors can be used to record scenes from the real world and thus, by reproducing and injecting the exact same scene from the virtual world into the real sensor, comparison between the real and the virtual dataset can be performed (by comparing the results of the HiL testing) as the only changing parameter is the dataset (sensor receivers and post-processing unit used in the real and virtual cases are the same).

#### 12.5.5.3 Virtual testing

For AI systems operating within complex environments, there are challenges in physically testing an adequate range of use case conditions (e.g. weather effects, behaviour of other road users), and in particular to achieve a sufficient coverage of edge cases. In these particular situations, virtual test platforms can be used to simulate all the desired variations.

NOTE 1 Edge cases can be defined as scenarios with very specific and rare conditions, for example extreme weather conditions, sun glare, specific environment conditions (erased road marking), etc.

#### 12.5.6 Evaluation of the safety-related performance of the AI system

The performance of an AI system refers to the level of precision for prediction, the accuracy for classification, and the efficiency of an algorithm. Evaluation of the performance of an AI system is typically carried out with comparison of the system's output to the output from a benchmark, using a dataset which is proposed as, or has become, a standard dataset by which different solutions are evaluated. However, the performance of the

AI system is brittle in the sense that the AI system that performs well has generally either been tailored to solve particular problems, or trained on specific set of data relating to the problems in a particular domain. Therefore, the lack of universally accepted or formalized criteria to assess the safety-related performance of the AI system poses an additional hurdle for wide adoption or utilization of an AI technology. This subclause aims to provide some guidance that fits into the framework of evaluation process of the safety-related performance of the AI system.

NOTE 1 AI models, particularly ML-based models, manifest the characteristics of statistical models. Traditional functional safety requires the system to be predictable, and hence the AI model's behaviour can be predictable in a probabilistic sense. Predictability does not equate determinism, as it does in traditional software development. This implies that the AI system can contribute to a failure which is caused by software itself.

The following are examples of common causes of the negative impacts on safety-related performance due to AI errors:

- Inadequacy and uncertainty in the learning process: The learning process is instrumental for any ML-based system to generate accurate and reliable outputs. Insufficiencies in the training and test data, dynamically changing environments and unpredictable intentions of road users etc. can lead to unreliable learning results or misinterpretations by the AI system.
- Inappropriate cost function selection: The cost function is either less representative, or not affordable to re-evaluate in a consistent manner. This results in negative side effects or reward hacking.
- Inappropriate metrics: using metrics that do not suitably match the actual goals and priorities obscures the general system performance.
- Inconsistency between trained AI model and deployed AI model.
- Lack of benchmarks: Lack of universally adopted benchmarks which are reliable, transparent, standard and vendor-neutral results in performance differences between different parameters, even within the same application domain.

The applications of ML-based AI systems are usually categorized as follows:

- regression, where the task is predicting a continuous quantity;
- classification, where the task is predicting a discrete class label.

Clearly and unambiguously defined metrics are required to evaluate the performance of an AI system, which in turn implies the safety level of an AI system (e.g. using performance indicator as a pass/fail criterion of the system). A summary of the widely-adopted performance metrics for both categories are listed in [Annex H](#). The metrics included in [Annex H](#) are by no means an exhaustive list of what the industry is currently using. Other safety-related metrics can be derived, based on particular use cases and domain experts' knowledge and judgement, to evaluate the AI system from specific aspects of the system requirements.

NOTE 2 The performance metrics in [Annex H](#) are different from the loss functions. Loss functions are measures to quantify the model's performance during training process, while metrics are used to monitor and evaluate the performance of trained models in testing phase.

### 12.5.7 AI system safety validation

In contrast to verification, AI system safety validation refers to checking if the safety requirements allocated to the AI system (from the encompassing system) are met after the AI system is integrated into the encompassing system. AI system safety validation activities are usually done by the system integrator (e.g. original equipment manufacturer), where the validation target is defined separately. The AI system developer may need to support the activity.

For AI system safety validation, the individual methods listed in [12.5.3](#) can also be used. However, the focus is on the systematic exploration of all relevant scenarios within the input space and to examine abnormal

situations. Systematic random testing by first discretizing the input space is an example of such methods to argue relative completeness<sup>[52][53]</sup>.

NOTE If the verification of the encompassing system admits the usage of virtual techniques like simulation, then the safety validation of the AI system, once integrated into the encompassing system, can also be based upon virtual techniques, for instance simulation can be conducted to systematically explore relevant scenarios and identify corner cases or abnormal situations.

For DNNs, AI system safety validation using field testing (e.g. by operating a fleet of autonomous driving vehicles) can be done with the assistance of active learning methods or other methods for detecting out-of-distribution data. The underlying idea is that active learning methods try to infer if an input can be included in the training dataset by considering how different this input is with all existing training data.

As explained in [Clause 9](#), in addition to the standard SOTIF-related AI safety requirements that directly address the performance targets, the workflow can introduce additional AI safety requirements that concretize AI safety-related properties (e.g. robustness, interpretability). Safety validation of the AI system also considers validating the appropriateness of these additionally introduced requirements. The purpose of the validation is to ensure that no insufficiencies of the specification exist. In particular, the activity ensures that the quantitative thresholds being set in the requirements are appropriate (via methods such as statistical hypothesis testing) and can positively impact safety (by positively influencing the safety-related properties). Since these requirements are not exposed to the system integrators (OEM), validating these requirements is the task of the AI component/system provider.

## 12.6 Work products

**12.6.1 AI system verification report**, resulting from [12.3.1](#) to [12.3.5](#) and [12.3.7](#).

**12.6.2 Integrated AI system**, resulting from [12.3.6](#).

**12.6.3 AI system validation report**, resulting from [12.3.8](#).

## 13 Safety analysis of AI systems

### 13.1 Objectives

The objectives of this clause are:

- a) to identify safety-related faults and AI errors that can lead to the violation of AI safety requirements;
- b) to identify their potential causes;
- c) to support the definition of safety measures to prevent or control safety-related AI errors;

NOTE 1 These measures can include improving AI design, AI methods, dataset generation, updating AI safety requirements and related AI system development processes.

- d) to support the verification of AI safety requirements, through modification or identification of new AI safety requirements on data specifications and collection, design specifications, and test specifications.

NOTE 2 The objectives, scope, and level of granularity of the safety analysis can depend on the phases of AI safety life cycle.

NOTE 3 Safety analysis of the AI system complements the safety analysis in accordance with the ISO 26262 series and ISO 21448.

NOTE 4 Safety analysis of an AI system can be performed within the safety analysis of the encompassing system, e.g. an item or a vehicle.

NOTE 5 DFA is an important activity that follows the safety analysis of an AI system. DFA of AI systems is a new area of research and is not covered in this document. The reader can refer to ISO 26262-9:2018, Clause 7 for guidance on DFA, which can be applicable to AI systems.

This clause aims to provide confidence that the risk of violation of the AI safety requirement at the AI system level due to AI errors is sufficiently low, i.e. within the acceptable residual risk.

### 13.2 Prerequisites and supporting information

The following information shall be available at the initiation of the safety analysis activity:

- a) AI safety requirements, [Clause 9](#);
- b) input space definition (refined), [Clause 9](#);
- c) known insufficiencies of the AI system and the corresponding subdomains of the input space, [Clause 9](#);
- d) AI component or AI system architecture (refined), [Clause 10](#);
- e) dataset requirements specification, [Clause 11](#);
- f) dataset design specification, [Clause 11](#);
- g) dataset verification report, [Clause 11](#);
- h) dataset validation report, [Clause 11](#);
- i) dataset safety analysis report, [Clause 11](#);
- j) AI system verification report, [Clause 12](#);
- k) AI system validation report, [Clause 12](#).

NOTE 1 The AI component or AI system architecture (refined) can be used to determine the boundaries of the safety analysis.

NOTE 2 Safety analysis can be performed at different phases of an AI safety lifecycle. Therefore, during early phases of an AI system development, availability of the prerequisites can be limited.

NOTE 3 Safety analysis can be performed at different levels of integration, e.g. AI system, AI components, or AI models, or with different focus, e.g. architectural aspects, data aspects, or combination thereof.

### 13.3 General requirements

**13.3.1** Safety analysis techniques suitable for identifying the safety-related AI errors of the AI models in AI systems shall be applied.

**13.3.2** Safety analysis of the AI system shall identify the AI errors of the AI system and its components that have the potential to violate one or more AI safety requirements.

**13.3.3** Safety analysis shall identify the safety-related faults, potential functional insufficiencies and their potential underlying issues of the identified safety-related AI errors, if one or more AI safety requirements are violated due to the identified AI errors.

**13.3.4** Safety analysis results shall be used to identify prevention or mitigation measures to address the causes of the AI errors that are potentially violating one or more AI safety requirements.

**13.3.5** Safety analysis results shall be used to verify the completeness of the AI safety requirements.

## 13.4 Safety analysis of the AI system

### 13.4.1 Scope of the AI safety analysis

Safety analysis of AI systems includes a systematic identification of AI errors in an AI system and, in particular, functional insufficiencies and safety-related faults that can lead to the violation of an AI safety requirement. These AI errors can be related to:

- an AI component consisting of an AI model;
- an AI component not consisting of an AI model.

NOTE Safety-related faults and functional insufficiencies related to an AI component can be originated in that AI component or be due to the interaction of the AI component with other components within the AI system or outside of the AI system.

When AI models are in the scope of the safety analysis, safety analysis addresses the safety-related faults and functional insufficiencies of the AI models, their causes and their impact on vehicle behaviour. Safety analysis starts early during development. [Figure 13-1](#) shows a flowchart of a top-down safety analysis approach in an AI system as an example. In this example, safety analysis is started upon the observation of an undesired safety related behaviour at the vehicle level. However, in general, safety analysis can start from the level required determined by the team performing the analysis. If safety-related faults and functional insufficiencies are related to an AI component that does not contain an AI model, safety analysis can be performed following the requirements and recommendations of ISO 26262-9:2018, Clause 8 and ISO 21448. During the development of an AI system, other safety analysis methods, for example, bottom-up approaches can also be used to identify the faults and potential insufficiencies which might lead to AI errors.

In general, AI errors of an AI model are either due to issues in data specification and collection or issues in design and implementation or issues in requirement specification. Safety analysis to identify issues in data specification and collection and to define safety measures for prevention or control of safety-related issues is discussed in [11.4.3](#).

Safety analysis at the design phase identifies design-related issues which can contribute to AI errors violating an AI safety requirement. Safety analysis at the design phase is discussed in [Clause 10](#). If safety analysis identifies insufficiencies in an AI safety requirement specification, the requirement specification may need to be modified or one or more new requirements may need to be added. Requirement modification/addition follows the guidance provided for requirement modification/addition in [Clause 9](#).



### 13.4.2 Safety analysis based on the results of testing

If testing of an AI system, at any level, reveals the presence of AI errors, the results of the safety analysis are used to evaluate the impact of the AI errors on the conformity of the system under test to its AI safety requirements, to identify the causes of the safety-related AI errors and to define mitigation measures. Here, safety analysis activities consist of risk evaluation, root-cause analysis and risk mitigation as shown in the AI component testing workflow in [Figure 12-3](#).

- a) Risk evaluation: During this activity, the risk due to the failed test is evaluated to estimate the impact on safety. In general, if any of the AI safety requirements are violated due to an AI error, it can be concluded that safety is not achieved.

NOTE 1 A test can also fail due to the violation of non-safety requirements. In case multiple non-safety requirements are violated as the result of the failed test, the risk due the violated requirements is evaluated to assess the impact on safety.

- b) Root-cause analysis: In this step, underlying issues for the AI error(s) are identified. Issues in an AI model, in general, can be related to many areas including:
- AI safety requirements allocated to the AI component consisting of the AI model;
  - AI data including datasets;
  - AI model design.

Once a potential category of causes has been identified, a more detailed safety analysis related to that area can be performed to evaluate the cause in more detail. For this, the results of safety analysis performed during AI model design or dataset generation can be used. For safety analysis on datasets, see [Clause 11](#), and for safety analysis during AI model design, see [Clause 10](#).

- c) Risk mitigation: When the root causes of the issues are identified, prevention, detection and/or control measures regarding the identified root causes need to be defined. These risk mitigation measures include:
- modification/addition/removal of AI safety requirements;

NOTE 2 For supervised machine learning, the activities in root-cause analysis and the proposal of modification/addition/removal of AI safety requirements are further detailed in [9.5.3](#).

- changes in the AI model;

EXAMPLE Testing of an ML model developed for detecting pedestrians in an autonomous driving system reveals that the ML model does not detect pedestrians which are standing next to a traffic post. Safety analysis shows that the hyperparameters of the neural network are not selected optimally. The hyperparameters of the AI model are changed to mitigate the issue.

- changes in the dataset;
- modification of the AI development processes.

Mitigation measures are discussed in more detail in [Clause 9](#), [Clause 10](#), and [Clause 11](#). These risk mitigation measures then need to be implemented as part of the AI system development including requirement derivation, design and dataset creation according to [Clause 9](#), [Clause 10](#) and [Clause 11](#), respectively. This activity might require the creation of additional safety-related test cases.

### 13.4.3 Safety analysis techniques

Safety analysis techniques should provide adequate identification of hazards and their potential causes. The sufficiency of a safety analysis technique to model a system is argued by the following methods:

- proven-in-use-argumentation;

- critical review of the chosen technique, where pros and cons of the technique for safety analysis of the system is evaluated and its limitations are identified.

Since safety of AI systems is a relatively new topic, the proven-in-use-argumentation is challenging to apply. Safe application of AI is challenging, because AI introduces new classes of mechanisms for how risks can emerge and safety concerns. The concerns include inclusion of training data instead of system specifications, no clear design as system architecture, uncertainties and the explainability challenges in the models' outputs. Some of the salient features of AI systems that can impact the safety analysis are:

- AI systems can behave nonlinearly. Depending on their current state and context, they might react to the same inputs very differently. Additionally, smaller disturbances in the input can produce irregular outputs.
- In some cases, the environment that the AI system is deployed in is ever evolving. For example, in cases of highly-automated driving vehicles operating in the open context, new traffic participants can appear over the course of their operations.
- AI systems can produce complex interactions within its elements and with the environment. The models that result for such systems might introduce complex correlations.

Safety analysis techniques analyse the systems and its underlying assumptions. Some of the commonly used safety analysis methods are shown in [Table 13-1](#). These salient features of AI systems require a thorough understanding of the safety analysis method selected to analyse these systems. Some of the existing analysis techniques have been enhanced to model the AI systems,<sup>[54],[55]</sup> while newer modelling techniques have been introduced with stronger assumptions to model the AI systems.<sup>[56],[57],[58]</sup>

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

Table 13-1 — Safety analysis techniques

Safety analysis	Modelling assumptions	Advantages	Limitations
Fault tree analysis [54]	<ul style="list-style-type: none"> <li>— Independence of events</li> <li>— Bernoulli model</li> <li>— Simplified causal relation</li> <li>— Static temporal concept</li> </ul>	Based on Boolean algebraic concepts	Generally static
Failure mode and effects analysis [59]	<ul style="list-style-type: none"> <li>— Single point of failure</li> <li>— Simplified causal relation</li> <li>— Static temporal concept</li> </ul>	<ul style="list-style-type: none"> <li>— Documented process</li> <li>— Early design decision</li> <li>— Easy to implement</li> </ul>	<ul style="list-style-type: none"> <li>— Inability to determine complex failure mode</li> <li>— Cause-and-effect chains (might not be predictable)</li> </ul>
System theoretic process analysis [60]	<ul style="list-style-type: none"> <li>— Simplified causal relation</li> <li>— System fails in a certain pattern</li> </ul>	<ul style="list-style-type: none"> <li>— Models' interactions</li> <li>— Implementation is more comprehensive</li> </ul>	Limited keyword set
Event tree analysis [61]	Single point of initiations	Assessment of multiple faults and failure	Probability identification is difficult
Bayesian network/ causal Bayesian network [57]	<ul style="list-style-type: none"> <li>— Model is the best representation</li> <li>— Probability distributions are known</li> </ul>	<ul style="list-style-type: none"> <li>— Can model complex relations</li> <li>— Hybrid modelling is possible</li> <li>— Models multiple point initiations and failure</li> <li>— Can handle conditional independence concepts</li> </ul>	<ul style="list-style-type: none"> <li>— Difficult to model</li> <li>— Probability identification is difficult</li> </ul>
HAZOP [62]	Single point of initiations	<ul style="list-style-type: none"> <li>— Documented process</li> <li>— Early design decision</li> <li>— Easy to implement</li> </ul>	<ul style="list-style-type: none"> <li>— Single point of failure</li> <li>— Propagation of failures is not clear</li> <li>— Identification of causes is weak</li> </ul>

## 13.5 Work products

13.5.1 Safety analysis report, resulting from 13.3.1 to 13.3.5.

## 14 Measures during operation

### 14.1 Objectives

The objectives of this clause are:

- a) to define the process requirements to continuously assure AI safety after deployment;
- b) to use the measures defined in Clause 8 and/or additional measures for the identification of safety risk associated with the AI system during operation and measures to maintain AI safety during operation;
- c) to ensure responses are in place to address unacceptable safety risks associated with the AI system and ensure re-approval of the modified AI system before release.

## 14.2 Prerequisites and supporting information

The following information shall be available at the initiation of this phase:

- a) AI system definition (from external sources, e.g. the encompassing system development), including:
  - 1) interfaces with the encompassing system;
  - 2) assumptions on the use of the AI system;
- b) field data collected by the encompassing system;
- c) safety requirements on the AI system, [Clause 9](#);
- d) AI component or AI system architecture, [Clause 10](#);
- e) dataset requirements specification, [Clause 11](#);
- f) dataset design specification, [Clause 11](#);
- g) dataset maintenance plan, [Clause 11](#);
- h) known insufficiencies of the AI system and the corresponding subdomains of the input space, [Clause 9](#)
- i) results of verification and validation activities including known functional insufficiencies of the AI system (if available), [Clause 12](#);
- j) safety assurance argument, [Clause 8](#).

## 14.3 General requirements

**14.3.1** The process and its activities necessary to assure the AI safety and the validity of the assurance argument during operation shall be specified.

NOTE 1 This process can include the procedure to terminate the safety support and properly notify the user of AI system regarding this termination.

NOTE 2 These activities include the identification of safety issues of the AI system during operation and their resolution procedure.

**14.3.2** The on-board and off-board measures necessary to execute the specified activities in [14.3.1](#) shall be developed and implemented.

EXAMPLE Measures can include monitoring the operational status of the AI system, detecting safety-related errors, etc.

**14.3.3** The identified safety-related field events shall be evaluated and, if the risk is deemed unacceptable, countermeasures shall be taken to mitigate the risk.

**14.3.4** The effectiveness of the countermeasures shall be evaluated after their application during the operation phase. The countermeasures shall be modified if the residual risk is still unacceptable.

**14.3.5** The specified maintenance activities during operation shall be executed in order to continuously keep AI safety to a reasonable level.

EXAMPLE Field data collection, AI re-training, re-validation and re-approval, etc. can be executed in order to continuously maintain the AI safety.

## 14.4 Planning for operation and continuous assurance

### 14.4.1 Safety risk of the AI system during operation phase

Upon achieving recommendation for release, the residual risk is evaluated to be acceptable based on the evidence and assumptions generated during the development phase. However, post-deployment field risk evaluation can detect an elevated risk associated with the AI system due to hazards resulting from:

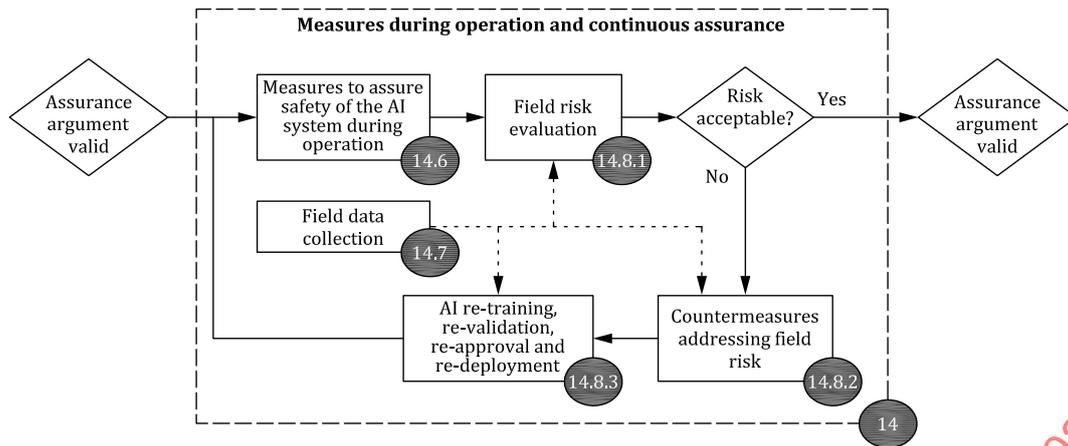
- a) development uncertainties, for example:
  - incorrect estimation of residual risk;
  - previously unknown hazardous functional insufficiencies;
  - incorrect estimation of the occurrence of AI related faults occurred during operation.
- b) incorrect unexpected operation-specific activities, for example:
  - during maintenance, e.g. retrofit camera or radar without recalibration or poor tolerance;
  - during update of the AI system, or update of external systems that interacted with the AI system, e.g. outdated software version, out of sync of component updates.
- c) changes in the operation environment, for example:
  - new traffic rules;
  - new traffic facilities;
  - new type of traffic participants;
  - changes of assumptions on the operating conditions.

NOTE Changes in the operation environment can introduce OOD samples and potentially cause OOD errors in the AI system.

Specific activities can be necessary during the operation phase to address these risks and assure a continuous level of AI safety.

### 14.4.2 Safety activities during the operation phase

[Figure 14-1](#) illustrates the flow of the activities of this clause to assure the safety of AI during the operation phase.



Key

 (sub)clause number(s)

 process flow

 relation

Figure 14-1 — Safety assurance during operation phase<sup>2)</sup>

Safety assurance during operation starts with applying measures to assure AI safety in the operation phase (14.6). Technical measures are applied to monitor the behaviour of AI systems during operation, and if anomalies or undesired safety-related behaviour at the vehicle level are detected, mitigation measures are taken (14.6.2). Some additional measures are introduced to address misuse related risk, e.g. user guidance (14.6.3). Field data will be collected during monitoring (14.7), in order to support afterwards risk evaluation, mitigation and AI system modification.

Risks identified in the field are evaluated (14.8.1). If the level of identified risk is acceptable, then no further activities are applied. Alternatively, if the risks are found to be unacceptable, countermeasures are determined based on the evaluation results (14.8.2). The AI system can go through re-training, re-validation and re-approval process, if necessary (14.8.3).

14.5 Continual, periodic re-evaluation of the assurance argument

Due to the complexity of the functionality to be implemented, the environment in which it is deployed, as well as the nature of the AI technologies themselves, some uncertainty in the assurance argument can remain.

This leads to a residual risk that the safety requirements allocated to the AI system are violated during operation. This residual risk can be related to previously unknown triggering conditions of residual insufficiencies in the AI system, inadequacies of the assurance argument or to changes within the operating context. A continual, periodic re-evaluation of the assurance argument can offset this emergent risk.

The residual risk can be offset by operational restrictions until sufficient evidence can be collected to increase confidence in the assurance argument. Examples of operational restrictions that can be applied include:

- restricting the set of operational conditions, thus reducing the risk of violating assumptions in the assurance argument;
- limiting the functionality of the AI system, thus reducing the severity of residual errors.

2) Safety assurance activities during operation phase rely on previous clauses, i.e. 14.6 refers to Clause 10, 14.7 refers to Clause 11, 14.8 refers to Clause 8, Clause 9, Clause 10, Clause 12 and Clause 13.

The re-evaluation of the assurance argument can be performed based on the criteria outlined in [8.7](#). In particular, evidence collected during operation can be re-used to provide additional support for claims of the assurance argument, as well as to identify potential defeaters to these claims.

**EXAMPLE** An assurance argument uses a set of assumptions on the input space to define the context of the assurance argument. A run-time anomaly detection identifies OOD inputs that were not considered within the set of assumptions. In reaction to this new information, the assurance argument is re-evaluated using a wider set of assumptions, which includes the identified OOD inputs.

Triggers for the re-evaluation of the assurance argument can include:

- periodic review;
- collection of observations that can be used as additional evidence in the assurance argument;
- analysis of reported field incidents;
- results of on-board and off-board monitoring;
- change in operational parameters or environment conditions;
- modification or maintenance of the encompassing system;
- changes in operating procedures.

## 14.6 Measures to assure safety of the AI system during operation

### 14.6.1 General

The intention of this subclause is to give guidance for applying measures to assure AI safety during operation and trigger potential updates to the AI system. This subclause also provides additional non-technical measures, for example user involvement to prevent misuse and assure the safe operation of AI, if possible.

### 14.6.2 Technical safety measures

Monitoring, detection and mitigation, which are used to evaluate the behaviour of AI systems against errors or insufficiencies, are measures used to assure safety of the AI system during operation. These measures can rely on on-board mechanisms and/or off-board mechanisms (e.g. cloud monitoring).

**NOTE 1** The architectural measures to assure safety of AI during the operation phase are defined in [Clause 10](#).

**NOTE 2** In contrast with the on-board mechanisms which are used to detect and mitigate abnormal behaviours of AI systems and the vehicle equipped with AI, the off-board mechanisms can detect abnormal behaviours with higher accuracy due to, for example, larger computing power and a more precise model.

**NOTE 3** The off-board measures (e.g. cloud monitoring) can also be used to monitor the general behaviours of all vehicles equipped with the AI systems. These off-board measures can then be used to support the evaluation of the overall risk after deployment of the AI systems.

Regarding monitoring and detection, the following events, including related context, can be reported with the purpose of finding insufficiencies or errors of AI system, if applicable:

- a) Input space related events:
  - detection of OOD data and data distributional shift;
  - detection of exiting from operating context;
  - detection of concept drift or changes in features (e.g. new objects, different behaviours, new or changed rules);
- b) Model behaviour related events:

- detection of abnormal behaviour;

NOTE 4 Some abnormal behaviours can be caused by rare input conditions and these behaviours can be evaluated as safe after detection. For example, while reversing, the AI model stopped at a shorter distance than specified from a parked car. Behaviour is detected by the system and logged. After analysis of the report, the behaviour is considered safe and an update is not necessary.

c) Output related events:

- detection of abnormal output;

NOTE 5 Exercise caution when implementing plausibility checks as these can lead to missed objects and safety concerns (e.g. rejecting humans taller than 7 ft can lead to mis-detection of a pedestrian carrying a flagpole or on stilts or having a child on their shoulders).

- detection of output bias;
- outputs with low confidence level;

d) Incidents/accidents analysis:

- incidents/accidents where the AI system was directly or indirectly involved are analysed to support the improvement of the AI system.

NOTE 6 For detected errors or insufficiencies of the AI system which can lead to a hazard, the risk can be mitigated by measures within the AI system or the encompassing system. For example, switching to non-AI system or executing a manoeuvre that results in a minimal risk condition.

As the insufficiencies of the AI system can influence the behaviour of the encompassing vehicle system, any abnormal behaviours or emergency events of the vehicle may also imply or influence on insufficiencies of the AI system, for example:

- function degradation;
- take-over request;
- emergency manoeuvre;
- transition to a minimal risk condition;
- collision or near-collision event;
- contradiction between AI system and non-AI system.

EXAMPLE The AI system and non-AI systems, which are both used for decision making, may provide diametrically opposed results under an unprotected left-turn scenario, for instance, one for “yield” and the other for “not yield”.

Besides triggering modification activities, errors or insufficiencies of the AI system identified during operation may indicate the weaknesses in the development and safety assurance process, architectural measures or incorrectness of their usage assumptions, thus modification of these measures may be needed.

### 14.6.3 Safe operation guidance and misuse prevention in the field

The user of the AI system can lack understanding of its capabilities which results in misuse due to overconfidence in the AI system. To prevent overconfidence, users are made aware of the limitations of AI systems via, for example, user training or relevant information through the human-machine interface.

EXAMPLE The user is trained to correctly use the AI system and be informed of scenarios in which the AI system is intended for use, considering the performance limitations of the AI system within these scenarios.

Another possible prevention of misuse are technical measures (e.g. warning or degradation or disablement of services) that are triggered when the AI system is misused by the user during operation.

## 14.7 Field data collection

The intention of this subclause is to introduce field data collection as a supplementary data source for AI system maintenance to improve dataset integrity, distribution and usage (see [Clause 11](#)).

NOTE 1 Field data collection is related to AI systems whose safety can be affected by field conditions. An autonomous driving system that makes use of AI technologies is a typical case and selected as example in this subclause.

The motivation to collect field data during operation includes, for example, addressing environment changes which may affect the behaviour of AI system, identifying and removing residual insufficiencies and collecting additional training data. The quality of the collected field data needs to be ensured and the data needs to be transmitted to relevant parties (e.g. manufacturers, suppliers and/or regulators) for use to support the update of the AI system if necessary. The following topics can be considered when collecting field data:

- a) competence management: to ensure the efficiency and quality of the field data collection, competence management measures can be applied to people responsible for field monitoring, data collection or data analysis. All systems involved in field monitoring activities are tested, validated and released to ensure required reliability level.
- b) data characteristics: the data characteristics of the field data can be defined depending on the planned usage of the data.

EXAMPLE 1 In some cases, a large number of images of a high resolution are needed in order to improve the 2D image perception performance of the AI system, such as classifying a certain type of traffic sign. In other cases, the AI-based image processing algorithm can depend on relationships between sequences of images over time. For such cases, a minimum length of video sequences along with other associated sensor data and the results of the current iteration of the AI system are required to improve performance.

When analysing the field data, the following data characteristics can be considered:

NOTE 2 The data characteristics given below are not exhaustive.

- data categories;

EXAMPLE 2 Data source (radar, LiDAR, camera or HD map).

- data content;

EXAMPLE 3 Vehicle identification number (VIN), images from front camera, parsing data from front perception, changes of control mode, received remote control command, operation status and HMI data.

- data format;

EXAMPLE 4 JPEG, PNG and BIN.

- data size.

- c) Data collection trigger and transfer: To ensure that the collected data is sufficient to identify, analyse and improve safety-related issues, clear data collection triggering criteria are defined including the triggering conditions, triggering interval, start time and end time, triggering priority according to different cases.

EXAMPLE 5 The triggering rules of field data collection can be:

- accident or incident: collision event involving the automatic driving vehicle equipped with AI system;
- functional termination: autonomous function failure/insufficiencies, terminated by the human taking over;
- exiting operational design domain (ODD): the specific objects are detected, such as the red light, stop signs, etc. which are not within the ODD scope for a highway pilot feature, the value reported by rain sensor exceeds the threshold for a feature designed for no or small rain weather;
- implausible events: the distance or speed jitter of the detected object exceeds a certain value, the distance deviation of the detected object is greater than a certain value measured by different sensors;

- other functional insufficiencies: the target motion predicted by the AI algorithm is inconsistent with the actual situation;
- diverging decisions of redundant diverse AI elements or between AI and non-AI based elements.

NOTE 3 Triggering rules can be updated over-the-air (OTA) in order to collect different kinds of data. To collect sufficient data for each event, a timing buffer can be considered, for example: recording starts from at least X s before the event to at least Y s after the event.

When transferring the data, the conditions that may affect the reliability of the transfer are considered, for example: the data transfer is interrupted by loss of power and can therefore lead to loss of data.

- d) data storage: to ensure the integrity of data storage, safety mechanisms are implemented where reasonably practicable, for example, adding data integrity protection. The operation conditions that may influence the data storage are also considered.

NOTE 4 The general data storage requirements used for AI data collection can also be used for field data storage.

- e) configuration information: to ensure correctness of data collection, the configuration information about the field data to be collected is specified, which may include access rights, tools and repositories, and aligned with the requirements for datasets (see [11.3](#)).

## 14.8 Evaluation and continuous development

### 14.8.1 Field risk evaluation

Based on field measures ([14.6](#)) and data collected ([14.7](#)), the accidents, anomalies and undesired safety-related behaviour at the vehicle level potentially related to AI systems can be manually or automatically reported to the manufacturers or service providers. The number of reported issues can be large during the early phase after deployment. To solve the reported issues efficiently and economically, the manufacturers or service providers investigate the causes and evaluate the field risk of the issues, to determine the proper reactive actions to be taken, such as recall or OTA update.

The field risk evaluation is different compared to the hazard risk evaluation during the development phase. In particular, field risk evaluation is based on the real consequence of issues occurred during the operation instead of assumptions or estimations made at development phase.

To objectively evaluate the effects of the issues, the probability of occurrence, the severity and the effectiveness of countermeasures addressing the risk of the existing issues can be considered. This is similar to the occurrence, severity and detection parameters used by FMEA method for a systematic evaluation of risk.

NOTE 1 Alternative risk evaluation methods to FMEA based on a systematic methodology and predefined criteria can also be applied.

Field risk evaluation can be based on the following factors:

- a) Evaluation of the probability of occurrence: as described in [Figure 6-12](#), safety-related issues of the AI system can be caused by random hardware faults and/or systematic factors (e.g. systematic faults or functional insufficiencies).
- For issues associated with random hardware faults, the occurrence considered is determined by the failure rate and the probability of exposure to a hazardous scenario which has been considered by the ISO 26262 series.
  - For issues associated with systematic faults or functional insufficiencies, the risk is mainly determined by the probability of the exposure to the critical situations or probability of triggering events.
  - As the quantity and location of the vehicles can be known at this phase, it is possible to provide the occurrence with higher accuracy than during development.

EXAMPLE The occurrence rate of the issue over a given time period can be predicted based on the failure rate of the component, the quantity of vehicles in the field and the probability of the vehicles facing the hazardous scenarios.

- b) Evaluation of the severity: the severity evaluation can be based on the method defined by ISO 26262-3, which recommends the abbreviated injury scale (AIS) ranking method.

NOTE 2 In addition, other issues that can cause loss due to cybersecurity risk, violation of traffic rules or serious customer complaints can also be considered.

- c) Evaluation of the detection and mitigation measures: the measures in [14.6](#) can help to detect and mitigate the risk of errors in the AI system. The potential controllability by the driver can also be considered as a mitigation of the issue, if the field data shows relevant evidence.

It is possible to give risk evaluations based on the factors a), b) and c) in a qualitative way or quantitative way (if rates are defined for each factor). The evaluation results will support the identification of the response actions to be taken.

NOTE 3 For serious accidents (e.g. fatalities), even if the occurrence is rated as low or detection as high based on predefined criteria, the rating criteria can be adjusted and risk can be considered differently.

### 14.8.2 Countermeasures addressing field risk

The safety development of AI systems does not end after the encompassing system release for operation. The field risks can be higher than expected in case the on-board measures cannot detect and mitigate all risks. If hazardous events occur in the field, the following additional countermeasures can be taken:

- issue investigation actions to determine the causes of risk, e.g. scenario reconstruction based on the data collected, especially for AI-related incidents or accidents;
- risk evaluation as introduced in [14.8.1](#);
- restrictions on context of use or functionality deactivation or replacement;
- update of the AI system, for example an OTA update, when unacceptable systematic faults or insufficiencies are identified;
- customer notification, which can be taken together with the restrictions on context of use or AI system update actions, or dedicated notification to address misuse risk, e.g. emphasizing the operation requirements to the passengers by placing a warning card in the robotaxi.

NOTE Depending on the urgency of identified field risks, immediate actions or long-term actions can be taken based on risk evaluation.

An appropriate issue management process is important to ensure the effectiveness of countermeasures, including incidents or accident reporting, issue investigation, risk evaluation and countermeasure management processes.

The effectiveness of the countermeasures taken are monitored and evaluated after implementation and adjusted, if the risk is still unreasonable.

### 14.8.3 AI re-training, re-validation, re-approval and re-deployment

The system can be incrementally developed on the basis of collected field data and countermeasures to compensate for the identified risks, making the system safer and more robust. The intention of this subclause is to introduce AI model re-training, re-validation, re-approval and re-deployment.

NOTE AI system update and re-approval involves the activities described in [Clause 7](#) to [Clause 13](#), if relevant.

- a) re-training

During operation, valuable field data can be collected. Together with the data from the previously trained model, this data can be used to re-train the new model with the expectation of better performance. Re-training can be achieved by fine-tuning the pre-trained model or by training from scratch.

- fine-tune: this refers to small adjustments to model parameters. The newly acquired field data can be used to fine-tune the released model. When fine-tuning, a small learning rate is used so as not to over-distort the existing model.

**EXAMPLE 1** For some DNN specific multi-task networks, often only one specific task head needs to be fine-tuned while freezing the backbone and other task heads. For example, only the detection head can be fine-tuned when input training data are labelled for detection.

- train from scratch: usually after a long period of operation, all parameters of the model can be randomly initialized to re-train the model from scratch. This approach is expected to get better performance than fine-tuning in the original model. However, compared to fine-tuning, this method requires more data volume, computing time and computing resources. Using a pre-trained backbone is a common method. Applying an existing backbone to train on the desired task can reduce the computational cost and speed up the convergence.

### b) re-validation

After re-training, the updated AI model is integrated into AI system and re-validated to provide evidence that the safety-related issues are solved and all relevant safety requirements are met.

**EXAMPLE 2** The datasets of known issues can be used to re-validate the updated AI system in virtual or real-world testing, or a combination of both, and to demonstrate the absence of safety performance degradation, if applicable.

### c) re-approval and re-deployment

After an update to the AI system, the safety assurance argument is re-evaluated (see [Clause 11](#)). Once re-approved, the AI system update can be deployed.

## 14.9 Work products

**14.9.1 Specification of the process and its activities for assuring AI safety during operation**, resulting from [14.3.1](#).

**14.9.2 Specification of the necessary off-board and on-board measures ensuring AI safety during operation**, resulting from [14.3.2](#).

**14.9.3 Field data and functional insufficiencies detected during operation**, resulting from [14.3.3](#).

**14.9.4 Evidence of the effectiveness of measures for ensuring AI system during operation**, resulting from [14.3.4](#).

**14.9.5 Evaluation report of functional insufficiencies detected during operation, and updated version of the safety assurance argument** if applicable, resulting from [14.3.3](#) and [14.3.5](#).

## 15 Confidence in use of AI development frameworks and software tools used for AI model development

### 15.1 Objectives

The objective of this clause is:

- a) to provide requirements and guidance to identify, mitigate and document possible sources of errors and inappropriate biases in the off-line processes, tools and principles used to develop, verify and deploy safety-related AI models.

## 15.2 Prerequisites and supporting information

The following information shall be available at the initiation of this activity:

- a) documentation of development processes and tools used within the AI safety lifecycle (from external sources);
- b) AI system-specific development measures and procedures (from [Clause 7](#) to [Clause 14](#)).

## 15.3 General requirements

**15.3.1** Processes, tools and work products used to develop safety-related AI models shall be analysed to identify, mitigate and document possible sources of errors.

**EXAMPLE** Errors can be caused by inappropriate biases in processes such as field data collection, labelling, sampling, tools such as data processing, deep learning frameworks, work products such as data, AI models.

**NOTE** The approaches discussed in ISO/IEC TR 5469:2024, 11.5.3 and ISO 21448:2022, D.2.5 can be used to analyse offline training processes.

**15.3.2** Confidence shall be demonstrated that software tools used to develop, verify and deploy safety-related AI models are suitable to support activities or tasks required by this document.

**NOTE** ISO 26262-8:2018, Clause 11 can be used to demonstrate confidence in the use of software tools.

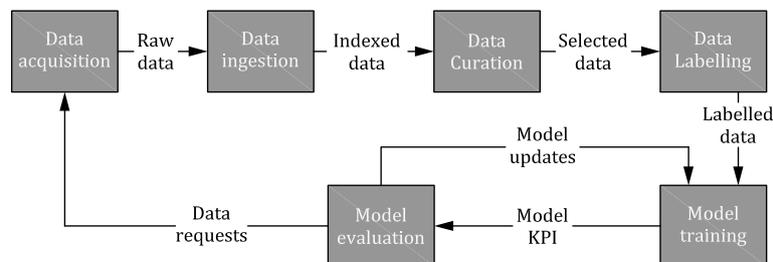
**15.3.3** Appropriate principles for data-driven AI models shall be applied to training and evaluation to ensure control or avoidance of safety-related faults in the AI models.

**NOTE** Design principles that govern software unit design and implementation at the source code level, such as enforcing a single entry and exit point in subprograms and functions, have traditionally been employed to attain the desired quality and robustness in conventional software. However, these principles are applicable solely to the software implementation aspect of data-driven AI models that have been trained and evaluated on data and do not achieve the quality and robustness in the training and evaluation aspects of data-driven AI models.

**EXAMPLE** The influencing factor classes listed in [Table 9-3](#) and their managing approaches elaborated in [Table 9-1](#) can be used as principles.

## 15.4 Confidence in the use of AI development frameworks

Using a robust process to develop AI models reduces the risk of introducing errors in the development of the AI system, thereby making the AI system safer. Specific analysis depends on each system, e.g. automated emergency braking systems and driver status monitoring systems. Typically, AI models are developed using a multi-step process such as the one given in [Figure 15-1](#). The offline training used in the process can be a source of errors in the final AI model. ISO/IEC TR 5469:2024, 11.5.3 proposes the use of a process failure mode and effects analysis (PFMEA) to analyse AI offline training. ISO 21448:2022, D.2.5 describes a similar analysis approach used for SOTIF issues.



**Figure 15-1 — Example offline multistep AI training process for PFMEA**

PFMEA is a well-known technique in the automotive industry<sup>[63]</sup>. PFMEA is an inductive method often applied to manufacturing processes. The analogy is that the offline training process is "manufacturing" an AI model and many of the benefits of a PFMEA apply. It is beyond current technology to trace AI's safety-related systematic issues to root causes, e.g. training and deployment errors, SOTIF issues, etc. Therefore, the overall integrity of training processes, which can be a source of errors in the final AI model, is analysed during AI development. The perspectives used in the safety analysis of systems, e.g. four influencing factor classes described in [Table 9-1](#), are connected to PFMEA.

A PFMEA finds failure modes in each element of the AI training processes. Their effects on subsequent processes and countermeasures to detect such failures are then reviewed. [Table 15-1](#) describes examples of potential failure modes and effects in the AI training processes of [Figure 15-1](#) which can result in a safety-related systematic issue, i.e. performance insufficiencies and safety-related systematic faults as included in [Figure 6-12](#).

**Table 15-1 — Example of potential failure modes and effects in PFMEA**

Process	Potential failure modes	Potential effects
<b>Data acquisition</b> Description: Process step for collecting data to be used in model training and model evaluation	Specific scenes are missing in test datasets.	The model has degraded performance in scenarios involving missing scenes.
	Test data coverage is biased.	In the model evaluation process, evaluation results are biased.
	Only a small number of routes are planned for data collection and the collected training and test datasets lack variation.	Due to lack of variation, the model training process results in low-performance models and test results are unreliable in the model evaluation process.
	Unintended data collection scenarios are used and the collected datasets have inappropriate attributes (meta labels), e.g. weather and time of the day.	The mixture of training data samples using data attributes in the model training process and scene-wise evaluation based on data attributes in the model evaluation do not work as intended.
<b>Data ingestion</b> Description: Process step for uploading collected data to servers used for off-line ML model training and model evaluation	Data is corrupted during upload from data collection vehicles to cloud storages.	Corrupted or lost data during training
<b>Data curation</b> Description: Process step for generating input datasets for further labelling	Curation recognises edge cases as outliers and unintentionally excludes them.	In the model training process, the trained models have performance degradation for these edge cases.
<b>Data labelling</b> Description: Process step for identifying and labelling objects within the datasets to be used for model training and model evaluation	Objects carried or pushed by pedestrians may be included or excluded in the bounding box, leading to inconsistent labelling.	In the model training process, the trained models perform differently for different labelled objects.
	Labellers have biases, e.g. omit labelling motorcycles.	The trained models have performance degradation due to biases.
	Displayed labels and recorded labels are different in data labelling tools.	In the model training process, the trained models learn the wrong labels.

Table 15-1 (continued)

Process	Potential failure modes	Potential effects
Model training Description: Process step for creating trained model from labelled data (e.g. <a href="#">Figure 11-2</a> )	Fixed random seeds are used during the development for debugging and these are left in the production code.	In the model training process, random number generators do not work appropriately and ML and hyperparameter optimization frameworks do not work as intended. As a result, the trained models have consistently low performance.
	Unintended training datasets and AI test datasets are loaded.	Within the model training process, the trained model is optimized to different contexts and has consistently low performance.
Model evaluation Description: Process step for verifying whether the model meets KPIs. A decision is then made to continue with more training, collect more data or end training	The harmonic mean of precision and recall was specified as an evaluation metric, but only recall is evaluated.	As a result, the trained models become recall oriented, i.e. many FP and few FN, which does not meet system requirements.
	Training datasets are leaked to AI test datasets. AI test data sets are not covering the input space definition in a suitable manner.	In the model evaluation process, the evaluation results are not reliable or are overestimated.

Each step of the process can be further broken down (e.g. model training broken down to flow of [Figure 11-2](#)) for a more detailed analysis.

The process analysis may begin as soon as one has a basic understanding of the considered process's inputs, outputs and internal architecture, even if that means proceeding without a complete requirements specification or a complete architecture specification of the process. This iterative process analysis may lead to the specification of additional requirements and process updates. The analysis completes by considering the fully refined architecture and requirements. Even though the analysis focuses on the process for the creation of the AI model, the process analysis can be iterative and may also influence architectural and design decisions.

Example information to start a PFMEA:

- process flow diagram(s) showing the entire AI model creation process flow;
- block diagram of individual process steps including internal process steps showing major components of the process step;
- boundary showing what is the scope of the analysis, the neighbouring process steps to the considered process element;
- identified purpose(s) of the individual steps;
- conceptual data flow(s) between the considered process step, its neighbouring steps and its internal components;
- list of tools used in the AI model creation process;
- training and evaluation requirements.

### 15.5 Confidence in the use of tools used to support the AI-safety lifecycle

The training of AI models often involves tools (e.g. data labelling tools, ML frameworks and hyperparameter optimization frameworks) to train or optimize models. Tools may also be used in the labelling and curation of data along with other steps in the training process. These tools are potential sources of training and deployment errors. ISO 26262-8:2018, Clause 11 can be used to ensure that the tools do not cause an unreasonable safety risk.

### 15.6 Principles for data-driven AI model training and evaluation

Training and (data-driven) evaluation form key parts of the development of data-driven AI models. The causes of insufficiencies of data-driven AI models are classified into the influencing factor classes (Table 9-1), as depicted in Figure 9-3. The certainty in influencing factor classes during the AI model development process ensures the development quality of data-driven AI models. Uncertainties in influencing factor classes can impact a multi-step process like the one given in Figure 15-2. Table 9-1 elaborates on the approaches to manage the certainty of influencing factor classes. These can be used as the principles for data-driven AI model training and evaluation.

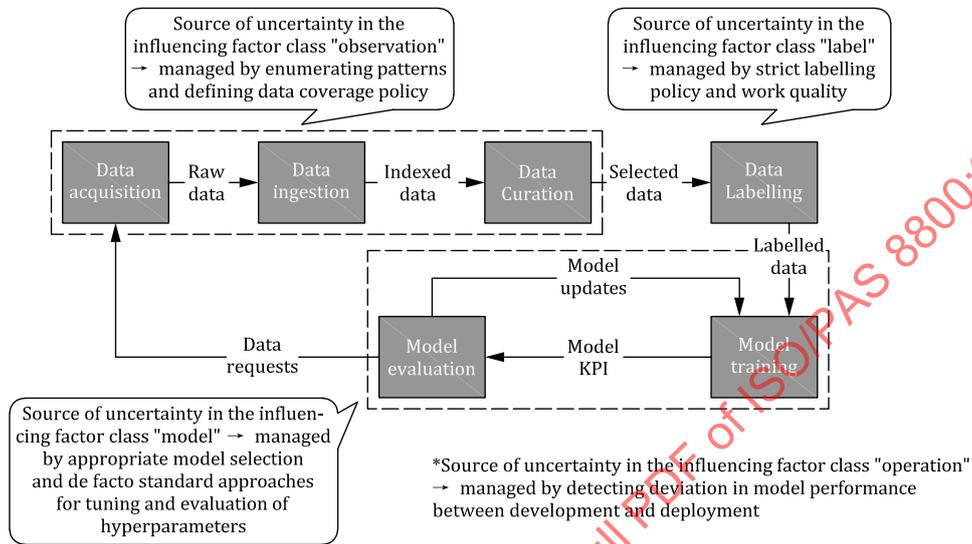


Figure 15-2 — Example offline multistep AI training process and influencing factor classes

### 15.7 Work products

15.7.1 Evidence for the analysis of the AI model creation processes, resulting from 15.3.1.

15.7.2 Evidence for the confidence in the software tools, resulting from 15.3.2.

15.7.3 Evidence for the execution of the AI model creation processes with the principles, resulting from 15.3.3.

**Annex A**  
(informative)

**Overview and workflow of this document**

**Table A-1 — Summary of the normative clauses of this document**

Clause	Objectives	Pre-requisites	Work products
7	<p>a) to define an AI safety lifecycle and its activities to ensure that contributing errors of the AI system do not lead to unreasonable risk of undesired safety-related behaviour at the vehicle-level;</p> <p>b) to ensure that overall and project specific safety management processes and activities are appropriate to ensure the safety of the AI system;</p> <p>c) to plan, initiate and conduct the AI safety activities.</p>	<p>a) the AI system definition (from external sources, e.g. the encompassing system development), including:</p> <ol style="list-style-type: none"> <li>1) the AI system functionality;</li> <li>2) the interfaces of the AI system with the encompassing system, including if applicable, the ASIL capability of the inputs to the AI system.</li> <li>3) the safety requirements allocated to the AI system, including if applicable:                             <ol style="list-style-type: none"> <li>i) the ASIL value of the safety requirements;</li> <li>ii) the acceptance criteria or validation targets derived in conformity to ISO 21448:2022, Clause 6 or 9.</li> </ol> </li> </ol>	<p><a href="#">7.6.1</a></p> <p><a href="#">7.6.2</a></p> <p><a href="#">7.6.3</a></p> <p><a href="#">7.6.4</a></p>

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

Table A-1 (continued)

Clause	Objectives	Pre-requisites	Work products
<p><a href="#">8</a></p>	<p>a) to develop an assurance argument demonstrating that the safety requirements allocated to the AI system are fulfilled;</p> <p>b) to evaluate whether the assurance argument reflects the actual residual risk of the AI system violating its safety requirements;</p>	<p>a) the AI system definition (from external sources), including:</p> <ol style="list-style-type: none"> <li>1) a specification of the safety requirements allocated to the AI system;</li> <li>2) a definition of the technical context within the encompassing system (e.g. definition of interfaces, conditions under which the AI system functionality is triggered, etc.);</li> <li>3) a specification of the input space;</li> </ol> <p>b) requirements on the assurance argument and work products for the AI system (from external sources). These requirements can be derived from the assurance argument of the encompassing system as well as safety management procedures from <a href="#">Clause 7</a>;</p> <p>The following information shall be available for the finalization of these activities:</p> <p>c) the work products of the AI safety life cycle;</p>	<p><a href="#">8.8.1</a></p> <p><a href="#">8.8.2</a></p>
<p><a href="#">9</a></p>	<p>a) to specify a complete and consistent set of safety requirements on the AI system, that are sufficient to ensure AI safety;</p> <p>b) to refine AI safety requirements based on learnings from development, verification and validation;</p> <p>c) to specify the limitations of an AI system over its input space to be escalated to its encompassing system development process.</p> <p>d) the specification of the necessary off-board and on-board measures, from <a href="#">Clause 14</a>.</p>	<p>a) AI system definition (from external sources, e.g. the encompassing system development), including:</p> <ol style="list-style-type: none"> <li>1) safety requirements allocated to the AI system;</li> <li>2) input space definition;</li> <li>3) functional requirements;</li> <li>4) impacted stakeholders;</li> <li>5) the interfaces of the AI system with the encompassing system, including if applicable, the ASIL capability of the inputs to the AI system;</li> <li>6) interfaces to the environment, if applicable.</li> </ol>	<p><a href="#">9.6.1</a></p> <p><a href="#">9.6.2</a></p> <p><a href="#">9.6.3</a></p>

Table A-1 (continued)

Clause	Objectives	Pre-requisites	Work products
<p><a href="#">10</a></p>	<p>a) to select and justify appropriate AI technologies for use in the AI system;</p> <p>b) to identify appropriate architectural and development measures to fulfil the safety requirements prior to deployment;</p> <p>c) to identify appropriate architectural measures to mitigate residual functional insufficiencies of the AI system revealed after deployment;</p> <p>d) to identify measures for ensuring the safety requirements of the AI system are fulfilled within its target execution environment.</p>	<p>a) safety requirements on the AI system, from <a href="#">Clause 9</a>;</p> <p>b) training and validation datasets, from <a href="#">Clause 11</a>;</p> <p>c) AI component or AI system architecture, if already existing;</p> <p>d) AI component or AI system development process, if already existing.</p>	<p><a href="#">10.6.1</a></p> <p><a href="#">10.6.2</a></p> <p><a href="#">10.6.3</a></p>
<p><a href="#">11</a></p>	<p>a) to define the dataset lifecycle of activities related to the gathering, creation, analysis, verification and validation, management, and maintenance of the datasets used in the development of the AI system;</p> <p>b) to identify the dataset insufficiencies that may impact the safety of the AI system;</p> <p>c) to identify the data-related safety properties that have a bearing on the safety of the AI system and that support dataset safety analysis;</p> <p>d) to define the countermeasures to prevent or mitigate dataset insufficiencies using dataset safety analysis methods at different steps in the dataset lifecycle;</p> <p>e) to define the data-related work products that support providing evidence of the safety of the AI system.</p>	<p>a) AI system definition, including:</p> <ol style="list-style-type: none"> <li>1) AI safety requirements, from <a href="#">Clause 9</a>;</li> <li>2) input space definition (refined), from <a href="#">Clause 9</a>;</li> </ol> <p>b) field data and functional insufficiencies detected during operation, from <a href="#">Clause 14</a>;</p> <p>c) safety analysis report, from <a href="#">Clause 13</a>.</p>	<p><a href="#">11.5.1</a></p> <p><a href="#">11.5.2</a></p> <p><a href="#">11.5.3</a></p> <p><a href="#">11.5.4</a></p>

Table A-1 (continued)

Clause	Objectives	Pre-requisites	Work products
<p><a href="#">12</a></p>	<p>a) to verify that the AI system fulfils its AI safety requirements;</p> <p>b) to validate that the safety requirements allocated to the AI system are achieved when integrating into the encompassing system;</p>	<p>a) safety requirements allocated to the AI system (from external sources, e.g. the encompassing system development);</p> <p>b) AI safety requirements, from <a href="#">Clause 9</a>;</p> <p>c) known insufficiencies of the AI system and the corresponding subdomains of the input space, from <a href="#">Clause 9</a>;</p> <p>d) input space definition (refined), from <a href="#">Clause 9</a>;</p> <p>e) AI component or AI system architecture, from <a href="#">Clause 10</a>;</p> <p>f) implemented AI component, from <a href="#">Clause 10</a>;</p> <p>g) dataset lifecycle, from <a href="#">Clause 11</a>;</p> <p>h) evidence for the outputs of the defined phases of the dataset lifecycle, from <a href="#">Clause 11</a>;</p> <p>i) evidence for the safety analyses of the dataset, from <a href="#">Clause 11</a>;</p> <p>j) dataset requirements specification, from <a href="#">Clause 11</a>.</p>	<p><a href="#">12.6.1</a></p> <p><a href="#">12.6.2</a></p> <p><a href="#">12.6.3</a></p>
<p><a href="#">13</a></p>	<p>a) to identify safety-related faults and AI errors that can lead to the violation of AI safety requirements;</p> <p>b) to identify their potential causes;</p> <p>c) to support the definition of safety measures to prevent or control safety-related AI errors;</p> <p>d) to support the verification of AI safety requirements, through modification or identification of new AI safety requirements on data specifications and collection, design specifications, and test specifications.</p>	<p>a) AI safety requirements, from <a href="#">Clause 9</a>;</p> <p>b) input space definition(refined), from <a href="#">Clause 9</a>;</p> <p>c) known insufficiencies of the AI system and the corresponding subdomains of the input space, from <a href="#">Clause 9</a>;</p> <p>d) AI component or AI system architecture (refined), from <a href="#">Clause 10</a>;</p> <p>e) dataset requirements specification, from <a href="#">Clause 11</a>;</p> <p>f) dataset design specification, from <a href="#">Clause 11</a>;</p> <p>g) dataset verification report, from <a href="#">Clause 11</a>;</p> <p>h) dataset validation report, from <a href="#">Clause 11</a>;</p> <p>i) dataset safety analysis report, from <a href="#">Clause 11</a>;</p> <p>j) AI system verification report, from <a href="#">Clause 12</a>;</p> <p>k) AI system validation report, from <a href="#">Clause 12</a>.</p>	<p><a href="#">13.5.1</a></p>

Table A-1 (continued)

Clause	Objectives	Pre-requisites	Work products
<p><a href="#">14</a></p>	<p>a) to define the process requirements to continuously assure AI safety after deployment,                      b) to use the measures defined in <a href="#">clause 8</a> and/or additional measures for the identification of safety risk associated with the AI system during operation and measures to maintain AI safety during operation,                      c) to ensure responses are in place to address unacceptable safety risks associated with the AI system and ensure reapproval of the modified AI system before release.</p>	<p>a) AI system definition (from external sources, e.g. the encompassing system development), including:                      1) interfaces with the encompassing system;                      2) Assumptions of the use of the AI system;                      b) field data collected by the encompassing system;                      c) safety requirements on the AI system, from <a href="#">Clause 9</a>;                      d) AI component or AI system architecture, from <a href="#">Clause 10</a>;                      e) dataset requirements specification, from <a href="#">Clause 11</a>;                      f) dataset design specification and maintenance plan, from <a href="#">Clause 11</a>;                      g) dataset maintenance plan, from <a href="#">Clause 11</a>;                      h) known insufficiencies of the AI system and the corresponding subdomains of the input space, from <a href="#">Clause 9</a>                      i) results of verification and validation activities including known functional insufficiencies of the AI system (if available), from <a href="#">Clause 12</a>;                      j) safety assurance argument, from <a href="#">Clause 8</a>.</p>	<p><a href="#">14.9.1</a>  <a href="#">14.9.2</a>  <a href="#">14.9.3</a>  <a href="#">14.9.4</a>  <a href="#">14.9.5</a></p>
<p><a href="#">15</a></p>	<p>a) to provide requirements and guidance to identify, mitigate and document possible sources of errors and inappropriate biases in the off-line processes, tools and principles used to develop, verify and deploy safety-related AI models.</p>	<p>a) documentation of development processes and tools used within the AI safety lifecycle (from external sources);                      b) AI system-specific development measures and procedures (from <a href="#">Clause 7</a> to <a href="#">Clause 14</a>).</p>	<p><a href="#">15.7.1</a>  <a href="#">15.7.2</a>  <a href="#">15.7.3</a></p>

## Annex B (informative)

### Example assurance argument structure for an AI system

#### B.1 General

This annex provides an example of how an assurance argument for the safety of an AI system based on the principles outlined in this document can be expressed using the goal structuring notation (GSN)<sup>[22]</sup>. The assurance argument structure is expressed as an argument pattern that is intended to be instantiated for a given AI system.

The assurance argument depicted in this annex is for illustrative purposes only and can be used as a starting point for AI system-specific assurance arguments. The argument is not necessarily complete and additional arguments and evidence may be required dependent on the AI system context and specific requirements.

Evidence can be referenced multiple times within the assurance argument.

Work products as defined within this document can contain multiple pieces of evidence as referenced in the assurance argument.

A description of the notation used can be found in Reference <sup>[22]</sup>.

#### B.2 Assurance argument pattern for supervised machine learning

The assurance argument pattern described here can be used to construct an assurance argument for an AI system that makes use of supervised machine learning algorithms (e.g. DNNs). Example applications to which this pattern can apply include the use of AI for image processing tasks such as classification or object detection or predictive maintenance of safety-critical components.

The top level of the assurance argument is depicted in [Figure B.2-1](#). Information that is to be replaced for an AI system-specific instantiation of the pattern are indicated using the following notation: {Instantiable element}. The goal of the argument (G.1) is to demonstrate that the AI system satisfies the requirements allocated to it within the overall system context. This context is defined in terms of:

- a set of assumptions on the input space (A1.1);
- a set of assumptions on the system context (A1.2);
- a definition of the functionality to be implemented by the AI system (C1.1);
- a definition of the safety requirements allocated to the AI system (C1.2).

The assurance argument also assumes that:

- quality management principles have been applied during the development of the AI system and its assurance argument (A1.3) that reduce the risk of systematic errors and increase the confidence in the assurance structure and evidence. The assurance argument is supported by a documented and repeatable development process.
- malfunctioning behaviour caused by random hardware faults or systematic faults are adequately addressed and confirmed through an additional argumentation, not described here. For example, by following the guidance of the ISO 26262 series (A1.4).
- development frameworks and tooling for the AI system do not impact AI safety (addressed by assumption A1.5, see [Clause 15](#)), which may be justified by the use of pre-qualified development frameworks and tools.

The argument is structured to demonstrate that all potential functional insufficiencies in the AI system have been prevented, minimised, or mitigated during the specification, design and operation of the AI system (S1). This strategy makes use of a set of causes of functional insufficiencies for the type of application and applied AI technology (C1.3). These causes can include those described within this document as well as AI system-specific causes that are identified based on safety analyses (see [Clause 13](#)).

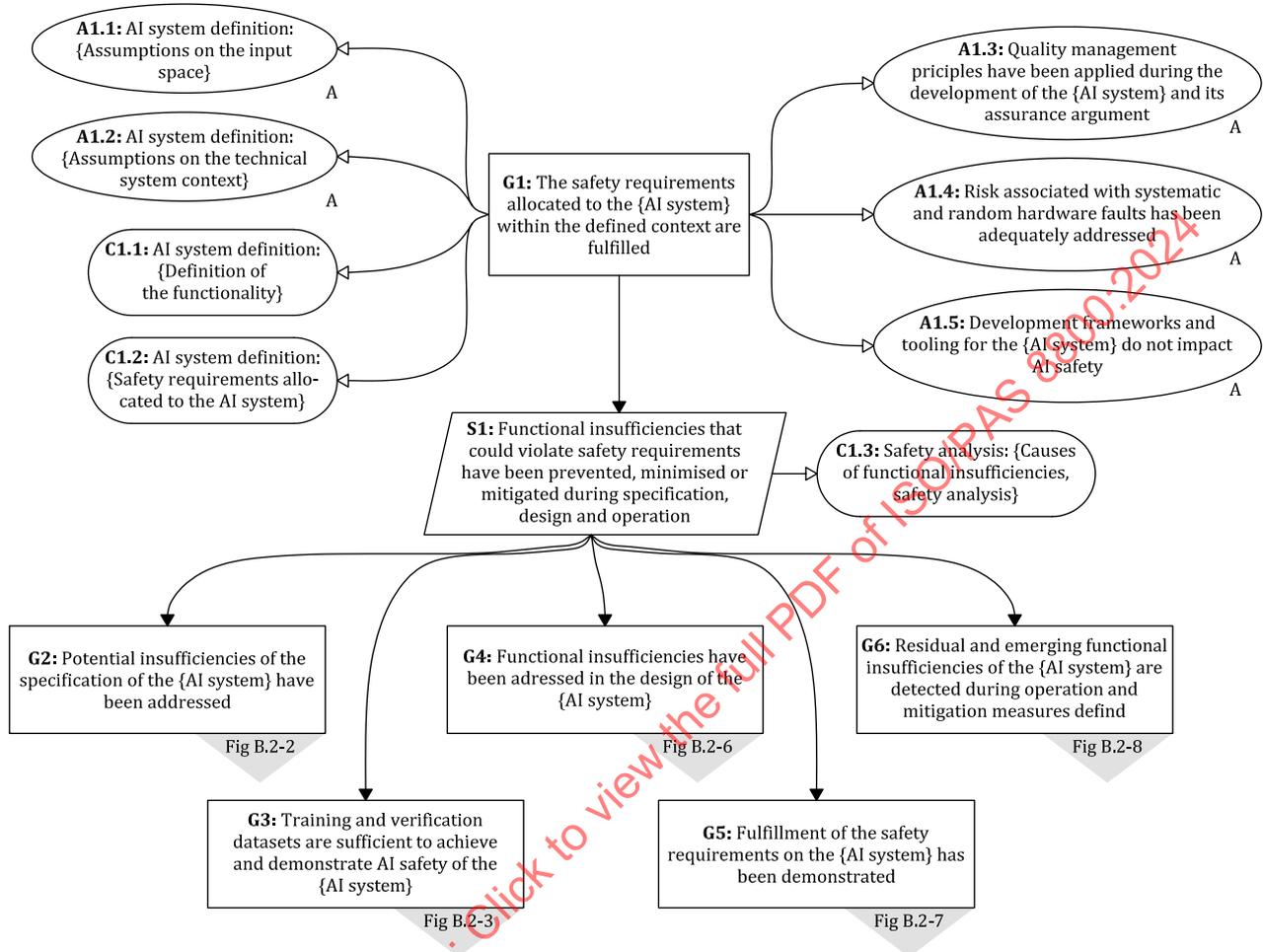


Figure B.2-1 — Assurance argument pattern for a supervised machine learning-based AI system

[Figure B.2-2](#) elaborates the claim G2 of the argument pattern that demonstrates that potential insufficiencies of the specification have been addressed as described in [Clause 9](#). This argument pattern consists of demonstrating:

- a sufficient understanding of the input space (G2.1);
- that the derived AI safety requirements are complete and consistent with respect to the safety requirements allocated to the AI system. This includes demonstrating that each individual AI safety requirement is well defined on the basis of safety-related properties of ML models (S2.2) as well as that the combination of all derived AI safety requirements are sufficient to fulfil the safety requirements allocated to the AI system (G2.2.1);
- the performance limitations of the AI system are sufficiently well defined that a safe behaviour at the system level can be ensured (G2.3).

The derivation of the AI safety requirements as well as the definition of residual performance limitations are supported by the use of safety analyses (see [Clause 13](#)) that determine the potential for safety-related functional insufficiencies and potential causes in the AI system.

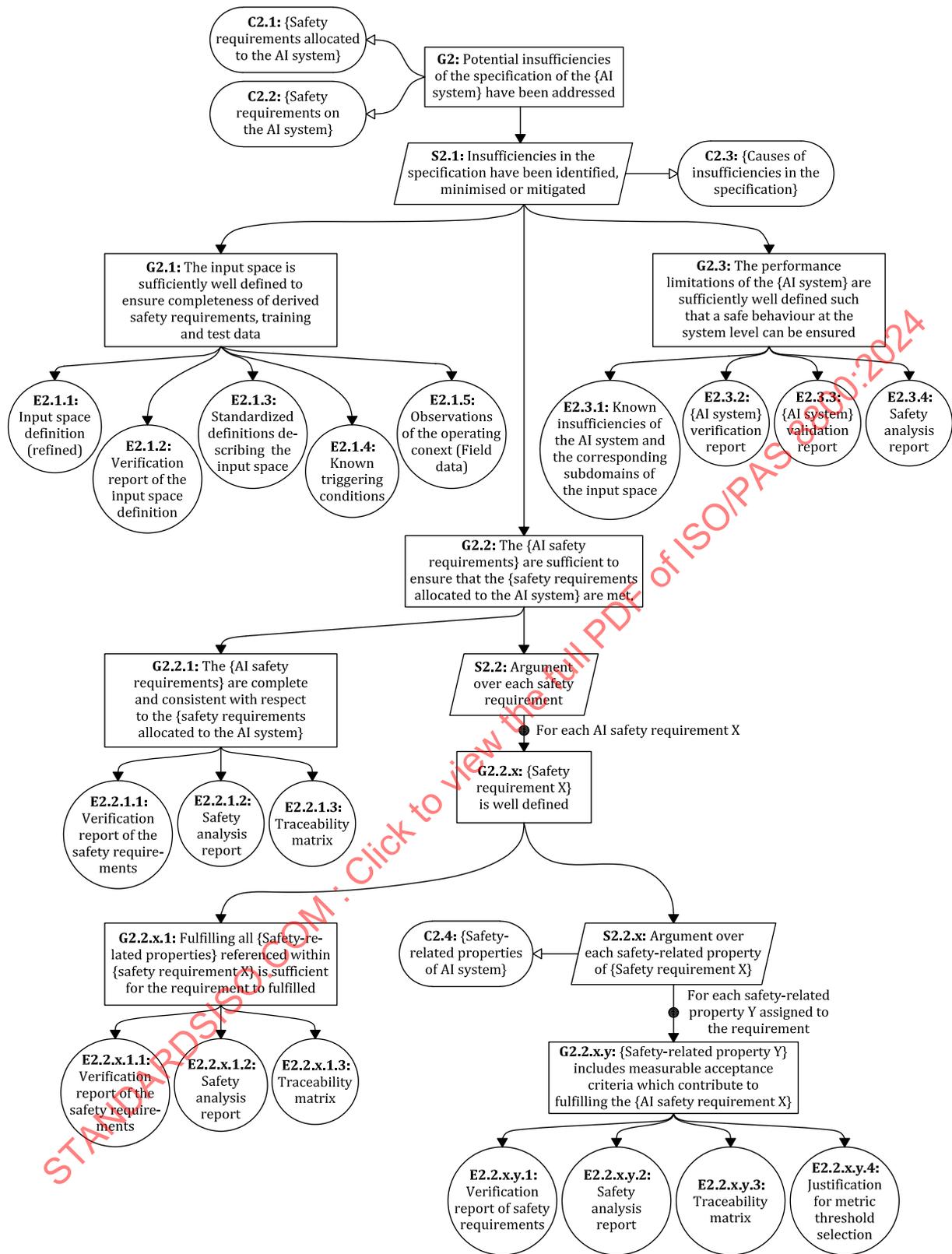


Figure B.2-2 — Assurance argument pattern for demonstrating potential insufficiencies of the specification of the AI system have been addressed

Figure B.2-3, Figure B.2-4 and Figure B.2-5 elaborate the claim G3 of the argument pattern that demonstrates that the datasets used for training and verification of the AI system are sufficient to achieve and demonstrate AI safety, as described in Clause 11. This claim is further refined as follows:

- the datasets consist of suitable selections of observations from the overall input space (G3.1);
- the integrity of the datasets is maintained throughout the data lifecycle (G3.2).

The assurance argument is supported by a set of safety-related properties of the datasets, which can be specific to the application and applied AI technology (see 11.4.3.2 for examples).

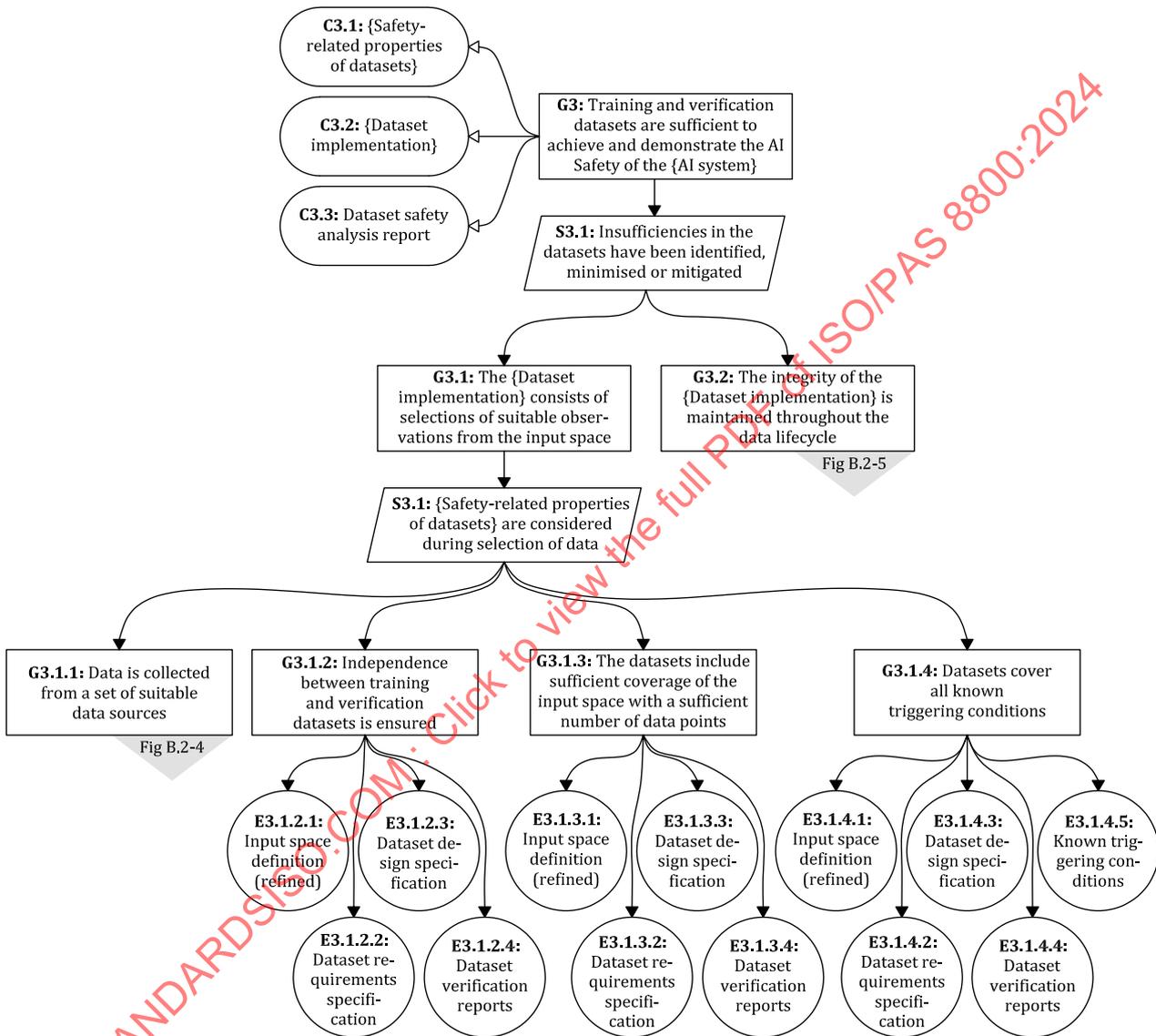


Figure B.2-3 — Assurance argument for the sufficiency of the datasets

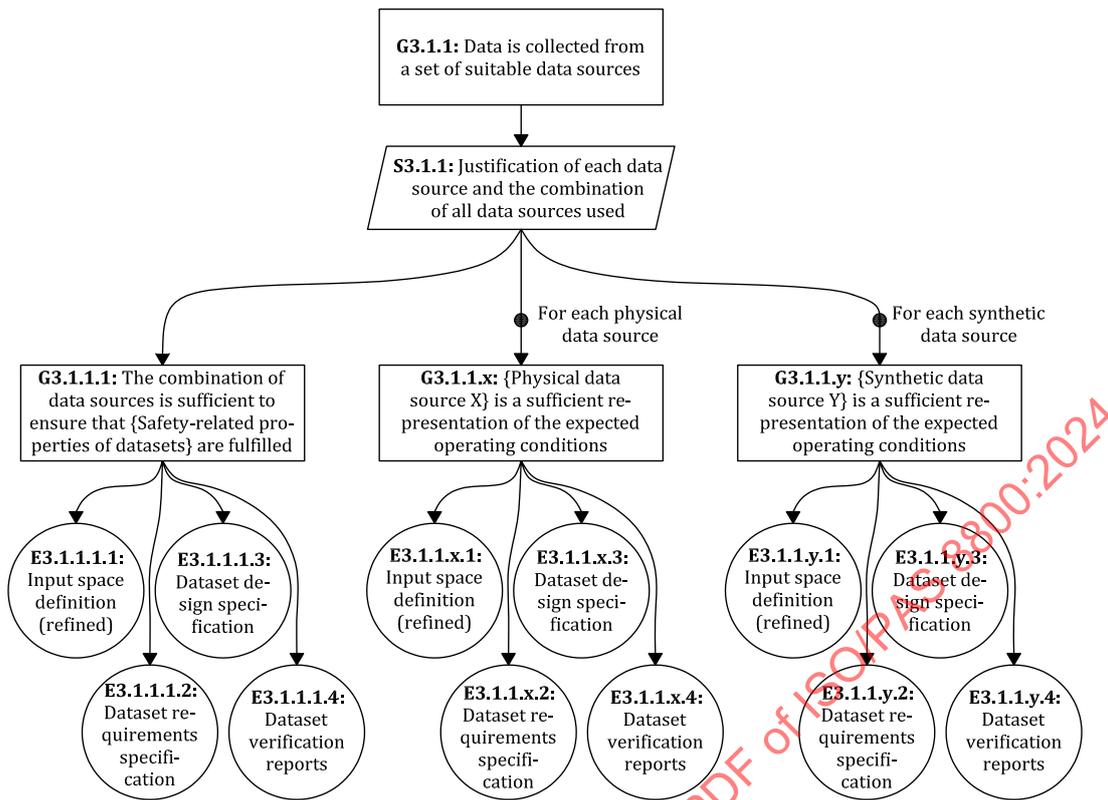


Figure B.2-4 — Assurance argument for claim G3.1.1

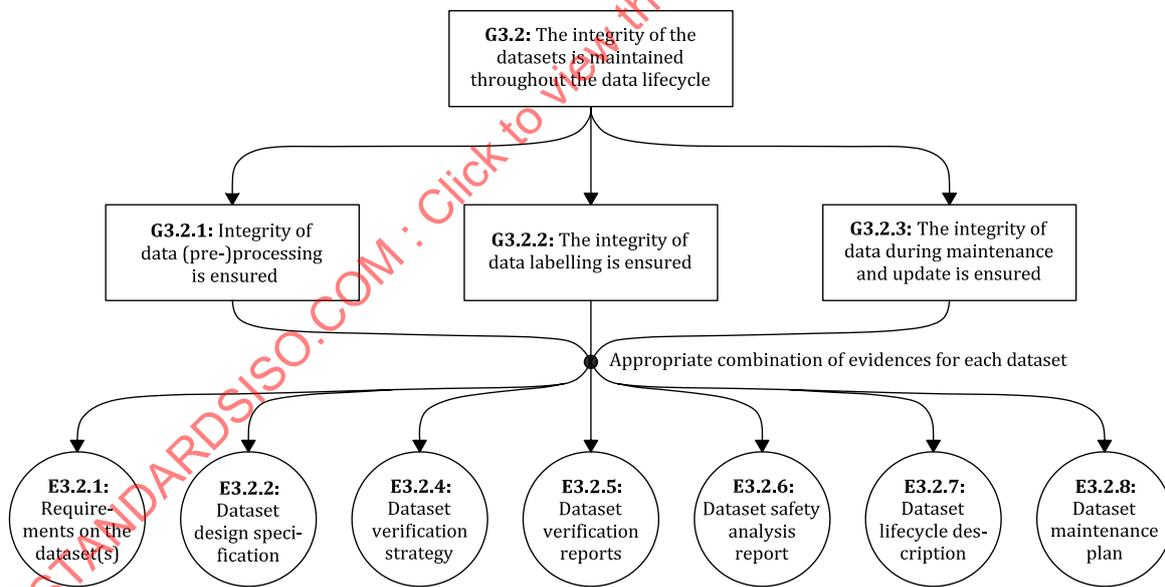


Figure B.2-5 — Assurance argument for claim G3.2

Figure B.2-6 elaborates the claim G4 that demonstrates that functional insufficiencies have been addressed in the design of the AI system, as described in Clause 10. This claim is made based on an understanding of the factors that influence the fulfilment of the AI safety requirements as well as the effectiveness of proposed development and architectural measures. The claim is further refined as follows:

- the chosen AI technology is inherently suitable for achieving the safety requirements allocated to the AI system (G4.1);

## ISO/PAS 8800:2024(en)

- development and architectural measures are chosen that ensure that the AI system meets its AI safety requirements (G4.2);
- architectural measures are identified to mitigate residual insufficiencies in the AI model (G4.3);
- the functional adequacy and sufficient performance are also ensured within its target environment (G4.4).

NOTE G4.4 is not elaborated further in this version of the document.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

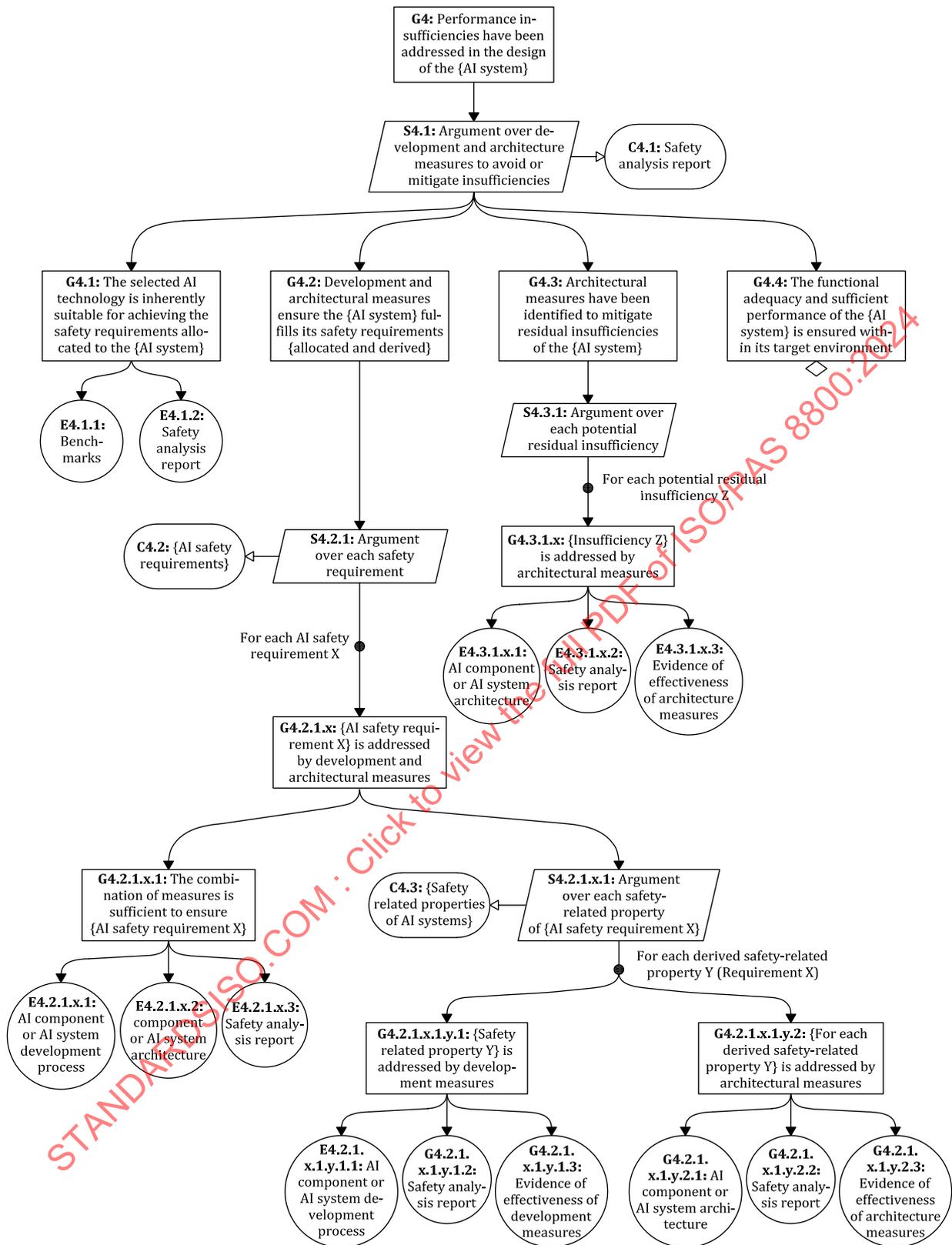


Figure B.2-6 — Assurance argument pattern that functional insufficiencies have been addressed during design

[Figure B.2-7](#) elaborates the claim that sufficient evidence exists that the safety requirements allocated to the AI system have been fulfilled as demonstrated through verification and validation as described in [Clause 12](#). This argument considers:

- the fulfilment of the safety requirements allocated to the AI system in its entirety (G5.1);
- the fulfilment of the derived AI safety requirements allocated to the individual components of the AI system (G5.2).

In each case, an argument is made over the appropriateness of the verification and validation strategy as well as the evidence used to evaluate each individual requirement.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

ISO/PAS 8800:2024(en)

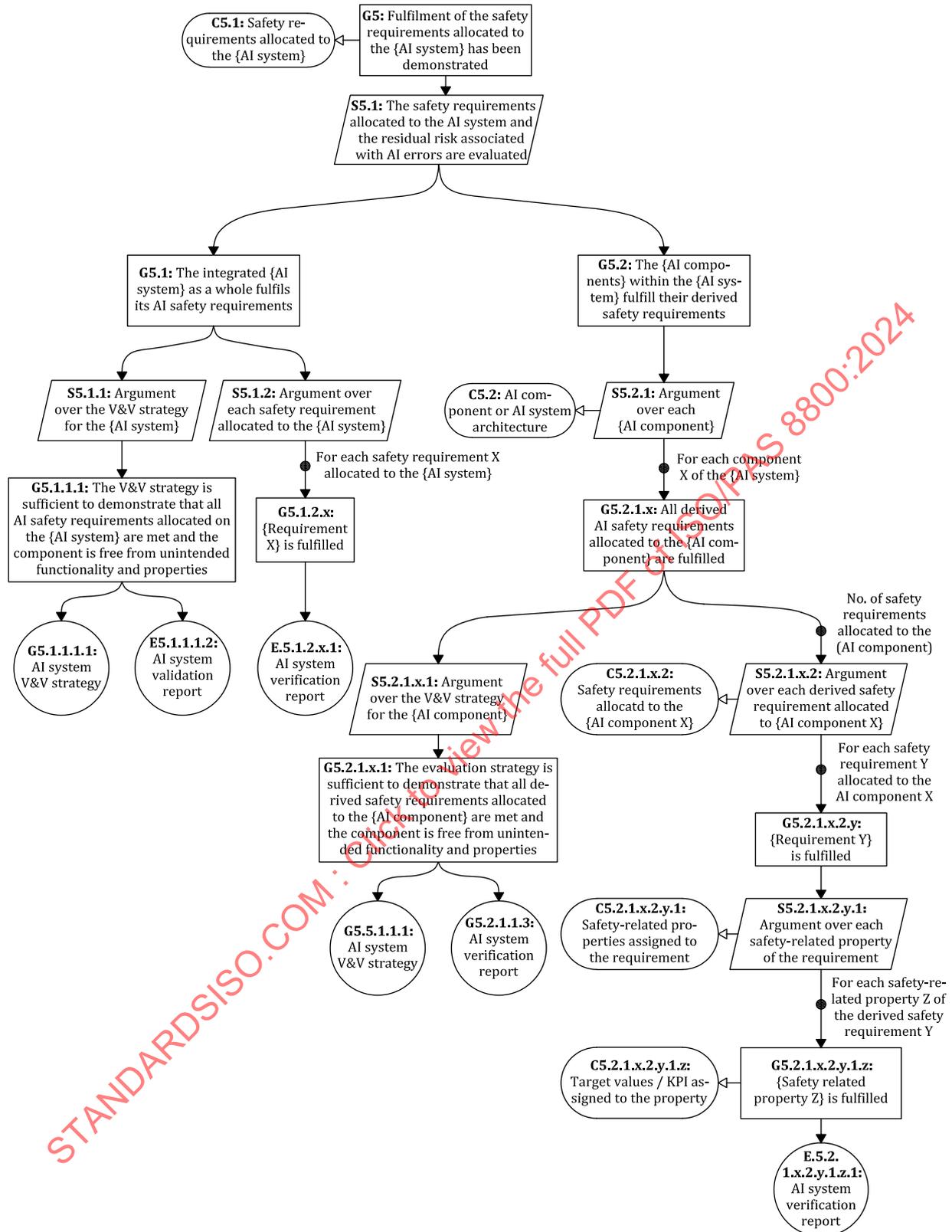


Figure B.2-7 — Assurance argument pattern that the fulfilment of the safety requirements has been adequately demonstrated

[Figure B.2-8](#) elaborates the claim that residual and emerging insufficiencies are identified during operation and mitigation measures are defined, as described in [Clause 14](#). This argument considers:

- the definition of effective operating procedures for the safe operation of the encompassing system based on known insufficiencies of the AI system (G6.1);
- the use of effective processes for continuous re-evaluation of residual risk (G6.2);
- that effective countermeasures are taken to address emerging insufficiencies (G6.3). This claim in the assurance argument can only be made after initial release of the AI system during re-evaluation of the overall safety assurance argument before deployment of changes.

In each case an argument is made over the effectiveness of the measures to control risk during operation.

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

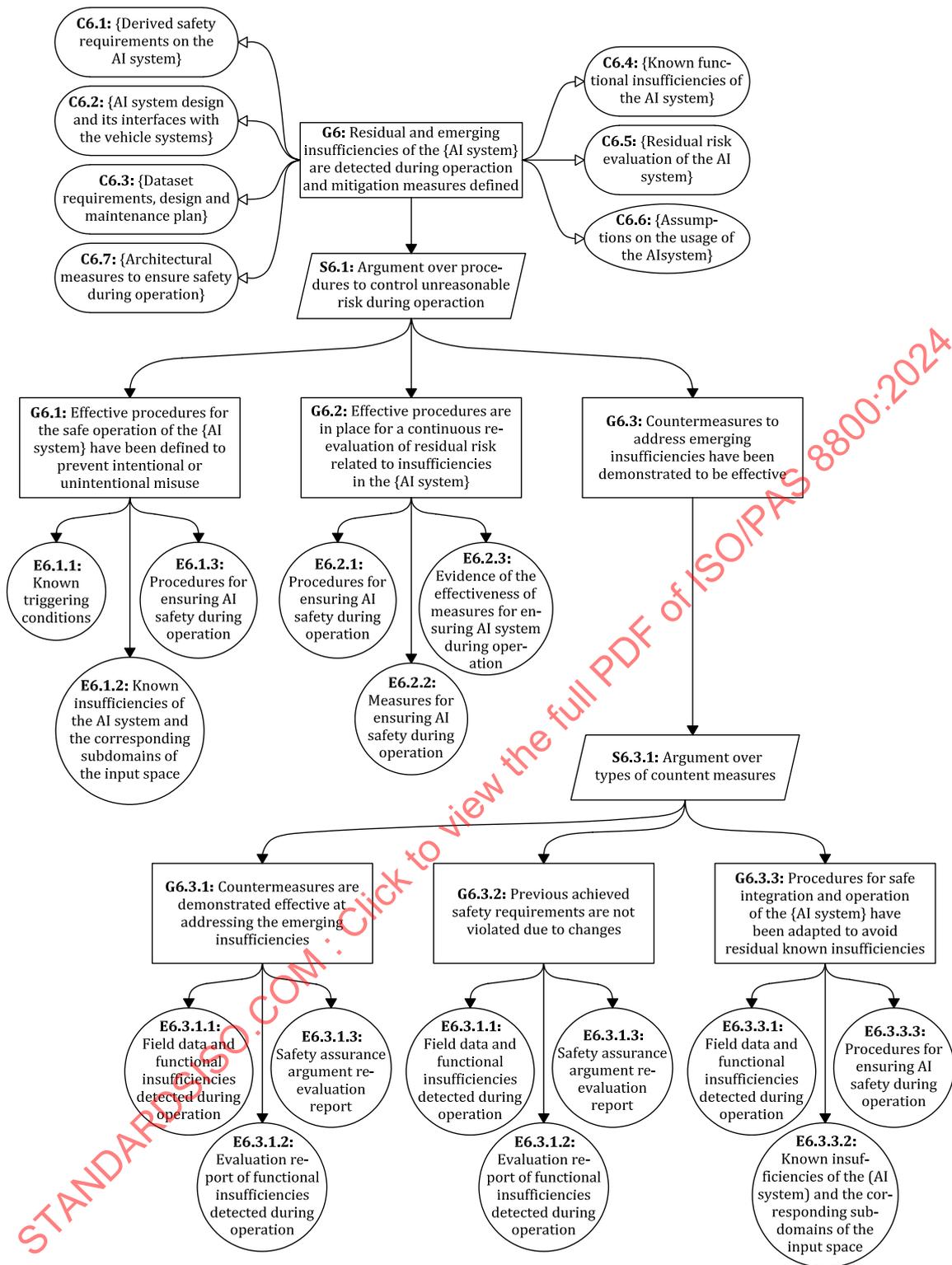


Figure B.2-8 — Assurance argument pattern that emerging and residual insufficiencies are identified and mitigated during operation

## B.3 Use of assurance claim points to increase confidence in the assurance argument

### B.3.1 General remarks on the use of assurance claim points

The GSN pattern outlined in [Clause B.2](#) reflects the basic structure of an assurance argument that safety requirements allocated to a supervised ML-based AI system are fulfilled. The GSN pattern reflects the objectives and requirements of this document.

For any given AI system, the strength of the assurance argument may depend on a number of factors related to the complexity of the task and its environment, the availability of sufficient training and test data and the types of AI technique used. These factors can lead to uncertainty and therefore diminished confidence in the argument.

As outlined in [8.7](#), an evaluation of the argument can include identification of defeaters based on the following types of assertions within the assurance argument<sup>[25]</sup>:

- asserted context: relationship between the claim, contextual information and assumptions;
- asserted evidence: relationship between the claim and the evidence supporting that claim;
- asserted inference: relationship between the claim and the strategies used to structure the sub-claims and evidence to support that claim.

Reference [\[22\]](#) provides the mechanism of assurance claim points (ACPs) to add further, deeper reasoning for particular relationships that would otherwise potentially undermine the confidence in the argument.

The reasoning linked to a particular ACP can be supplied in various forms. A separate GSN model for each ACP is one option, evaluation reports with links to further supporting evidence is another.

The following subclauses provide examples for ACPs to support each of the above types of assertion.

### B.3.2 Example assurance claim points to support assumptions or context: ACP-A2 for assumption A2

Referring to [Figure B.2-1](#), this subclauses addresses ACP-A2 related to the asserted context associated with ACP-A2 as illustrated in [Figure B.3-1](#).

The assumption A1.2 reads “{Assumptions on the technical system context}”. This refers to the technical integration of the AI-system into the encompassing system, i.e. this assumption refers to the interfaces to the other systems and sub-systems as part of the vehicle.

**EXAMPLE 1** ISO/IEC/IEEE 15289:2019, 10.28 mentions some of the properties which are subject to proper interface definitions: “systems or configuration items performing the interface (including human-system and human-human interfaces), standards and protocols, responsible parties, information or data records transmitted by the interface, interface operational schedule, and error handling”.

The information linked to ACP-A2 would demonstrate why the documented interface properties are considered complete in the sense that no safety-relevant interface property remains unspecified, i.e. none of the unspecified interface properties are able to interfere with the achievement of the safety requirements allocated to the AI-system.

**EXAMPLE 2** For a camera-based object detection and classification function implemented by an AI system, information regarding the resolution, depth of focus, quality (e.g. sensor noise), etc., of the camera providing the raw images is documented and analysed to ensure that it meets the assumptions made during the development and test of the AI system.

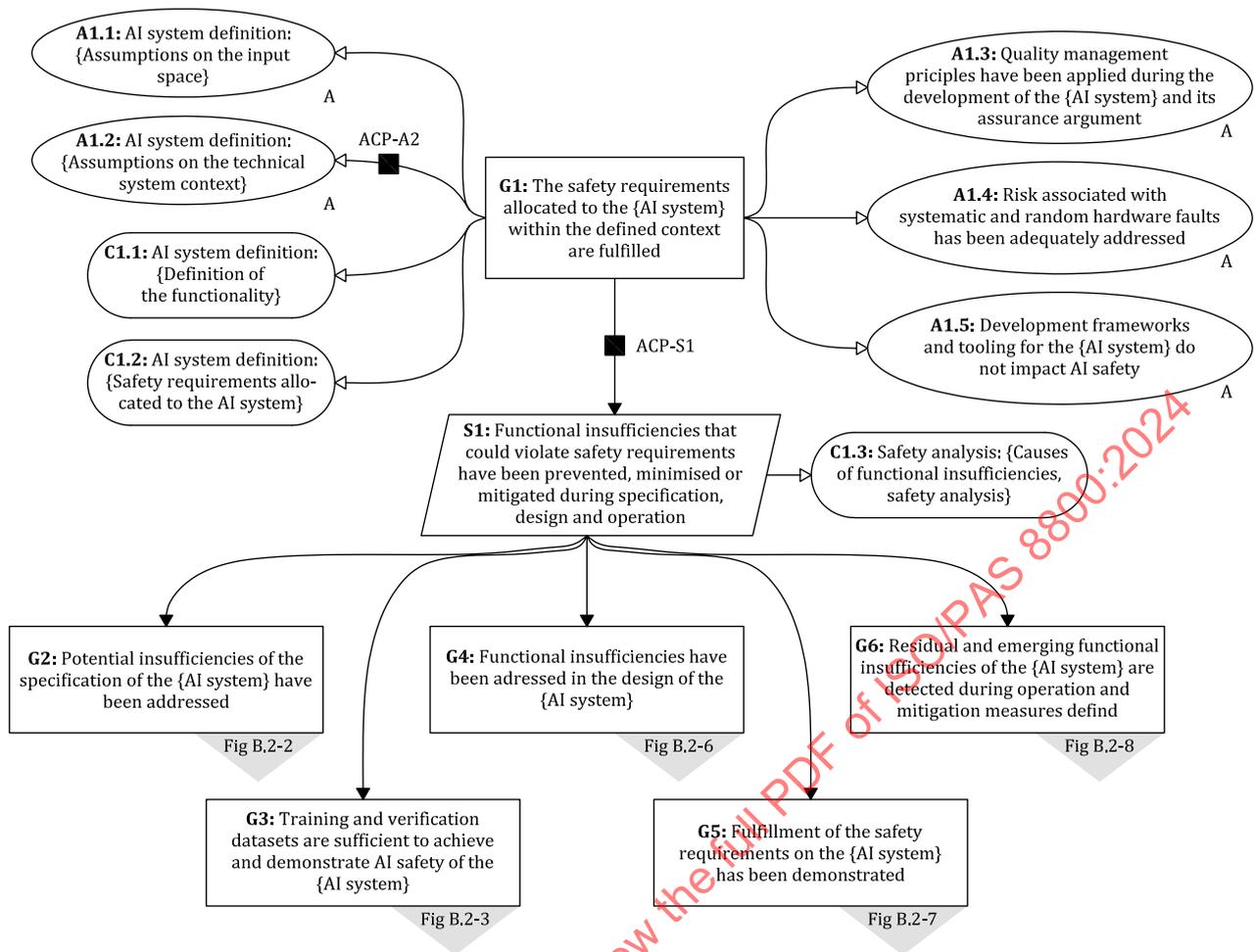


Figure B.3-1 — Example use of ACPs within the GSN assurance argument pattern

### B.3.3 Example assurance claim point to support inference: ACP-S1 for strategy S1

The strategy S1 reads “Functional insufficiencies that can violate safety requirements have been prevented, minimized or mitigated during specification, design and operation”. S1 is supported by sub-goals which reflect the various clauses of this document. The ACP-S1 in [Figure B.3-1](#) is inserted in order to strengthen the assertion that the argumentation strategy based on the set of hypothesized causes of insufficiencies derived from safety analysis is complete and sufficient to demonstrate that the safety requirements allocated to the AI system have been met.

This can be achieved by providing further information on previously (successful) applications of the strategy. Alternatively, this can be achieved by provisioning an evaluation procedure that supervises the strategy and that can flag slipped, untreated functional insufficiencies. Additional forms of reasoning can include reference to effectiveness of the safety analysis approach (see [Clause 13](#)) to identify potential insufficiencies and their causes that may otherwise not have come to light during the development and test of the AI system.

### B.3.4 Example assurance claim point to support evidence: ACP-E5

This subclause provides an example of an assurance claim point to support the assertion directly related to evidence. [Figure B.3-2](#) illustrates the ACP-E5 in the context of the GSN provided previously in [Figure B.2-7](#). This portion of the argument pattern relates to how a combination of evidence demonstrates that individual safety requirements are met.

Evidence is asserted to show the achievement of each AI safety requirement. Such evidence can be a collection of test results documented in one or more test reports. In the case that the safety requirement, test results and test reports are closely aligned with each other, no additional argumentation may be required.

However, in some cases the alignment between test cases and requirement might need additional justification. In other cases, the challenges might be associated with the testability of the requirement itself. Additional reasoning, such as traceability from requirements to test cases, or a description of the approach to indirectly verify a requirement, might be added to the argument using an assurance claim point such as ACP-E5.

This additional reasoning may address both the integrity of the evidence (e.g. have the results of the tests been collected and analysed without loss of critical information?) as well as its validity (e.g. have sufficient tests been performed to ensure a high level of statistical confidence?).

As in the previous subclauses, the assurance claim point can be linked to yet another separate GSN or be backed by some argument in natural language or other supporting analyses.

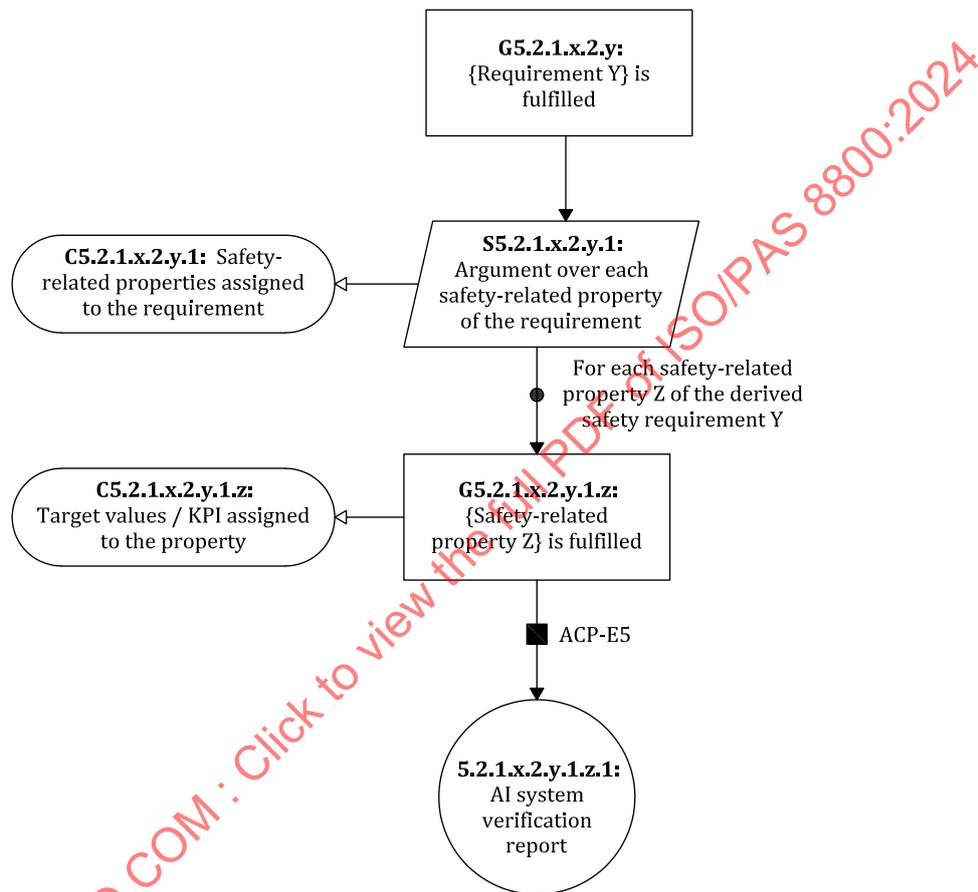


Figure B.3-2 — Example use of ACP to reason about the asserted evidence

## Annex C (informative)

### ISO 26262 gap analysis for ML

#### C.1 General

This annex presents the results of a gap analysis of the ISO 26262 series related to ML. The gap analysis is in the form of an example tailoring and guidance for ISO 26262-4 and ISO 26262-6. The analysis did not find significant gaps related to ISO 26262-1, ISO 26262-2, ISO 26262-3, ISO 26262-5, ISO 26262-7, ISO 26262-8 and ISO 26262-9.

#### C.2 ISO 26262-4:2018 Tailoring and Guidance for ML

[Table C.2-1](#) presents the example tailoring and guidance for the requirements of ISO 26262-4:2018, Clause 8 related to ML. The requirements of ISO 26262-4:2018, Clause 8 that are not listed in the table are considered to not need any additional tailoring or guidance for ML. Different tailoring can be applied to different AI technologies.

**Table C.2-1 — ISO 26262-4:2018, Clause 8 example tailoring/guidance for ML**

Subclause	Requirement	Proposed tailoring/guidance for ML
8.4.3.3	<p>The safety validation at the vehicle level, based on the safety goals, the functional safety requirements and the intended use, shall be executed as planned using:</p> <p>a) the safety validation procedures and test cases for each safety goal including detailed pass/fail criteria;</p> <p>b) the scope of application. This may include issues such as configuration, environmental conditions, driving situations, operational use cases, etc.</p>	<p>Additional guidance:</p> <p>intended use includes representative input space definition (e.g. operating environment, input domain, conditions of use)</p> <p>a) explicitly includes safety-related KPIs</p>
8.4.3.4	<p>An appropriate set of the following methods shall be applied:</p> <p>a) repeatable tests with specified test procedures, test cases and pass/fail criteria;</p> <p>b) analyses;</p> <p>c) long-term tests, such as vehicle driving schedules and captured test fleets;</p> <p>d) operational use cases under real-life conditions, panel or blind tests, or expert panels;</p> <p>e) reviews.</p>	<p>Guidance</p> <p>b) may be limited (e.g. simulation only)</p> <p>e) is typically not applicable for ML validation</p>

C.3 ISO 26262-6:2018 Tailoring for ML

Table C.3-1 presents the example tailoring and guidance for the requirements of ISO 26262-6 related to ML. Where noted, the tailoring/guidance is related to NN models only. The requirements of ISO 26262-6 that are not listed in the table are considered not to need any additional tailoring or guidance for ML.

Table C.3-1 — ISO 26262-6:2018 example tailoring/guidance for ML

Subclause	Requirement/method	Proposed tailoring/guidance for ML
5.4.3, Table 1	1a Enforcement of low complexity 1b Use of language subsets 1c Enforcement of strong typing 1d Use of defensive implementation techniques 1e Use of well-trusted design principles 1f Use of unambiguous graphical representation 1g Use of style guides 1h Use of naming conventions 1i Concurrency aspects	Tailoring For ML applications, Table 1 applies unchanged for use case independent elements (e.g. CUDA C++ libraries). For the use case dependent elements (i.e. the models), 1c, 1d, 1g, 1i with “o” for all ASILs (justification see ISO/IEC TR 5469:2024 Tables A.3 and A.4). NOTE Use case independent elements refer to elements that behave the same independent of the use case (i.e. CUDA libraries fulfil the same purpose independently if the trained models use case is in autonomous driving or predictive maintenance). In contrast, use case dependent elements like neural network models are dependent on the specific use case, which changes their properties.
6.4.1	The software safety requirements shall be derived considering the required safety-related functionalities and properties of the software, whose failures can lead to the violation of a technical safety requirement allocated to software	Guidance in the form of additional considerations Requirements in the form of: a) KPIs; b) data attributes; c) dataset requirements (e.g. input space definition) specification for training datasets, AI validation datasets, and AI test datasets The ML implementation not meeting its KPIs is an ISO 26262 issue. Incorrect or insufficient KPIs are a SOTIF concern.
6.4.4	The hardware-software interface specification initiated in ISO 26262-4:2018, Clause 6, shall be refined sufficiently to allow for the correct control and usage of the hardware by the software, and shall describe each safety related dependency between hardware and software	Guidance in the form of examples EXAMPLE 1 Software is specified to run on one CPU on a multi-CPU system. EXAMPLE 2 NN specified to run on a GPU

Table C.3-1 (continued)

Subclause	Requirement/method	Proposed tailoring/guidance for ML
6.4.7	<p>The software safety requirements and the refined requirements of the hardware-software interface specification shall be verified in accordance with ISO 26262-8:2018, Clauses 6 and 9, to provide evidence for their:</p> <ul style="list-style-type: none"> <li>a) suitability for software development;</li> <li>b) conformity to and consistency with the technical safety requirements;</li> <li>c) conformity to the system design; and</li> <li>d) consistency with the hardware-software interface.</li> </ul>	<p>Guidance in the form of additional considerations</p> <ul style="list-style-type: none"> <li>e) adequate coverage of the input space of the software.</li> </ul> <p>Adequate coverage of the input space typically involves:</p> <ul style="list-style-type: none"> <li>— sufficient labels that comprehend the entire labelling space (e.g. labels for emergency vehicles)</li> <li>— data with multiple views (e.g. multiple examples of emergency vehicles)</li> </ul> <ul style="list-style-type: none"> <li>f) address the handling of out-of-distribution inputs</li> </ul>
7.4.3	<p>In order to avoid systematic faults, the software architectural design shall exhibit the following characteristics by use of the principles listed in Table 3:</p> <ul style="list-style-type: none"> <li>a) comprehensibility</li> <li>b) consistency</li> <li>c) simplicity</li> <li>d) verifiability</li> <li>e) modularity</li> <li>f) encapsulation</li> <li>g) maintainability</li> </ul>	<p>Guidance</p> <ul style="list-style-type: none"> <li>1) software components implemented using ML are difficult to verify. A heuristic component is preferred over a ML component assuming the function can acceptably be implemented using a heuristic component.</li> <li>2) NN models are considered as individual units. The principles of Table 3 typically cannot be met for NN applications. Usually, they are generated from higher level languages using tools. The design principles therefore are applied to the code that generates the NN model and tool qualification are applied to the generator. This is similar to the usage of code generation in normal SW development.</li> </ul>
7.4.4	<p>The software architectural design shall be developed down to the level where the software units are identified.</p>	<p>Guidance</p> <ul style="list-style-type: none"> <li>1) an individual NN may consist of many nodes and layers but is typically considered as one unit. It may not be possible to express an NN software design at any level lower than the individual NN level.</li> <li>2) a SW unit can be an NN so long as suitable interfaces can be defined and requirements allocated to those units.</li> <li>3) an architecture description for an NN model, e.g. in ONNX, can be created. Nevertheless, explainability based on the architecture may be low however a justification for the choice of the structure can be provided, e.g. motivated by ablation studies.</li> </ul>

Table C.3-1 (continued)

Subclause	Requirement/method	Proposed tailoring/guidance for ML
7.4.7	If a pre-existing software architectural element is used without modifications in order to meet the assigned safety requirements without being developed according to the ISO 26262 series, then it shall be qualified in accordance with ISO 26262-8:2018, Clause 12.	Guidance For pre-existing ML based software when the specification characteristics such as dataset attributes, KPIs, input space definition and output metrics are articulated, the verification should ensure that the specification characteristics are sufficiently met.
7.4.13	An upper estimation of required resources for the embedded software shall be made, including: a) execution time b) storage space c) communication resources	Guidance in the form of additional considerations d) parallel computation resources
8.4.3	The software unit design shall be described using the notations listed in Table 5 in order to avoid systematic faults and to ensure that the software unit design achieves the following properties: a) consistency b) comprehensibility c) maintainability d) verifiability	Guidance in the form of additional considerations Additionally, use the derived AI safety-related properties for the given AI system as appropriate, for example: e) interpretability (see ISO/IEC TR 5469 for definition); f) explainability (see <a href="#">Annex D</a> and ISO/IEC TR 5469 for definition); g) predictability (see <a href="#">Annex D</a> for definition); h) specificity (see ISO/IEC TR 5469 for definition); i) generalisation (see <a href="#">Annex D</a> and ISO/IEC TR 5469 for definition); j) domain shift (see ISO/IEC TR 5469 for definition); k) robustness-safeness (see ISO/IEC TR 5469 for definition); m) diversity (see ISO/IEC TR 5469 for definition); n) confidence (see ISO/IEC TR 5469 for definition).
8.4.4	The specification of the software units shall describe the functional behaviour and the internal design to the level of detail necessary for their implementation.	Guidance For the case of a unit containing an NN, the structure of the NN (e.g. number of nodes, layout, interconnects and activation function) and hyperparameters and training methods of NN (e.g. learning rate) are part of the specification of the software unit.

Table C.3-1 (continued)

Subclause	Requirement/method	Proposed tailoring/guidance for ML
8.4.5	<p>Design principles for software unit design and implementation at the source code level as listed in Table 6 shall be applied to achieve the following properties:</p> <ul style="list-style-type: none"> <li>a) correct order of execution of sub-programs and functions within the software units, based on the software architectural design;</li> <li>b) consistency of the interfaces between the software units;</li> <li>c) correctness of data flow and control flow between and within the software units;</li> <li>d) simplicity;</li> <li>e) readability and comprehensibility;</li> <li>f) robustness;</li> <li>g) suitability for software modification;</li> <li>h) verifiability.</li> </ul>	<p>Guidance</p> <ul style="list-style-type: none"> <li>1) for NN units, a) and c) may not apply since the NN is considered as one function and the order of execution of individual nodes is not guaranteed.</li> <li>2) for h), the structure of the network can be verified, for example, by inspection that the correct structure is implemented.</li> <li>3) for AI models trained using data, the influencing factors of <a href="#">Table 9-1</a> may be considered as additional design principles: observation certainty, label certainty, model certainty and operation certainty</li> </ul>
9.4.2	<p>The software unit design and the implemented software unit shall be verified in accordance with ISO 26262-8:2018, Clause 9 by applying an appropriate combination of methods according to Table 7 to provide evidence for:</p> <ul style="list-style-type: none"> <li>a) conformity to the requirements regarding the unit design and implementation in accordance with <a href="#">Clause 8</a>;</li> <li>b) the conformity of the source code to its design specification;</li> <li>c) conformity to the specification of the hardware-software interface (in accordance with 6.4.4), if applicable;</li> <li>d) confidence in the absence of unintended functionality and properties;</li> <li>e) sufficient resources to support their functionality and properties; and</li> <li>f) implementation of the safety measures resulting from the safety-oriented analyses in accordance with 7.4.10 and 7.4.11.</li> </ul>	<p>Guidance</p> <ul style="list-style-type: none"> <li>a) the NN model software verification report documents the test result KPI and the dataset used for the testing.</li> </ul> <p>Tailor d) to</p> <ul style="list-style-type: none"> <li>d) confidence in the absence of unintended functionality and properties (unintended functionality is primarily a SOTIF concern for systems modelled using ML).</li> </ul>

Table C.3-1 (continued)

Subclause	Requirement/method	Proposed tailoring/guidance for ML
9.4.2 Table 7	1a Walkthrough 1b Pair-programming 1c Inspection 1d Semi-formal verification 1e Formal verification 1f Control flow analysis 1g Data flow analysis 1h Static code analysis 1i Static analyses based on abstract interpretation 1j Requirement-based test 1k Interface test 1l Fault injection test 1m Resource usage evaluation 1n Back-to-back test between model and code, if applicable	Tailoring 1a through 1i for ML "o" since often infeasible to do effectively Guidance 1) for ML, 1j feasible for only some properties such as invariants and equivariants 2) for ML, 1l fault injection has limited applicability 3) for ML, 1n applicable when comparing off-line versus optimized code versions
9.4.3	To enable the specification of appropriate test cases for the software unit testing in accordance with 9.4.2, test cases shall be derived using the methods as listed in Table 8.	Guidance For ML, unit test cases can be selected from test dataset
9.4.4	To evaluate the completeness of verification and to provide evidence that the objectives for unit testing are adequately achieved, the coverage of requirements at the software unit level shall be determined and the structural coverage shall be measured in accordance with the metrics as listed in Table 9. If the achieved structural coverage is considered insufficient, either additional test cases shall be specified or a rationale based on other methods shall be provided.	Tailoring Requirement NA for NNs, includes NA for Table 9 For NNs without separate program statements, branches or decision logic this requirement does not apply. An example of where this requirement is still applicable is conditional computation in neural networks, which is sometimes implemented to reduce latency and save energy (i.e. only part of a net is activated). It is possible to select inputs as unit tests to cover all branches.
9.4.5	The test environment for software unit testing shall be suitable for achieving the objectives of the unit testing considering the target environment. If the software unit testing is not carried out in the target environment, the differences in the source and object code, as well as the differences between the test environment and the target environment, shall be analysed in order to specify additional tests in the target environment during the subsequent test phases.	Guidance in the form of an example EXAMPLE A NN model is trained using fp32 math, but the on-line inferencing uses int8 for throughput and bandwidth savings. Off-line unit testing of the model uses int8 to match the target environment.

Table C.3-1 (continued)

Subclause	Requirement/method	Proposed tailoring/guidance for ML
10.4.2 Table 10	1a Requirements-based test 1b Interface test 1c Fault injection test 1d Resource usage evaluation 1e Back-to-back test between model and code, if applicable 1f Verification of the control and data flow 1g Static code analysis 1h Static analyses based on abstract interpretation	Tailoring 1c) For NNs, targeted SW fault injection might only be appropriate at certain interfaces. HW fault injection on the target environment can test the response to permanent and transient faults. 1g) For NNs "o" 1h) For NNs "o" Justification g) and h) Static code analysis aiming to verify functionality of NN does not scale beyond small networks

STANDARDSISO.COM : Click to view the full PDF of ISO/PAS 8800:2024

## Annex D (informative)

### Detailed considerations on safety-related properties of AI systems

This Annex provides a list of properties of AI systems that are considered desirable/necessary from a safety perspective. These properties are conceptual, and the list is based on past AI development experience and is not exhaustive. See [Table D-1](#).

Safety-related properties can be quantitative in nature as well as qualitative. As a result, they are not always completely achievable. For example, while the robustness property indicates that a model is either robust or not, a DNN model for classifying objects in an open world is never 100 % robust against all possible insignificant input changes. The choice of safety-related properties relevant to the AI system should be validated through safety analysis to ensure their contribution to the system's safety, and target thresholds should be provided with justification.

A safety-related property may or may not apply depending on use cases, systems, AI models, etc. For example, while a self-driving vehicle's actions, such as acceleration and steering, can be controllable, the outputs of a DNN model for object detection in the perception pipeline of the vehicle are not.

The scope of a safety-related property of AI systems refers to the entity to which the property is attributed. For example, the organization can effectively and efficiently update an AI model whenever necessary. In this context, the overall system is considered the whole product, e.g. the vehicle.

NOTE One or more KPIs are typically defined to characterize each safety-related property. The safety requirement specifies the acceptable threshold value for these KPIs.

**Table D-1 — Safety-related properties of AI systems**

Property	Description	Scope
AI robustness	Ability to maintain an acceptable level of performance under the presence of semantically insignificant, but reasonably expected changes to the input (see <a href="#">3.1.12</a> for the definition of AI reliability) NOTE AI robustness focuses on foreseeable/relevant perturbations (type of perturbation as well as amplitude of perturbation) which can occur in the real world, to avoid defining unnecessary safety requirements.	Model, system
AI generalization capability	Ability of a model to adapt and perform well on the previously unseen data during inference	Model, system
AI reliability	Ability to maintain all functionalities for a specified period (see <a href="#">3.1.12</a> for the definition of AI reliability)	Model, system
AI resilience	Ability to quickly recover from an incident (see <a href="#">3.1.13</a> for the definition of AI resilience)	(Overall) system, organization
AI controllability	Ability of an external agent to overwrite the behaviour or output of an AI system	(Overall) system, organization