

© 2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Accepted to be published in *2025 IEEE International Automated Vehicle Validation Conference (IAVVC)*, Baden-Baden, Germany, 2025.

Cite as:

M. Nolte, N. F. Salem, O. Franke, J. Heckmann, C. Höhmann, G. Stettinger, and M. Maurer, “What’s Really Different with AI ? – A Behavior-based Perspective on System Safety for Automated Driving Systems,” in *2025 IEEE International Automated Vehicle Validation Conference*, to be published.

BibTeX:

```
@inproceedings{nolte_salem_AI_safety_2025,
  author = {Nolte, Marcus and Salem, Nayel Fabian and Franke, Olaf and Heckmann, Jan and H\"ohmann,
    Christoph and Stettinger, Georg and Maurer, Markus},
  booktitle = {12},
  title = {{What's} {Really} {Different} with {AI} ? -- {A} {Behavior-based} {Perspective} on {System} {
    Safety} for {Automated} {Driving} {Systems}},
  address = {Baden-Baden, Germany},
  year = {2025},
  publisher = {IEEE. to be published},
}
```

# What's Really Different with AI? – A Behavior-based Perspective on System Safety for Automated Driving Systems

Marcus Nolte<sup>\*✉</sup>, Nayel Fabian Salem<sup>\*✉</sup>, Olaf Franke<sup>†</sup>, Jan Heckmann<sup>‡</sup>, Christoph Höhmann<sup>§</sup>,  
Georg Stettinger<sup>¶✉</sup>, Markus Maurer<sup>\*✉</sup>

<sup>\*</sup>*TU Braunschweig*  
*Institute of Control Engineering, Braunschweig, Germany*  
{m.nolte, n.salem, m.maurer}@tu-braunschweig.de

<sup>†</sup>*MAN Truck & Bus SE, Munich, Germany*  
olaf.franke@man.eu

<sup>‡</sup>*Deutsche Bahn Regio AG, Berlin, Germany*  
jan.heckmann@deutschebahn.com

<sup>§</sup>*Mercedes-Benz AG, Stuttgart, Germany*  
christoph.hoehmann@mercedes-benz.de

<sup>¶</sup>*Infineon Technologies AG, Neubiberg, Germany*  
georg.stettinger@infineon.com

**Abstract**—Assuring safety for “AI-based” systems is one of the current challenges in safety engineering. For automated driving systems, in particular, further assurance challenges result from the open context that the systems need to operate in after deployment. The current standardization and regulation landscape for “AI-based” systems is becoming ever more complex, as standards and regulations are being released at high frequencies.

This position paper seeks to provide guidance for making qualified arguments which standards should meaningfully be applied to (“AI-based”) automated driving systems. Furthermore, we argue for clearly differentiating sources of risk between AI-specific and general uncertainties related to the open context. In our view, a clear conceptual separation can help to exploit commonalities that can close the gap between system-level and AI-specific safety analyses, while ensuring the required rigor for engineering safe “AI-based” systems.

**Index Terms**—artificial intelligence, safety, automated vehicles

## I. INTRODUCTION

Alongside the ongoing deployment of Automated Driving Systems (ADS) and recent technological progress in Artificial Intelligence (AI), huge research efforts are currently geared toward safety assurance for “AI-based” automated driving systems. These efforts are sidelined by the release of a plethora of different standards and regulations for both, general AI-based applications and AI-based automated driving systems.

The fast-paced standardization and regulation landscape poses challenges to all stakeholders involved in the development of automated vehicles: Product-compliance laws typically demand a development aligned with the current state of the art. Depending on the actual regulation, this state of the art refers to combinations of research, standardization, industry best practices. Consistent regulation and standardization is

critical here to support a safe industrial realization of ADS, as contradicting definitions of what an “AI-based system” is or what “safety” means can lead to competitive disadvantages in more tightly regulated markets such as the European Union.

Three main challenges follow from this which are related to a) extracting relevant demands from closely related standards that consider ADS, “AI-based” systems or both, b) providing compelling arguments, why a standard might *not* be applicable to the developed system, and c) defining processes that efficiently address relevant contributions to the state of the art from such standards, e.g. by establishing methods for safety analyses that allow addressing multiple standards at once. These tasks are further complicated by different definitions in current standardization and regulation for what constitutes an “AI-based” system.

In the context of safety assurance for “AI-based” automated driving systems, it is crucial to clearly differentiate possible sources of risk. Previous work in this context (e.g. [1–3], cf. also Fig. 1) has argued the importance of tracing back relevant risks to different sources of uncertainty that affect a system of interest. Such uncertainty can, e.g., include knowledge gaps related to the open world in which a system needs to operate or uncertainty stemming from the world’s complexity (e.g. such as unpredictability of long-term interactions between traffic participants). This type of uncertainty reflects in general challenges related to the definition of safe system-level behavior, regardless of how a system is implemented. It will, however, also have an impact on specific “AI-related” challenges in the sense that this type of uncertainty requires addressing the representativeness of training data or assessing the generalization capabilities of data-driven algorithms in contexts that are unknown at design or training time.

---

M. Nolte and N. F. Salem contributed equally to this work.

Further uncertainty stems from the complexity of the system itself and the (“AI-”) algorithms that are applied to implement system functionality. It becomes crucial to differentiate what risks are related to uncertainty caused by emergent system properties that are related to complex *system architectures* (i.e., high numbers of interrelated system components) and what risks are related to the specific nature of e.g. the black-box properties of (*network*) *architectures* in data-based algorithms. To provide guidance in the current regulation and standardization landscape, this position paper will address the challenges related to differing “AI” definitions (section II).

Furthermore, we find that current research, standards and regulation are separating assurance challenges too stringently into the categories of “AI-related” and “non-AI-related”. Often, explicit differences between “AI-based” and “classic” systems [4–6] are emphasized.<sup>1</sup> We will compare challenges for the safety assurance of automated driving systems operating in open contexts with current challenges regarding safety assurance for “AI-based” systems.

Based on this comparison, we critically discuss the question: “What’s *really* different with AI?”, when it comes to assessing causes of risk related to artificial intelligence and system safety for automated vehicles (section III). We highlight the importance of providing answers to very specific questions, such as: *What risks are specific to the implementation of AI algorithms? How do these risks impact established assurance practices in the automotive industry?* As a partial answer to these questions, we discuss behavior-based aspects of safety assurance that help establish engineering rigor for autonomous systems which operate in an open context (section IV).

The paper will conclude with a set of recommendations and a call to action when considering the identified commonalities regarding system safety for “AI-based” systems (section VI). The recommendations are meant to provide a *starting point* to establish rigorous traceability from system-level safety aspects to AI-specific safety aspects that is currently e.g. demanded by ISO/PAS 8800:2022 [4] without providing specific methodical guidance. This paper will provide according high-level suggestions in section IV. Further research initiatives on their refinement and validation will still be required.

## II. THE TROUBLES WITH AI DEFINITIONS

For dealing with AI-related assurance challenges, a clear definition of what is considered as an “AI-based” system is paramount. Regarding regulation and standardization, varying definitions may influence the scope of regulations and standards. The differing notions of *artificial intelligence* in current standards and regulation may hence pose hurdles when navigating the current standardization and regulation landscape.

The EU AI Act (Regulation (EU) 2024/1689, [7]) defines an *AI System* as “[...] a machine-based system that is designed to

operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.” This is an extremely broad and partially blurry definition: The degree of autonomy or adaptiveness remain undefined, while any sufficiently complex rule-based system can be considered to infer predictions or make decisions based on its inputs.

Similar, broad definitions are adopted by several ISO standards related to “AI Systems”: ISO/IEC 22989 [8] follows similar concepts as the EU AI Act. The standard defines an “AI system” as an “engineered system that generates outputs such as content, forecasts, recommendations or decisions for a given set of human-defined objectives” [8, Def. 3.7]. ISO/IEC 5338 [9] and ISO/IEC 23894 [10] follow the same definition. Similar challenges as for the EU Act definition arise here, as “human-defined objectives” extend the definition to any type of system that follows human-defined rulesets.

ISO/IEC TR 24028 [11] and ISO 29119 [12] define artificial intelligence in a slightly more narrow fashion as the “capability of an engineered system to acquire, process and apply knowledge and skills” [11, Def. 3.4], with an *engineered system* being a “combination of interacting elements organized to achieve one or more stated purposes” [11, Def. 3.38]. *Knowledge* here refers to “facts, information [...] and skills acquired through experience or education” [11, Def. 3.4, NOTE 1]. This is, again, problematic, as a) a definition of *skills* is not given, and b) the applied knowledge definition attributes human-like properties to a technical system.

ISO/PAS 8800 [4], which is of direct relevance for the design and assurance of automated driving systems in the context of “AI-based systems”, adopts a restricted definition, as it defines an “AI system” as an “item or element that utilises one or more AI models” [4, Def. 3.1.17], where an AI model is a “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” [4, Def. 3.1.7].

The idea of defining an AI model as a construct that can infer results without being completely defined by human knowledge is helpful in so far as this definition provides a clear differentiating property: This comprises all data-driven (or “machine-learning-based”) algorithms, but also more old-fashioned expert systems which may perform logical inference on ontological knowledge-bases, while every system that generates predictions, data or information on a mere set of pre-defined rules would not fall under the given definition. For the remainder of the paper, we will use this definition to differentiate between “AI-based” (as in the definition of ISO/PAS 8800) and “classic” systems.

Given the challenges that come with the above-mentioned broad and differing definitions for “AI-based” systems, in the following, we will summarize related work that can help to scope assurance-related challenges for automated driving systems that contain AI models in the sense of ISO/PAS 8800.

<sup>1</sup>From a Safety Engineering perspective, systems that learn in an unsupervised fashion at runtime (“online methods” as per ISO/PAS 8800) represent a whole additional category of challenges. We will focus on systems trained offline for this paper. ISO/PAS 8800 explicitly states [4, p. 145] that online methods for retraining to e.g. learn distribution shift can be subject to “other requirements”.

### III. MAPPING OUT ASSURANCE CHALLENGES

In the following, we will provide a literature overview in three categories: First, we will give an overview of assurance challenges that are related to risks emerging from the fact that automated vehicles are complex systems exposed to an open context. Second, we will compare and contrast the collected challenges with (alleged) challenges that can be found in selected state-of-the-art approaches for the assurance of AI-based systems. Finally, we will provide additional context for the arguments in section IV by addressing particular (safety) assurance challenges related to defining safe system behavior that are relevant in all aforementioned cases.

#### A. Selected Assurance Challenges related to an Open System Context

As discussed in the introduction, assurance challenges related to automated driving systems are closely related to uncertainty that comes with the requirement to operate in an open world. Assuring system properties (i.e., providing “grounds for justified confidence that a *claim* [...] has been or will be satisfied” [13, Def. 3.1.1]) such as safety becomes challenging: Not all situations that the engineered system will encounter over its lifetime can be foreseen by the developers at design time.

Burton and Herd [2] give a visual explanation of uncertainty-related challenges (cf. Fig. 1a): A technical *system* perceives the *environment* through *observations* and infers decisions or actions based on these observations. The environment is subject to the aforementioned open-context-related uncertainty (e.g., unpredictability, incomplete knowledge about the world, inherent complexity). The observations that a technical system makes are subject to technology-related uncertainty (e.g., measurement noise, limited fields of view, resolution limitations). The system is subject to uncertainty, as it is complex in itself, must infer decisions based on incomplete information, and may apply non-deterministic algorithms. [2, p. 04]

Eventually, according to Burton and Herd [2], the aforementioned manifestations of uncertainty result in *assurance uncertainty*, i.e., in uncertainty about the confidence that a certain claim has been or will be satisfied. In the context of safety assurance (i.e., the process of assuring a claim about the “absence of unreasonable risk” [14, Part 1, Def. 3.132] for a system), *risk* is strongly related to uncertainty (according to ISO 31000, risk is the *effect of uncertainty on objectives* [15, Def. 3.1]). Hence, the assurance uncertainty according to [2] can be rephrased in terms of a residual risk (cf. Fig. 1a) that can be mitigated, but not fully eliminated, as at least the uncertainty related to the open world can never be eliminated.

For AV safety assurance, it must be proven that the residual risk is *acceptable* or *reasonable* with respect “to valid societal or moral concepts” [14, Part 1, Def. 3.176]<sup>2</sup>. Strategies that

can be used to mitigate such residual risk below such a threshold must target the root causes of uncertainty. These include methods for obtaining an understanding of the real world, e.g. by diligently analyzing technological requirements for data acquisition, or generating sufficient knowledge about the systems’ internal and external behavior (cf. section IV).

Finally, the uncertainty related to the open context entails the need for diligent field monitoring [1]. Due to the boundaries of human knowledge, developers cannot foresee every possible situation that a system can encounter in the field. If any evidence is found that development decisions do not hold and that this leads to an underestimation of the risk that the system poses to surrounding road users, this must trigger reevaluation processes and a rollout of required system updates. In severe cases, a good safety culture demands that these processes can ground the fleet until the issue is mitigated (cf. [18, Sec. 8.3.4.1, pp. 286f.]).

#### B. Selected Assurance Challenges for AI-based Systems

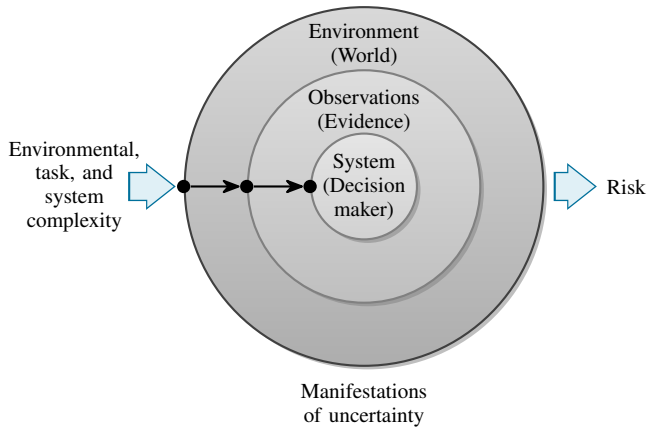
Assurance challenges for AI-based (as per ISO/PAS 8800) systems have been frequently discussed in the literature [2, 5, 19–21]. Main themes that can be identified in all related publications are: The *black box* character of AI models when compared to classic model-based white or gray box systems, *robustness* challenges regarding consistent outputs under minor input variations as well as *explainability* challenges.

These challenges are not fully independent of each other, but are closely connected to the black box property mentioned above: The AI models that constitute an AI-based system according to ISO 8800 are highly complex and nonlinear estimators, consisting of up to trillions of interacting parameters. This presents a degree of complexity that is, without applying proxy methods, not human comprehensible and turns the inner workings of those estimators into a factual black box for engineers. This hinders the structural assurance of AI elements, compared to the assessment of system architectures, as it is the state of the art for “non-AI” system elements. Further, the black box character leads to a lack of explainability [5, 20] regarding *why* an AI model generates an output given a certain input.

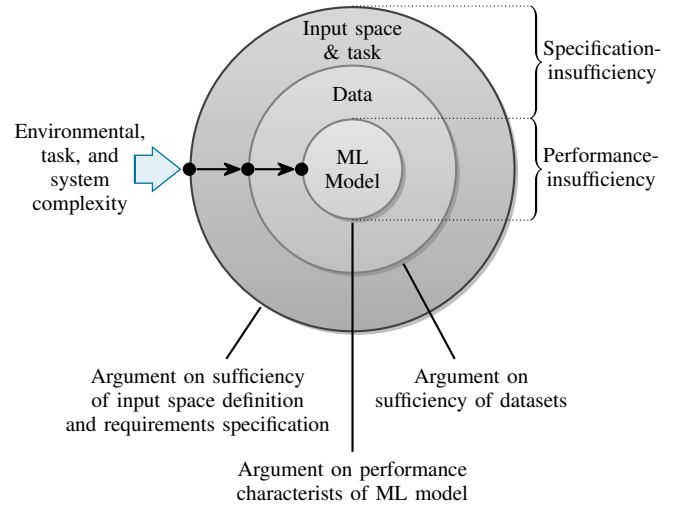
Regarding robustness challenges, the high nonlinearity in the AI models makes up the core strength of machine learning approaches: Given a suitable architecture, sufficient training and a sufficient number of parameters, any nonlinear input-output mapping can be approximated sufficiently close [22]. However, these nonlinearities can lead to significant changes in the output data (classifications, estimations, etc.) given only slight variations of the input data. This is the reason for the discussed robustness challenge [20] – and presents the main challenge for assuring properties regarding the input-output relations of AI models.

In this respect, very specific research questions must be answered, to tackle the aforementioned challenges, such as: “How can the translation from latent feature representations in black box models to explainable feature representations

<sup>2</sup>This is a logical inversion of the definition for *unreasonable* risk as per ISO 26262-1:2018 [14, Def. 3.176]. Many technical standards such as ISO 26262 [14] interpret risk in relation to the occurrence of *physical* harm. However, broader definitions of harm can include harm of stakeholder values [3, 16], including e.g. flawed interactions with emergency responders [17].



(a) General sources of uncertainty, related to an open world. Adapted from [2]: Modified by replacing assurance uncertainty with risk.



(b) Specific sources of uncertainty that cause assurance challenges related to AI-based systems. Redrawn from [2].

Fig. 1: Conceptual representations of uncertainty impacting systems (generic and AI-specific) operating in an open context according to [2].

be improved and robustified?” (cf, “Explainable AI” [23]) or: “How can the robustness of AI models be increased?”

Burton and Herd [2] discuss different aspects of uncertainty which are closely related to AI component characteristics. For this, the authors establish a corresponding model (Fig. 1b) to the one which they use for discussing the impact of uncertainty on general complex autonomous systems (Fig. 1a). According to [2], the uncertainty related to the open world requires a diligent definition of the input space of AI models – in our view, this requirement is no discriminator between AI-based systems and “classic” systems that must operate in an open context.

The inner rings in Fig. 1b are, on the contrary, specific AI-related challenges: While arguing safety over the sufficiency of datasets is closely related to a diligent understanding of the input space, properties of datasets, such as class labels, class distributions or coverage arguments on the sufficiency of data to cover the input space, are very specific to AI-based systems.

Further, Troubitsyna *et al.* [20] emphasize that safety is a system property. The authors argue that there is an urgent need for research that relates performance metrics for AI models to system-level safety metrics. Such a connection is required to enable assurance arguments on the performance characteristics of an ML model (inner circle in Fig. 1b) in a system context. This requires a close collaboration between safety and AI engineers.

Regarding the connection between system-level safety metrics and AI-specific performance metrics, ISO/PAS 8800 [4] assumes that safety analyses have been conducted and system-level requirements have been formulated, as a precondition for initiating the AI safety lifecycle. This underpins that further work is needed to connect both levels of abstraction.

Finally, Troubitsyna *et al.* [20] argue that more engineering

rigor must go into the design and assurance of AI-based systems. ISO/PAS 8800 [4] aims to address this challenge by providing process guidelines, including processes for a diligent monitoring of AI performance over the entire system lifecycle. This translates to life-cycle monitoring requirements, as discussed for the general operation of ADS in an open context (section III-A).

### C. Intermediate Conclusion

Koopman [24, slide 6] states that “Machine Learning breaks the Vee” in the sense that the application of data-based algorithms breaks engineering rigor. Traditionally, engineering rigor allows building a strong prior (in a Bayesian sense) that a system is safe, which can be confirmed or refuted by evidence generated in the verification and validation process [18].

Considering the argumentation in this paper, it should have become clear that it is not only machine learning what breaks the established verification and validation concepts in rigorous engineering processes. Major contributions to related assurance challenges stem from the uncertainty of an open system context. This uncertainty is the root cause of assurance challenges that affect ADS in general. In the Bayesian analogy given in [24], the challenge becomes to build well-formed and well-informed priors. In consequence, this means that, regardless of the type of system, methods are required which help to reduce uncertainty where this is possible (e.g., through diligent analyses of the operational environment and expected ADS behavior), and which enable a diligent documentation of development assumptions with respect to the properties of the operational context. Additionally, methods are required that allow to validate, verify, and, if required, update such assumptions over the entire system lifecycle [1, 20].

As ISO/PAS 8800 does explicitly not make the connection between system-level analyses or requirements and AI-specific

analyses, further guidance is required for a targeted exploitation of commonalities and to make efficient connections between system-level and AI-specific analyses.

#### D. Selected Assurance Challenges for Safe System-Level Behavior

The challenge of connecting system-level metrics to AI-specific performance metrics is further complicated by the fact, that effective system-level metrics for ADS (safety) assurance are still an area of active research [20, 25–28]. Besides the *definition* of these system-level metrics, their *traceability* through the design and verification & validation process remains a challenge. Specifically, the meaningful decomposition of system-level metrics to V&V pass-fail criteria remains challenging [21, 29, 30]. Note that these challenges are fully independent of how the system is implemented.

Specific concepts that contribute to establishing traceability are discussed in both, the AI community [31], and the safety assurance community [3, 21, 27, 30, 32]. Related concepts are the *Operational Design Domain* [31, 33], the required *system behavior* [30–32, 34] as well as *system capabilities* or *behavioral competencies* [3, 21, 29, 34].

In the context of safety assurance, the definition of the Operational Design Domain, a corresponding behavior specification as well as corresponding sets of required behavioral competencies are artifacts that allow to capture key assumptions about how the automated vehicle interacts with the open world. Such assumptions can be captured in terms of the elements which belong to the ODD, which characteristics of the ODD the vehicle needs to recognize, and how the vehicle should react to and interact with the ODD. A concept that allows to specify such interactions already at a very abstract level are sequences of *maneuvers* [30, 32, 34].

The definition of the ODD, the respective desired behavior<sup>3</sup>, and the corresponding behavioral competencies comes with its own set of challenges: A diligent definition of these artifacts requires a consequent derivation of behavioral requirements and relevant ODD elements from stakeholder needs. Business goals, but also normative sources such as laws, standards or societal expectation are examples of origins of stakeholder needs. For a traceable definition of artifacts and a documentation of justified assumptions, such normative sources, (semi-)formal representations of such sources are required [32].

## IV. BENEFITS OF BEHAVIOR-BASED SAFETY ANALYSES

The idea of making assumptions about the open context explicit by providing traceable definitions of the ODD, behavior, and behavioral competencies is a direct consequence of applying established systems engineering practice. In the following, we will summarize corresponding concepts that can help to initiate the safety life cycle for automated driving systems.

<sup>3</sup>ODD and behavior specification are highly iterative processes, as both can be mutual sources of constraints or extension needs.

TABLE I: Summary of terms used in section IV-B. Arrows should be read as a “subsumable under” relation.

Operational Level Terminology	
ADS Safety	Systems Engineering
Societal Expectations, Regulation,* Technical Standards*	Stakeholder Needs, Stakeholder Requirements
Use Case & Functional / Abstract Scenarios	Operational Use Cases & Scenarios
Operational Domain (OD), Target Operational Domain (TOD), Operational Design Domain (ODD)	Operational & System Context(s)
Behavioral Competencies & (System) (Cap)abilities	Capabilities

\* Note that, depending on the level of detail provided by a standard or regulatory document, these can also provide system or implementation-level requirements.

Operational Concept

#### A. Systems Engineering Concepts: Operational Domain & Operational Concept

Modern Systems Engineering puts a particular focus on diligent problem space analyses. The problem space is typically differentiated from the solution space. The problem space is focused on understanding what properties (or capabilities) a system needs. The solution space is focused on designing the system itself. In contrast to the solution domain, which would include technical and physical architectures and technology-specific solutions for realizing a system, it is good practice to keep problem-domain analyses solution and technology neutral.

From a Systems Engineering perspective, such considerations are summarized in the *operational concept*, which describes system characteristics and how the system shall be operated from the perspective of relevant stakeholders. As the focus is on system operation, i.e., its interaction with other entities and the world, the layer of abstraction that considers questions of stakeholder needs, use-cases, scenarios, system behavior, and required capabilities is called the *operational domain*<sup>4</sup>. [35]

#### B. Mapping to Automated Driving Systems

These Systems Engineering concepts can help to provide additional structure to the assurance challenges related to automated driving systems. When considering [2] and ISO/PAS 8800 [4], all analyses related to defining an operational concept contribute to the definition of a systems’ input space and task and capture assumptions regarding the open context. Table I provides a summary of the following steps.

1) *Capturing Stakeholder Needs & Requirements*: In Systems Engineering, the first step for defining an operational concept is the elicitation of stakeholder needs and the translation to actionable stakeholder requirements. For traceable ADS (safety) assurance, this implies that sources of stakeholder needs (laws, standards, societal expectations, etc.) must be captured and represented in a structured way. [32] gives examples of

<sup>4</sup>As the term can be confused with the ADS Operational Design Domain, we will only refer to the *operational concept* as an artifact in the following.



how such a structured representation of legal requirements for system behavior can be achieved.

2) *Defining Use-Cases & System Context*: Following the elicitation of stakeholder needs and requirements, use cases are formulated that present a first abstract representation of *what* the ADS should do. Use cases subsume sets of scenarios in similar settings [36].

To detail use cases, so-called system or operational contexts are defined, which specify those entities with which the system must interact. An (iterative) ODD definition determines, which entities can be part of a system context and which can, by that, be part of a use case and of a scenario.

Regarding AI-based systems, the system context is an important contribution to the input space definition: The elements of the system context are, e.g., sources of requirements for must-have class labels.

3) *Early Scenario-Based Safety Analyses and Behavior Specification*: The formulated use cases are typically detailed in scenarios. System-level safety analyses and the definition of desired system behavior can already be conducted in compliance with ISO 26262 [14] or ISO 21448 [37] by using functional (or abstract scenarios) without an immediate need to provide parameter ranges or choosing parameter values [32].

This enables the behavior-based formulation of safety goals and safety requirements, leaving open what kind of insufficiency actually causes a system-level hazardous event [38] (i.e. E/E failures [14], sensor-related performance insufficiencies [37], AI-related performance insufficiencies [4] or specification insufficiencies [37]). In other words, hazard and risk assessments can be performed on the behavior level. Safety-domain-specific analyses can follow, e.g. applying the failure mode models provided by ISO 26262, ISO 21448 or ISO/PAS 8800.

This concept follows the idea of ISO/PAS 8800, by explicitly separating between system-level and AI-specific analyses. The resulting behavior specification provides an artifact that can be used to trace back additional AI-specific requirements, e.g. in the form of additional data labels that are required to comply with the behavior specification.

4) *Definition of Behavioral Competencies*: The way that system capabilities or behavioral competencies have been described in the literature [29, 34, 39], they can provide a common starting point for decomposing system-level safety indicators into component-level performance metrics. Capabilities are Systems Engineering concepts to capture required potentials of a system to achieve an outcome with a specified performance [39]. The AVSC consortium [29] applies the concept of *behavioral competencies* similarly.

Behavioral competencies are related to a behavior specification, as they provide a way of reasoning *why* a system can show the required behavior. As the concept is a combination of a desired potential (requirement), outcome (behavior), and performance, the definition of behavioral competencies can provide guidance for a traceable definition of pass-fail criteria from system-level safety indicators, as well as the definition of meaningful performance criteria for AI components, as discussed in [20].

## V. EXAMPLE CASE STUDY

To illustrate the concepts discussed in this section, we will provide a short, hypothetical case study that largely connects previously published work. Specified safety goals, behavioral competencies, and requirements hence only serve illustrative purposes. They are not reflecting existing or developed services or products by the involved industry partners.

### A. Stakeholder Needs

For this case study, we assume stakeholder needs that require traffic-code-compliant behavior of the vehicle and a net risk that is “lower” than the risk caused by human drivers in comparable Operational Design Domains.

### B. Scenario and Context Definition

To discuss the connection between system-level and AI-specific safety analyses, we consider a scenario that has been discussed in [40]. Safety Goals and requirements at the behavioral level for this scenario have been specified in [3]: An automated vehicle approaches a row of parked vehicles. Between two of the parked vehicles, a pedestrian is about to step onto the street and into the path of the automated vehicle.<sup>5</sup> For the sake of the example, we assume that the ODD comprises straight inner-city roads with one lane in each driving direction and parking lanes parallel to these driving lanes. Pedestrians and stationary vehicles of various vehicle types are also included in the ODD.

For this single scenario, the *operational context* instantiates one pedestrian, a number of stationary vehicles, as well as the road with a parking lane as described above.

1) *Safety Analyses*: A *hazard analysis* conducted at the behavior level can yield the following results:<sup>6</sup>

- A *hazard* in the discussed scenario is the *potential injury of vulnerable road users*.
- The discussed scenario becomes a *hazardous scenario* in the SOTIF sense if the ego vehicle shows hazardous behavior in the given scenario. Such behavior can, e.g., be represented by a *follow lane maneuver* at an *inadequate speed* that prevents timely deceleration. Hence the maneuver can cause a collision with the pedestrian.
- A hazardous event that instantiates the hazard in the scenario is the *collision of the ego vehicle with the pedestrian*.

A *safety goal* that must be fulfilled in the scenario would be: *Collisions with vulnerable road users must be prevented*.

### C. SOTIF-Specific Analysis Results

Continuing the analysis in a SOTIF scope, it is possible to identify a *functional insufficiency*: the ego vehicle is *not able to correctly predict the behavior of occluded pedestrians*. The corresponding *specification insufficiency* is the missing

<sup>5</sup>A corresponding use case can be formulated as *passing parked vehicles*, comprising a suite of similar scenarios.

<sup>6</sup>For the case study, we omit a concrete risk assessment: exact numbers for assigned risks are irrelevant for showcasing the interplay between behavior-level and AI-specific analyses.

specification of *occluded areas* which can contain pedestrians and which can occur in the given Operational Design Domain. The resulting triggering condition would be the *presence of a pedestrian in an occluded area* in the given scenario.

Conducting these SOTIF analyses provides insights into how the risk of a collision with an occluded road user in the given scenario can be mitigated. Risk mitigation can be achieved by a) including occlusions in the specification of a world model and b) defining according behavioral competencies for the automated driving system. The AVSC, e.g. defines the behavioral competency *Responding to vulnerable road users (VRUs)* [29, p. 9]. Our scenario would extend this to *Responding to occluded vulnerable road users*. In the process, these behavioral competencies would be specified by according (first functional and later technical) safety requirements.

1) *Connecting to AI-specific Risks*: Focusing on the connection to AI-specific risks, we can compare the discussion of the scenario with Fig. 1a. In this example scenario the cause for the risk is an insufficient definition of the environment. This is, at this time, independent of any specific system implementation.

If we assume that environment perception and decision making are performed by ML models, Fig. 1b shows that AI-specific risks in the scenario stem from an insufficient definition of the input space for these ML models. One possible risk mitigation measure that can be implemented in this context is the inclusion of occlusion scenarios (and possibly even the labeling of occlusions) in the datasets used for training the ML models. In the sense of the ISO/PAS 8800 [4], this results in *dataset requirements* such as: *The dataset for training the ML model must contain labeled occluded areas*.

#### D. Discussion

The hypothetical case study is just a rough textual example of single aspects regarding the transfer of traditional safety engineering approaches to AI safety. All current standards emphasize the need for traceability between the specification, test results and safety performance indicators collected in the field. To establish this kind of traceability, rich metamodels connecting behavior and ODD-specification will be required.

First steps in this direction exist, e.g., with the A.U.T.O. ontology created in the VVMethods project [41] or our own previous work connecting Systems Engineering concepts with AD-Domain specific and behavior-related ontologies [32, 42]. However, a comprehensive approach that has been fully connected to AI-specific needs is still yet to be established.

Tooling that can reduce the effort of integrating model-based approaches into existing development processes will also be crucial: AI, function, and safety experts will need easily accessible interfaces to apply the ontologies and metamodels in their daily work without creating the burden of learning new description languages that do not target the core of their development activities.

## VI. CONCLUSIONS & RECOMMENDATIONS

This paper has emphasized the importance of identifying the reason for assurance challenges in AI-based automated driving

systems. We mapped out and highlighted the importance of differentiating which challenges are directly related to the nature of AI elements, and which are more related to the complexity of dealing safely with the open world that automated vehicle systems are deployed into.

The answer to the question in the title “*What’s really different with AI?*” can be summarized along the lines of: There are specific risks related to performance insufficiencies of AI-based system components (cf. ISO 8800 [4]). However, principles of engineering rigor, which stem from traditional safety and systems engineering, are a necessary condition for building safe AI-based systems.

In our view, this engineering rigor is rooted in diligent analyses for defining a system’s operational concept. This can include, e.g., capturing Operational Design Domains, finding traceable models for ADS behavior in traffic, and defining the required behavioral competencies supported by scenario-based safety analyses. These analyses can provide solid grounds for defining those system-level requirements and metrics as well as their decomposition to AI-related performance metrics, which is explicitly demanded by the current ISO 8800. Again, the behavior-based approach outlined above is a *starting point*. The ADS safety assurance community has recently provided ideas and guidance [3, 21, 28, 30–32] contributing to such a behavior-based perspective. Continuous efforts are required to further establish these concepts and consequently link them to AI assurance processes.

To summarize recommendations, we think that: a) It is crucial to assess, which benefits can be drawn from standards with broad AI definitions in what stage of the development process. b) A strong focus should be put on diligent problem space analyses, and on exploiting joint technology-neutral safety analyses as far as possible. This can provide an efficient, common starting point for SOTIF, Functional Safety, and AI safety analyses. c) Further research and standardization is needed for methods, ontologies and/or metamodels for the definition of behavior, as well as the definition of behavioral competencies. Finally, d), We would like to emphasize the conclusions by Troubitsyna *et al.* [20] that further research is required regarding the decomposition of system-level safety metrics to component-level and AI-related performance metrics.

#### ACKNOWLEDGEMENT

This paper is the product of ongoing discussions in the “Focus Field Safety and Risk” at the “German Round Table Autonomous Driving” initiated by the Federal Ministry on Digital and Transport. We would like to thank Steffen Müller, Stefan Liening, and their teams at the BMDV DK20 for providing the frame for these discussions. Further, we would like to thank the anonymous reviewers for their extremely constructive improvement suggestions, in particular regarding the case study and the reader guidance in section IV and the addition of a mini case study in section V.



## REFERENCES

- [1] S. Burton and J. A. McDermid, "Closing the Gaps: Complexity and Uncertainty in the Safety Assurance and Regulation of Automated Driving," Technical Rep. 2023.
- [2] S. Burton and B. Herd, "Addressing Uncertainty in the Safety Assurance of Machine-Learning," *Frontiers Comput. Sci.*, vol. 5, Apr. 6, 2023. DOI: 10.3389/fcomp.2023.1132580.
- [3] M. Nolte, "Werte- und fähigkeitsbasierte Bewegungsplanung für autonome Straßenfahrzeuge – Ein systemischer Ansatz," (in German), Ph.D. dissertation, TU Braunschweig, 2025.
- [4] *Road Vehicles — Safety and Artificial Intelligence*, ISO Publ. Avail. Spec. 8800:2024.
- [5] R. Ashmore, R. Calinescu, and C. Paterson, "Assuring the Machine Learning Lifecycle: Desiderata, Methods, and Challenges," *ACM Comput. Surv.*, vol. 54, no. 5, 111:1–111:39, May 25, 2021. DOI: 10.1145/3453444.
- [6] R. Schnitzer, L. Kilian, S. Roessner, K. Theodorou, and S. Zillner, *Landscape of AI Safety Concerns – A Methodology to Support Safety Assurance for AI-Based Autonomous Systems*, Dec. 18, 2024. DOI: 10.48550/arXiv.2412.14020. arXiv: 2412.14020[cs].
- [7] *REGULATION (EU) 2024/1689 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL*, Aug. 26, 2022.
- [8] *Information Technology — Artificial Intelligence — Artificial Intelligence Concepts and Terminology*, ISO/IEC Standard 22989:2020.
- [9] *Information Technology — Artificial Intelligence — AI System Life Cycle Processes*, ISO/IEC Standard 5338:2023.
- [10] *Information Technology — Artificial Intelligence — Guidance on Risk Management*, ISO/IEC Standard 23894:2023.
- [11] *Information Technology — Artificial Intelligence — Overview of Trustworthiness in Artificial Intelligence*, ISO/IEC Tech. Rep. 24028:2020.
- [12] *Software and Systems Engineering — Software Testing*, ISO/IEC Tech. Rep. 29119:2020.
- [13] *Systems and Software Engineering — Systems and Software Assurance — Part 1: Concepts and Vocabulary*, ISO/IEC/IEEE Standard 15026-1:2023.
- [14] *Road Vehicles — Functional Safety*, ISO Standard 26262:2018.
- [15] *Risk Management — Guidelines*, ISO Standard 31000:2018.
- [16] N. F. Salem, S. Le Page, J. Millar, P. Junietz, M. Nolte, R. Graubohm, and M. Maurer, "Safety and Risk – Why Their Definitions Matter," in *Handbook Assisted Automated Driving*, H. Winner, K. Dietmayer, L. Eckstein, M. Jipp, M. Maurer, and C. Stiller, Eds., 4th ed., Heidelberg: Springer Nature, 2025, (in press).
- [17] P. Koopman and W. Widen, "Redefining Safety for Autonomous Vehicles," in *2024 Int. Conf. Comput. Saf., Rel., Secur. (SAFECOMP)*, A. Ceccarelli, M. Trapp, A. Bondavalli, and F. Bitsch, Eds., ser. Lecture Notes Comput. Sci. Vol. 14988, Florence, Italy: Springer, Cham, pp. 300–314. DOI: 10.1007/978-3-031-68606-1\_19.
- [18] P. Koopman, *How Safe Is Safe Enough? Measuring and Predicting Autonomous Vehicle Safety*, 1st ed. Pittsburgh, PA: Carnegie Mellon Univ., 2022, 352 pp.
- [19] K. Czarnecki and R. Salay, "Towards a Framework to Manage Perceptual Uncertainty for Safe Automated Driving," in *2018 Int. Conf. Comput. Saf., Rel., Secur. (SAFECOMP)*, B. Gallina, A. Skavhaug, E. Schoitsch, and F. Bitsch, Eds., ser. Lecture Notes Comput. Sci. Västerås, Sweden: Springer, Cham, pp. 439–445. DOI: 10.1007/978-3-319-99229-7\_37.
- [20] E. Troubitsyna, I. J. Alvarez, P. Koopman, and M. Trapp, "Methods and Tools for the Engineering and Assurance of Safe Autonomous Systems," Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2024, ISSN: 2192-5283 Issue: 4 Volume: 14, pp. 23–41. DOI: 10.4230/DAGREP.14.4.23.
- [21] G. Stettinger, P. Weissensteiner, and S. Khastgir, "Trustworthiness assurance assessment for high-risk AI-based systems," *IEEE Access*, vol. 12, pp. 22 718–22 745, 2024. DOI: 10.1109/ACCESS.2024.3364387.
- [22] K. Hornik, M. Stinchcombe, and H. White, "Multilayer Feedforward Networks Are Universal Approximators," *Neural Netw.*, vol. 2, no. 5, pp. 359–366, Jan. 1989. DOI: 10.1016/0893-6080(89)90020-8.
- [23] K. Sokol and P. Flach, *Interpretable Representations in Explainable AI: From Theory to Practice*, Dec. 23, 2023. DOI: 10.48550/arXiv.2008.07007. arXiv: 2008.07007[cs,stat].
- [24] P. Koopman, "L145 Challenges in Autonomous Vehicle Safety Assessment," US DOT Workshop, US DOT Workshop, online, Mar. 2024.
- [25] C. Neurohr, L. Westhofen, M. Butz, M. H. Bollmann, U. Eberle, and R. Galbas, "Criticality Analysis for the Verification and Validation of Automated Vehicles," *IEEE Access*, vol. 9, pp. 18 016–18 041, 2021. DOI: 10.1109/ACCESS.2021.3053159.
- [26] Y.-H. Chen, J. M. Scanlon, K. D. Kusano, T. L. McMurtry, and T. Victor, *Dynamic Benchmarks: Spatial and Temporal Alignment for ADS Performance Evaluation*, Oct. 11, 2024. DOI: 10.48550/arXiv.2410.08903. arXiv: 2410.08903[cs].
- [27] R. Galbas, M. Nolte, U. Eberle, H. Hungar, H. Mosebach, N. F. Salem, H. Schittenhelm, J. Reich, T. Kirschbaum, L. Westhofen, PEGASUS VVM Consortium, R. Galbas, M. Nolte, U. Eberle, H. Hungar, H. Mosebach, N. F. Salem, H. Schittenhelm, J. Reich, T. Kirschbaum, and L. Westhofen, "VV Methods Safety Assurance Position Paper," Zenodo, Jun. 15, 2024. DOI: 10.5281/ZENODO.11669422.
- [28] L. Fraade-Blanc, F. Favarò, J. Engstrom, M. Cefkin, R. Best, J. Lee, and T. Victor, *Being Good (at Driving): Characterizing Behavioral Expectations on Automated and Human Driven Vehicles*, Feb. 12, 2025. DOI: 10.48550/arXiv.2502.08121. arXiv: 2502.08121[cs].
- [29] Automated Vehicle Safety Consortium (AVSC), "AVSC Best Practice for Evaluation of Behavioral Competencies for Automated Driving System Dedicated Vehicles (ADS-DVs)," AVSC00008202111, Nov. 2021.
- [30] *Behaviour Taxonomy for Automated Driving System (ads) Applications – Specification*, BSI Standard BSI Flex 1891 v1.0:Jan. 2025.
- [31] G. Price, "A Behavioural Safety Centric Approach for E2E ADS," Presentation, The26262Club, online, Mar. 2025.
- [32] N. F. Salem, M. Nolte, V. Haber, T. Menzel, H. Steege, R. Graubohm, and M. Maurer, "An Ontology-Based Approach Toward Traceable Behavior Specifications in Automated Driving," *IEEE Access*, vol. 12, pp. 165 203–165 226, 2024. DOI: 10.1109/ACCESS.2024.3494036.
- [33] P. Weissensteiner, G. Stettinger, S. Khastgir, and D. Watzel, "Operational Design Domain-Driven Coverage for the Safety Argumentation of Automated Vehicles," *IEEE Access*, vol. 11, pp. 12 263–12 284, 2023. DOI: 10.1109/ACCESS.2023.3242127.
- [34] M. Nolte, G. Bagschik, I. Jatzkowski, T. Stolte, A. Reschka, and M. Maurer, "Towards a Skill- and Ability-Based Development Process for Self-Aware Automated Road Vehicles," in *2017 IEEE Int. Conf. Intell. Transp. Syst. (ITSC)*, Yokohama, Japan: IEEE, pp. 739–744. DOI: 10.1109/ITSC.2017.8317814.
- [35] D. D. Walden and International Council on Systems Engineering, Eds., *INCOSE Systems Engineering Handbook*, 5th ed., Hoboken, NJ: John Wiley Sons Ltd, 2023.
- [36] S. Ulbrich, A. Reschka, T. Menzel, F. Schuldt, and M. Maurer, "Defining and Substantiating the Terms Scene, Situation and Scenario for Automated Driving," in *2015 18th IEEE Int. Annu. Conf. Intell. Transp. Syst. (ITSC)*, Las Palmas, Spain: IEEE, pp. 982–988.
- [37] *Road Vehicles — Safety of the Intended Functionality*, ISO Standard 21448:2022.
- [38] N. F. Salem, T. Kirschbaum, M. Nolte, C. Lalitsch-Schneider, R. Graubohm, J. Reich, and M. Maurer, "Risk management core – towards an explicit representation of risk in automated driving," *IEEE Access*, vol. 12, pp. 33 200–33 217, 2024, tex.publisher: IEEE. DOI: 10.1109/ACCESS.2024.3372860.
- [39] C. S. Wasson, *System Engineering Analysis, Design, and Development: Concepts, Principles, and Practices*. Hoboken, NJ: John Wiley Sons Inc, 2005.
- [40] R. Graubohm, N. F. Salem, M. Nolte, and M. Maurer, "On Assumptions with Respect to Occlusions in Urban Environments for Automated Vehicle Speed Decisions," in *2023 IEEE 26th Int. Conf. Intell. Transp. Syst. (ITSC)*, citation-key: graubohm2023, Bilbao, Spain: IEEE, pp. 738–745. DOI: 10.1109/ITSC57777.2023.10422457.
- [41] L. Westhofen, C. Neurohr, M. Butz, M. Scholtes, and M. Schuldes, "Using Ontologies for the Formalization and Recognition of Criticality for Automated Driving," *IEEE Open J. Intell. Transp. Syst.*, vol. 3, pp. 519–538, 2022. DOI: 10.1109/OJITS.2022.3187247.
- [42] M. Nolte and M. Maurer, "Towards Closing the Gap between Model-Based Systems Engineering and Automated Vehicle Assurance: Tailoring Generic Methods by Integrating Domain Knowledge," presented at the 16. Uni-DAS e.V. Workshop Fahrerassistenz und automatisiertes Fahren, Irsee: Uni-DAS e.V., 2025.