



Requirements and software engineering for automotive perception systems: an interview study

Khan Mohammad Habibullah¹ · Hans-Martin Heyn¹ · Gregory Gay¹ · Jennifer Horkoff¹ · Eric Knauss¹ · Markus Borg² · Alessia Knauss³ · Håkan Sivencrona³ · Polly Jing Li⁴

Received: 4 July 2023 / Accepted: 3 December 2023 / Published online: 24 January 2024

This is a U.S. Government work and not under copyright protection in the US; foreign copyright protection may apply 2024

Abstract

Driving automation systems, including autonomous driving and advanced driver assistance, are an important safety-critical domain. Such systems often incorporate perception systems that use machine learning to analyze the vehicle environment. We explore new or differing topics and challenges experienced by practitioners in this domain, which relate to requirements engineering (RE), quality, and systems and software engineering. We have conducted a semi-structured interview study with 19 participants across five companies and performed thematic analysis of the transcriptions. Practitioners have difficulty specifying upfront requirements and often rely on scenarios and operational design domains (ODDs) as RE artifacts. RE challenges relate to ODD detection and ODD exit detection, realistic scenarios, edge case specification, breaking down requirements, traceability, creating specifications for data and annotations, and quantifying quality requirements. Practitioners consider performance, reliability, robustness, user comfort, and—most importantly—safety as important quality attributes. Quality is assessed using statistical analysis of key metrics, and quality assurance is complicated by the addition of ML, simulation realism, and evolving standards. Systems are developed using a mix of methods, but these methods may not be sufficient for the needs of ML. Data quality methods must be a part of development methods. ML also requires a data-intensive verification and validation process, introducing data, analysis, and simulation challenges. Our findings contribute to understanding RE, safety engineering, and development methodologies for perception systems. This understanding and the collected challenges can drive future research for driving automation and other ML systems.

Keywords Requirements engineering · Software quality · Software development methodologies · Driving automation systems · Autonomous driving

✉ Khan Mohammad Habibullah
khan.mohammad.habibullah@gu.se

✉ Gregory Gay
greg@greggay.com

✉ Jennifer Horkoff
jennifer.horkoff@gu.se

Hans-Martin Heyn
hans-martin.heyne@gu.se

Eric Knauss
eric.knauss@gu.se

Markus Borg
markus.borg@cs.lth.se

Alessia Knauss
alessia.knauss@zenseact.com

Håkan Sivencrona
hakan.sivencrona@zenseact.com

Polly Jing Li
polly.jing.li@kognic.com

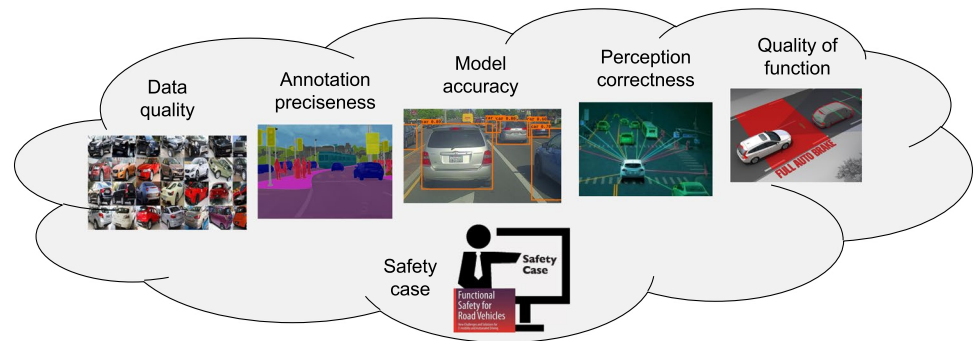
¹ Department of Computer Science and Engineering,
Chalmers and University of Gothenburg, Gothenburg,
Sweden

² Department of Computer Science, Lund University, Lund,
Sweden

³ Zenseact AB, Gothenburg, Sweden

⁴ Kognic AB, Gothenburg, Sweden

Fig. 1 Conceptual model of quality transitions from data collection to the quality of the automotive function



1 Introduction

Driving automation systems, including both autonomous driving (AD) and advanced driver assistance systems (ADAS), are software systems designed to augment or automate aspects of vehicle control [1]. Driving automation systems have long been a domain of interest. However, the increased capabilities and usability of machine learning (ML) have subsequently improved the capabilities of—and interest in—such systems. Research advances have improved comfort and safety, and reduced fuel and energy consumption, emissions, and travel time [1].

Driving automation system functionality depends on the correctness and the integrity of perception systems that blend ML-based models and traditional signal processing¹. The usage of ML for perception relies on a large quantity and high quality of data. Data quality, context, and attributes—as well as annotation quality—have a significant impact on the resulting system quality. However, it is difficult to make direct connections between data, annotation, ML model quality, and the resulting functional quality of a perception system (e.g., between the boxes in Fig. 1). The inherent uncertainty of ML—coupled with high requirements on data quality and coverage—creates substantial requirements, systems, and software engineering challenges in perception system development [2].

Requirements engineering (RE) is an important foundational element of quality assurance and safety engineering. RE plays a critical role in perception system development by enabling explicit capture of safety and quality requirements, supporting communication, recording functional expectations, and ensuring that standards are followed. Additionally, systems and software engineering play a critical role in the successful development and deployment of perception systems by enhancing real-time decision-making [3], supporting adaptability and continuous learning, facilitating

complex system integration [4], maximizing performance, ensuring dependability and safety [5], encouraging cross-disciplinary collaboration [6], and advancing ethical and responsible development methods [7].

Recent research has explored RE challenges for ML systems, e.g., [8, 9], as well as systems and software engineering challenges [10–12]. However, such challenges have not been thoroughly explored in the context of perception systems for driving automation systems. Addressing this gap is necessary to advance practices in both this domain and in the broader context of RE for ML systems.

To explore important engineering topics and challenges for perception systems, we have conducted an interview study with 19 domain experts from five companies working in various driving automation systems roles. We analyzed interview data using thematic coding to produce eight major themes: perception, requirements engineering, systems and software engineering, AI and ML models, annotation, data, ecosystem and business, and quality.

This paper is an extension of previous work [13]. The initial article focused specifically on the RE themes from the thematic analysis, encapsulating RE topics and challenges discussed by the participants. In this paper, we extend the analysis and discussion of the RE theme to include findings from two additional themes—systems and software engineering and quality. For both themes, we also explore topics and challenges for driving automation systems development that were raised in the interviews.² These two themes, in particular, add relevant insights for practitioners and, additionally, enrich our understanding of RE practices and challenges in this domain (e.g., requirements and quality are tightly interconnected). In addition, we include a more extensive related work and discussion section, including an outline of future directions in research and practice for driving automation systems and other ML systems.

¹ In this paper, we focus specifically on ML-based perception systems for driving automation systems, but often use the term *perception systems* as shorthand.

² Another recent article has also used the same interview data, but focused on the annotation, data, and ecosystems and business themes [14].

Related to RE, our findings indicate that practitioners have difficulty breaking down specifications for the ML components. In practice, individuals report that they use scenarios, operational design domains (ODDs), and simulations as part of RE. Practitioners experience RE challenges related to uncertainty, ODD detection, realistic scenarios, edge case specification, traceability, creating specifications for data and annotations, and quantifying quality requirements.

In terms of quality, practitioners consider performance, reliability, robustness, safety, and user comfort as important quality attributes. In the context of driving automation systems, safety is particularly critical. Practitioners establish safety goals, often in negotiation with component suppliers. To ensure safety, practitioners must comply with evolving safety and AI standards—which are challenging and costly to meet. They must also manage trade-offs between safety and other qualities. Safety cases are a critical element of ensuring that the safety goals are met. Quality assurance is performed by tracking critical key performance indicators (KPIs) during the execution of catalogs of scenarios. Quality assurance is complicated by the non-determinism and data requirements of ML and the realism of simulation.

From a systems and software engineering perspective—though practitioners work with traditional and agile methods—ML complicates the overall development process. Current agile methods are insufficient for the needs of large-scale ML because practitioners lack appropriate data quality methods as part of their overall development methodology [15]. Furthermore, the addition of ML leads to a data-intensive verification and validation (V & V) process with challenges related to data quality, statistical analysis, and simulation.

By exploring the views and challenges of practitioners on RE, quality, and software and systems engineering for ML-enabled perception systems, we provide valuable insights for practitioners working in this safety-critical domain. Additionally, our findings contribute to improving RE, and systems and software engineering knowledge more broadly, in other domains reliant on ML.

2 Related work

In this section, we review related work in requirements engineering for machine learning systems and for automotive and driving automation systems. We give an overview of work on quality for machine learning and software development methods for machine learning, as these are key themes of focus in our interview study.

2.1 Requirements engineering for machine learning

Recent research has focused on how RE could or must change in the face of rising use of ML. Systematic mapping studies on RE for ML identified new contributions in this area, including approaches, checklists, guidelines, quality models, classifications and evaluations of quality models, taxonomies, and quality requirements [16–18]. Pei et al. reviewed literature on RE for ML, went through a collaborative requirements analysis process, and provided an overview of RE processes for ML applications in terms of cross-domain collaboration [19]. They provided an example case of an industrial data-driven intelligence application, discussed in relation to the provided requirements analysis process. Ahmad et al. performed a systematic mapping study to find articles on current RE for AI approaches and identified available frameworks, methodologies, tools, and techniques used to model requirements and found existing challenges and limitations [20]. They identified 43 primary studies and found several challenges and limitations of existing RE for AI practices, for example, that current RE processes are not adequately adaptable for building AI systems. The authors emphasized that new techniques and tools are needed to support RE for AI.

Further papers have identified RE-related challenges for ML and AI. A high number of AI solutions fail or do not make it to production due to missing or bad RE processes. Therefore, Maalej et al. discussed six aspects that need careful consideration and tailoring to the AI context that include acceptable levels of quality requirements, data- and user-centered prototyping, expanding RE to focus on data, embedding responsible AI terminology into the engineering workflows, trade-off analysis for responsible AI, and requirements as foundation for quality and testing of AI [21]. Gjorgjevikj et al. discussed the challenges in applying conventional RE practices to ML systems and proposed best practices and adjustments to RE concepts [22].

Ahmad et al. investigated current approaches for writing requirements for AI/ML systems, identified tools and techniques to model requirements for AI/ML, and pointed out existing challenges and limitations in this area [23]. Belani et al. identified and discussed RE challenges for ML- and AI-based systems and reported that identifying NFRs throughout the software lifecycle is one of the main challenges [8]. Heyn et al. used three use cases of distributed deep learning to describe AI system engineering challenges related to RE [24], including context, defining data quality attributes, human factors, testing, monitoring, and reporting. In further study, Heyn et al. identified several challenges related to training data specification (e.g., unclear design domain, missing guidelines for data selection, and unsuitable safety standards) and run-time monitoring for ML models,

with challenges relating to RE (e.g., lack of explainability for ML decisions, missing conditions for run-time checks, and overhead for monitoring solution) [25].

Other studies move toward proposed solutions. Farrell et al. identified key characteristics of ML-based software requirements such as confidence, accuracy, average value, robustness, data-driven learning, and quality aspects, providing a foundation for developing a taxonomy of requirements for such software [26]. Villamizar et al. propose a catalogue of 45 concerns that should be considered in specifying ML systems, covering five different relevant perspectives for such systems, such as objectives, user experience, infrastructure, model, and data [27]. Islam et al. presented a requirements process (RESAM) that integrates knowledge from different sources, such as discussion forums, domain experts, and formal product documentation, to discover and specify requirements and design definitions that contribute to the construction of effective deep learning anomaly detectors. They evaluated their process in a case study and demonstrate that it guides the construction of effective anomaly detection models that support explainability [28].

2.2 Requirements engineering for automotive and driving automation systems

Significant research has been performed on RE for vehicles. Liebel et al. identified challenges in automotive RE with respect to communication and organization structure [29]. Pernstl et al. stated that RE is one of the areas most in need of improvement at automotive original equipment manufacturers (OEMs) and also identified the ability to communicate via requirements as important [30]. Allmann et al. also noted requirements communication as a major challenge for OEMs and their suppliers [31]. Mahally et al. identified that requirements are the main enablers and barriers of moving toward Agile for automotive OEMs [32].

Research has also looked specifically at RE for AD, e.g., providing an overview of AD RE techniques [33]. Riberio et al. identified AD RE challenges addressed by the literature and identified the languages and description styles used to describe AD requirements, with special attention given to NFRs [34]. Heyn et al. investigated challenges with context and ODD definition in ML-enabled perception systems [35], including a lack of standardization for context definitions, ambiguities in deriving ODDs, missing documentation, and lack of involvement of function developers while defining the context. Ågren et al. identified six aspects of RE that impact automotive development speed, moving toward AD [36].

In further driving automation systems work relating to RE, Zhang et al. conducted a systematic mapping study in the context of driving automation systems and introduced a taxonomy for critical scenario identification methods

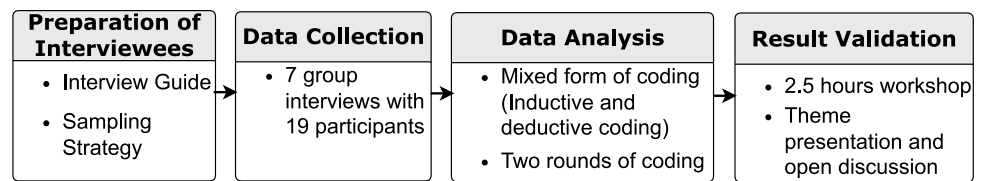
including encompassing the problem definition of the solution and the assessment of the established scenarios [37]. They also discussed challenges considering the perspectives of coverage, practicability, and scenario space explosion. Luo et al. proposed a hierarchical safety assessment approach to quantitatively analyze the quality trade-offs, violation severity of safety requirements, and distinguish safer autonomous driving systems configurations based on the requirements violations comparison in a hierarchical way, following requirements importance [38]. Zhang et al. presented a data-driven engineering process that includes hierarchical requirements engineering to link the operational design domain with the requirements and semi-automated generation of datasets for leveraging future application of ML in automated driving in industry [39].

2.3 Quality assurance for machine learning

Although quality for ML can be interpreted in a narrow sense, i.e., basic model performance, work exists which has focused on ML quality in a broader sense. Felderer et al. discussed terminology for quality assurance for AI systems, defining concepts and characterizing AI systems into artifact type, process, and quality characteristics [40]. They also discussed challenges in quality assurance such as lack of specifications and defined requirements; the need for validation data and test input generation; difficulty defining expected outcomes as test oracles; and baselines for AI-based systems. Furthermore, different challenges and opportunities related to quality requirements for machine learning systems are reported and discussed in [8, 41–43].

From the perspective of traditional quality assurance, the Japanese industry has collectively proposed a set of recommendations for the quality assurance of AI systems (e.g., in the Consortium of Quality Assurance for AI-based Products and Services) and the second iteration of these standards, which includes a list of quality evaluation criteria, a list of cutting-edge methods, and explanations of each of the five representative domains that are proposed in [44]. The research project PEGASUS (Project for the Establishment of Generally Accepted quality criteria tools and methods as well as Scenarios and Situations) focused on the release of highly automated driving functions. In this project, 17 partners from research and industry worked together with the aim to develop a complete toolchain to include criteria and measures for the evaluation of functions and for driving automation systems quality levels, with test catalogues, central methods for driving automation systems development, and processes for establishing safety, and to release highly automated driving functions [45].

From the perspective of technical standards, the automotive industry is aware of adjustments in machine learning-based technology demands in terms of technical expertise,

Fig. 2 Overview of the interview study

development paradigms, and cultural approaches. However, there is still a significant gap between the availability of technical standards and certification capacity. Currently, the automotive industry is governed by several standards. However, existing work has argued that these standards are not suitable for machine learning-based driving automation systems [46, 47]. Although further certifications of autonomous systems (e.g., SOTIF) are developed and are advancing, these efforts only cover some of the existing challenges [48].

Other work has focused specifically on data quality in relation to machine learning. Jain et al. discussed the importance of data quality and stated that the effort required to iteratively debug a machine learning pipeline in order to enhance model performance can be reduced by evaluating the quality of the data using intelligently defined metrics transformation operations [49]. The authors also survey the important data quality related approaches discussed in literature—highlighting their strengths and similarities and discussing their applicability to real-world challenges.

Further work has looked at AI in terms of risks. Poth et al. presented a systematic methodical approach (the evAIA method evaluates AI approaches) that evaluates risks of the machine learning model using a questionnaire specifically for AI products and services [50].

2.4 Software and systems methods for machine learning

Current systems and software development methods often do not account well for machine learning-enabled systems. Giray points out a lack of techniques to support machine learning system development as part of a systematic literature review, reporting that the non-deterministic nature of machine learning systems complicates SE aspects of engineering machine learning systems that include a lack of mature tools and techniques to support machine learning systems development and verification [51].

Given the rise of machine learning-enabled software, researchers have explored or introduced a number of methods and challenges for machine learning and AI system development. Hesenius et al. provided a structured engineering process framework named EDDA (engineering data driven applications) that bridges existing gaps, supports data-driven application development, and ensures the required quality levels for critical components of machine

learning systems [52]. Amershi et al. conducted a case study where the authors described how various Microsoft software teams developed software applications with customer-focused AI features—integrating existing Agile software engineering process with AI-specific workflows [12].

Further research looked into the challenges of engineering driving automation systems. Key collaboration challenges were identified in developing and deploying machine learning systems through interviews with 45 participants from 28 organizations [53]. The authors reported on common collaboration points and challenges from the perspective of requirements, data, integration, and team patterns and found the majority of the challenges center around communication, documentation, engineering, and process. In addition, safety criticality extends the decision-making, development, and related environmental perception [54]. This complexity does not harmonize with conventional safety engineering; hence, the application of concepts for intelligence is required to resolve the complexity.

3 Methodology

Our study is guided by the following research questions:

- **RQ1:** What requirements engineering topics of interest and challenges are encountered by the developers of perception systems for driving automation systems?
- **RQ2:** What quality topics and challenges are encountered by the developers of perception systems for driving automation systems?
- **RQ3:** What software and systems engineering topics and challenges are encountered by the developers of perception systems for driving automation systems?

We refer to a topic of interest as something that practitioners currently practice or are curious about or would like to learn more about. A challenge on the other hand refers to an obstacle or difficulty that practitioners encounter and must overcome in order to successfully develop perception systems for driving automation systems.

To address these questions, we conducted seven group interviews with 19 expert participants from five companies that are currently working with ML-based perception systems for driving automation systems. Figure 2 gives an overview of the interview study.

Table 1 Overview of the conducted interviews, with the focus of the work conducted by the participants and the roles of the participants (same interviews reported by Heyn et al. [14])

Interview	Field of work	Participants
A	Object detection	Product owner
B	Autonomous Driving	Product owner, test engineer, ML engineer, software developer
C	Vision systems	System architect, product owner, requirement engineer, deep learning engineer
D	AD and ADAS	System engineer, manager AD
E	Testing and validation AD	System architect, two product owners, compliance officer, data scientist
F	Data annotations	AI engineer, data scientist
G	Autonomous Driving	System safety engineer

3.1 Data collection

We used semi-structured group interviews with a set of pre-determined open-ended questions. The use of semi-structured interviews ensured that all participants addressed the same questions, while still allowing the freedom to follow-up with additional questions on particular topics.³

The interviews were conducted between December 2021 and April 2022 via Microsoft Teams, and each lasted between 1.5 and 2 h. We recorded all interview sessions with the permission of all participants, then transcribed, and anonymized the recordings for analysis. At least three researchers were present in each interview, with two particular researchers in all interviews to maintain consistency.

A summary of the interviews and the participants who took part is shown in Table 1. We chose participants who possess experience with ML, perception systems for driving automation systems, software and systems engineering, RE, or data science, or who were working in the driving automation systems industry. The sampling method was a mix of purposive, convenience, and snowball sampling. We sent open calls to the Swedish automotive industry and our known contacts; then, we asked the interviewees for further contacts. Our participants work with different aspects of driving automation systems.

We started by asking for demographic information about the participants. We then showed them Fig. 1, asking for their feedback and using the figure to ground further discussions about how functional requirements relate to requirements on data and data annotation. We asked further questions about their requirements documentation, safety issues, and quality. Although we carefully chose interview participants, the opinions of the individual interviewees do not necessarily reflect the overall opinion of their companies. Due to the sensitive nature of information provided by interview participants and their respective companies, we are unable

to disclose the raw interview data or specific details about ways of working. Finally, in a 2.5-h workshop with roughly 20 participants, many of whom were interviewees, we presented and discussed our findings with illustrative quotes.

3.2 Data analysis

We applied thematic analysis, following the guidelines by Saladana [55]. We used a mixed form of coding, where we started with a number of high-level deductive codes based on the interview questions, and then we started inductive coding, adding new codes while going through the transcripts. At least three of the researchers worked together to code each of the transcribed interviews.

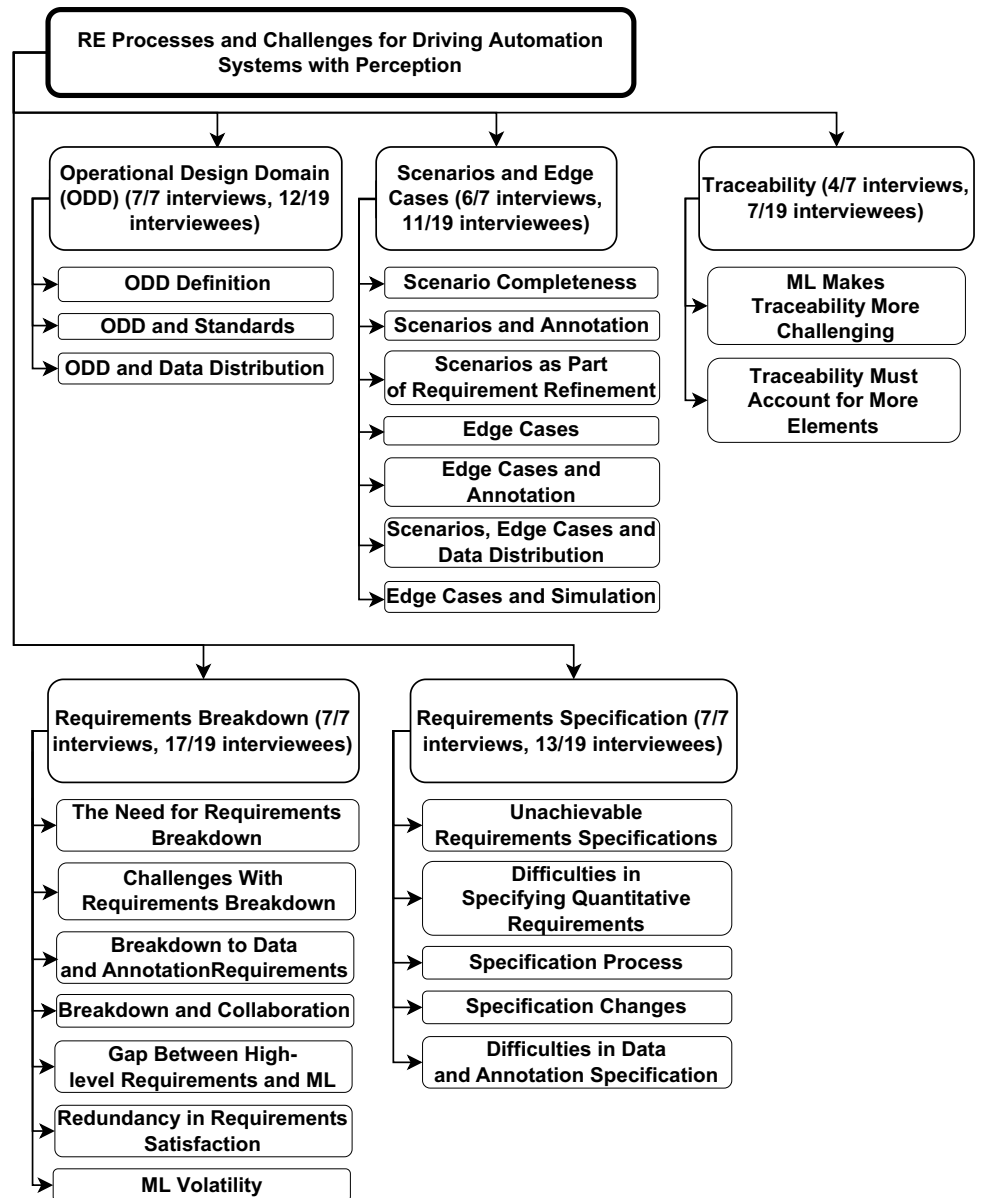
We observed saturation after five interviews, as not many new inductive codes emerged. In a second round of coding, a new group of at least two researchers per interview reviewed the interview transcripts and verified the codes. Finally, we used pattern coding to identify emerging themes and sub-categories. The final codes of each interview and the assignment of the statements of the interviewees to the sub-categories were reviewed by an additional independent researcher.

To illustrate our points, we use a number of interview quotes. For increased anonymity, participants are assigned a random identifier, such that P1 does not necessarily match to interview A.

As noted in Sect. 1, the results of the theme Requirements Engineering have been previously published [13]. This study enriches the findings on the RE theme, by reporting on the quality and systems and software engineering themes. The ecosystem and business, data, and annotation themes have been reported as well by Heyn et al. [14]. Although the article focuses on different themes, the qualitative topics covered in that article and our study here have some overlap, particularly in topics related to data and annotation. However, here, the topics of data and annotation are approached from an RE perspective, while the other article takes an ecosystems and process view on topics and challenges related to perception systems in driving automation systems.

³ The interview guide can be found at: <https://doi.org/10.7910/DVN/HCMVL1>.

Fig. 3 Mind map illustrating identified RE topics and challenges (RQ1) for driving automation systems with perception



4 Results: requirements engineering (RQ1)

Based on the thematic analysis, we divide the RE theme into sub-themes—“Operational Design Domain (ODD),” “Scenarios and Edge Cases,” “Requirements Breakdown,” “Traceability,” and “Requirements Specification”—and important topics within each sub-theme. The sub-themes and topics are summarized in Fig. 3. We also note how many interviewees discussed each sub-theme. Our findings reflect both RE topics and challenges, addressing RQ1.

4.1 Operational design domain (ODD)

An ODD is a description of a domain that a driving automation system will operate in—e.g., the road or weather

conditions. As part of RE, one needs to define not only requirements, but assumptions about the domain, context, and scope of operation. Operational context and scope for perception systems are particularly important as the intensity of hazards depends upon the current ODD. ODD-related topics came up in all interviews and were discussed by 12 of the 19 participants.

ODD definition: ODDs should be captured as part of the requirements specification. Several interviewees mentioned ODD detection—where the system detects that a certain ODD is currently applicable for a driving automation system function—and ODD exit detection—when the ODD is no longer applicable. ODD detection requires information on what to detect and detection accuracy. For

example, on highways, a driving automation system needs to detect different dynamic objects than in urban areas.

ODD and standards: Interviewees state that ODDs are critical, and therefore, it is desirable to follow a standard or process for specifying and defining ODDs. This need has been recognized and new initiatives for the definition of ODD exist, e.g., the interviewees mention the PAS-1883 standard, and we are aware of other standards (e.g., ISO 21448/SOTIF) that include ODDs.

ODD and data distribution: One interviewee stated that data distribution requirements are highly influenced by ODDs. For example, camera data can be classified according to descriptions in the ODD, and this mapping can reveal missing data, driving further data collection. As it is not feasible to collect data in all possible contexts, it is necessary to have an efficient sampling process covering the most common ODDs.

“If the performance of the model is not good enough in some part of the ODD, for instance, during the night or snow weather and so on, then we can select more samples from those areas.” P16

Another interviewee pointed out that although ODDs drive data collection, collecting certain types of data required by the ODD can still be very difficult.

“... mining for specific use cases. For instance, it is not easy to collect data that contains animals in it. You need some way to mine and find those specific frames which will be sent for annotations and then be used during training.” P16

4.2 Scenarios and edge cases

Several interviewees described how scenarios are crucial as part of the requirements specification process. In this context, scenarios describe specific operational paths and conditions for a vehicle, and one ODD may include a number of scenarios. As such, although there are links to scenario-based requirements methods [56], there are also clear differences. Scenarios and edge cases came up in six of the seven interviews and were discussed by 11 of the 19 participants.

Scenario completeness: It is important that perception systems perform correctly and that the vehicle handles failures in as many scenarios as possible. As such, scenarios can help in requirements derivation.

“If we refer to the classic system engineering process, I think nowadays it’s quite hard ... we are trying to use the scenario to derive the requirements. If we ... see the features or the distribution of the scenarios based on the data from the real world. Then we can derive the high-level requirements based on that data, the scenario database.” P4

However, when using scenarios to capture requirements, one interviewee stressed the difficulty of defining and assessing coverage. For example, when defining the scenarios for pedestrian children, how should the bounding box be defined, and how can coverage of a wide enough variety of children be ensured?

Scenarios and annotation: Even if all important scenarios are reflected in training data, annotation errors may result in unsafe behavior—e.g., a perception system may recognize a human as a tree during a snowy or rainy day.

“We’ll pick out some scenarios that we feel (are) likely not correct, for instance, if it’s a rainy night, then maybe the annotator is not annotating (people) as accurately as in the day.” P8

Scenarios as part of requirement refinement: Our results show that testing through scenarios, as part of a test-driven development process, enables iterative requirements refinement. Engineers iteratively refine their expectations of correct behavior by examining scenarios and capturing observations from simulation or in the field.

“... we have to learn through testing, so probably it will start with some rough set of requirements, some obvious set of requirements. Then we will, through real-world testing, discover and learn exactly how we want to behave.” P2

Edge cases: Interviewees stated that, in addition to normal scenarios, it is crucial and challenging to deal with edge cases. The interviewees used subtly different terms, such as edge cases, rare cases, and cases that occurred very infrequently. We use the term “edge cases” for simplicity. These cases may be missed by studying data distributions, but are very critical to ensure safety.

“The cars ... will end up in situations that no one could predict, that we’ve never seen before, and somehow we need, even in this situation, one individual car needs to perform better than a human driver, and human drivers are real good at handling edge cases. The neural networks will not do that.” P13

Edge cases and annotation: Edge cases cause issues by creating confusion among annotators. Data from edge cases are often annotated inconsistently. The topic of annotation is explored in more detail by Heyn et al. [14].

“We label whether a vehicle is in our lane or not. But how should you? You can think of so many corner cases when you are out driving. When you are doing a lane change. Which lane are you in then, and how would you then place all the other vehicles or lane lines? Maybe there are double lane lines and which

is valid and which is not? This leads to a lot of confusion among annotators.” P17

Scenarios, edge cases, and data distribution: One interviewee pointed out that scenarios, and especially rarer edge cases, are important for driving data collection efforts as part of having an effective data distribution. How well edge cases are covered can be an important development metric.

Edge cases and simulation: Interviewees stated that collecting data points for particular scenarios from the real world is necessary, but is particularly difficult for edge cases. This makes simulation challenging, as for safety-critical edge cases, practitioners have difficulty safely gathering enough data to run realistic simulations. This makes the process of iterative requirements refinement, as described previously, difficult for requirements associated with edge cases.

4.3 Requirements breakdown

Requirements breakdown can involve both refining and decomposing requirements. Requirements breakdown was brought up as a topic in all interviews and was discussed by 17 of the 19 participants.

The need for requirements breakdown: We see evidence that a traditional requirements breakdown is followed for perception systems. At least one participant spoke of splitting the problem to reduce complexity, as in standard development.

Another participant described an architectural-oriented breakdown. Based on a safety goal as a top requirement, a high-level architecture should realize that requirement, and then, an iterative refinement process is started to break down the high-level architecture into a logical architecture with smaller and smaller pieces.

Others describe the importance of separation of high-level requirements from technical requirements to have an upper layer that is resilient to change.

“To me, at least the function level will be the same in 100 years because there’s no need that you change it. If your function doesn’t change, because today you satisfy that function by combustion engine, in the next 50 years by electric, and in the next, I don’t know, 100 years by something more intelligent ... By changing your technical system level specifications, you still can satisfy your function.” P19

Challenges with requirements breakdown: Participants commented on the challenges of connecting high-level requirements to low-level requirements and general challenges with requirements breakdown in this context.

“...we do believe that it’s necessary to connect the top-level requirements or the quality of the function, and to

map that to quantitative or performance requirements on, for example, perception, precision, and control.” P13

Another interviewee describe the breakdown of requirements over components to requirements over other components, but commented that—although this breakdown is hierarchical—there is a need to work from both directions, not just top-down, but also bottom-up.

Several other interviewees report that traditional requirements breakdowns cannot be easily applied.

“For sure, we will not start with the classical software approach, where you start with some requirements and then keep breaking those down and through the V-Model because it will be impossible to capture the behavior of autonomous vehicle with requirements.” P2

Breakdown to data and annotation requirements:

Interviewees explained that although linking functional requirements to system accuracy is often possible, breaking functional requirements into data and annotation requirements is more difficult.

“Working with system level requirements, I can look at function requirements and figure out roughly what kind of accuracy we need ... That does not necessarily mean that I can tell how precisely annotation has to be, because I need to know how the software works to figure that out. Another translation needs to happen where I gave my requirements to the developers and they have to figure out what kind of accuracy they need from the data to meet the system requirements and with so many translations on the way, it is easy for things to get lost somewhere.” P6

Capturing requirements over data and annotation, linking to feature requirements, is particularly a challenge as the input space of the problem is large; thus, it is hard to capture a finite and complete set of these requirements.

Breakdown and collaboration: Challenges arise when teams collaborate to specify quality requirements. Often this involves collaboration with other teams and people. Frequent and direct interaction with the stakeholders can reduce this difficulty and help engineers to identify the requirements. In this case, stakeholders have internal roles in the perception system development.

“I think it is a lot of interaction with direct stakeholders in the end ... because the direct consumers of whatever you are producing know exactly what they need to fulfill their own requirements from their own stakeholders. So the negotiation across these interfaces is where the most interaction happens.” P9

Gap between high-level requirements and ML: When breaking down high-level requirements to very specific requirements on the ML-based perception system, results show that traditional RE practices are able to be applied up to a certain point—even though challenging. However, the breakdown for the ML based components is particularly challenging. As such, there are boundaries within the system where requirements methods change.

“If we talked about some other requirements or specifications not for the AD stack. ... those things still can follow the traditional way for critical system. ... if we distinguish those two parts, ... for the black box or part or AD business part, it’s hard to follow, but for the rest we still can leverage the classic knowledge.” P4

We see that it is difficult to specify requirements for the whole perception system. However, there are often still requirements—in terms of various performance metrics—at a high-level, for example, high-level metrics like safety, performance, functionality, or traffic comfort metrics (P4).

Redundancy in requirements satisfaction: One interviewee described how requirements are allocated to ensure redundancy in the solution. This topic of redundancy came up in several interviews.

“We typically try to break down the problem to come up with redundant solutions. You would have one algorithm using one sensor, which has some capacity to detect the pedestrian, and then use another algorithm and another algorithm in parallel. And you use another sensor and ... decompose the problem such that ... it’s very unlikely that all of them would miss this pedestrian. That’s a way to try and get reasonable requirements on every perception component.” P6

ML ‘volatility: One interview pointed out, due to dependencies between components and the volatile nature of ML, changes in the ML model can cause drastic changes in other parts of the system.

“Maybe one problem we have with ML is that, if there are things slightly off, it cannot just lead to a slight degradation, but to complete degradation of the entire system.” P17

4.4 Requirements traceability

Seven interviewees, across four interviews, brought up points related to traceability in perception systems.

ML makes traceability more challenging: Known requirements traceability challenges are exacerbated by the use of ML and associated data. Interviewees described that when systems or modules fail to meet particular key performance indicators (KPIs), tracing the source of the issue

is difficult due to the combination of ML models and traditional code. Traceability was discussed in four out of our seven interviews and by seven out of 19 participants.

“I think what is important at the end is the KPIs on the rightmost features of the figure (Figure 1). Then if you want to track down why it is not working, it’s not very easy to find which module is not working as supposed to, or maybe it works, but in a combination of something else, it creates some kind of strange behavior.” P14

Traceability must account for more elements: Several interviewees mentioned that it is important that traceability be maintained not just between code and requirements, but also with ML elements—e.g., models and datasets—that determine the overall functionality. An interviewee (P8) pointed out that they need to keep track of exactly which datasets were used to train the model, potentially to show the general public.

Typically, trace links would link to typical elements like requirements and safety goals, in order to understand motivations when something goes wrong, but as part of DAS development, elements like safety goals should also link to scenarios.

4.5 Requirements specification

Aspects of documentation and requirements specification were discussed in all interviews and by 13 of 19 participants.

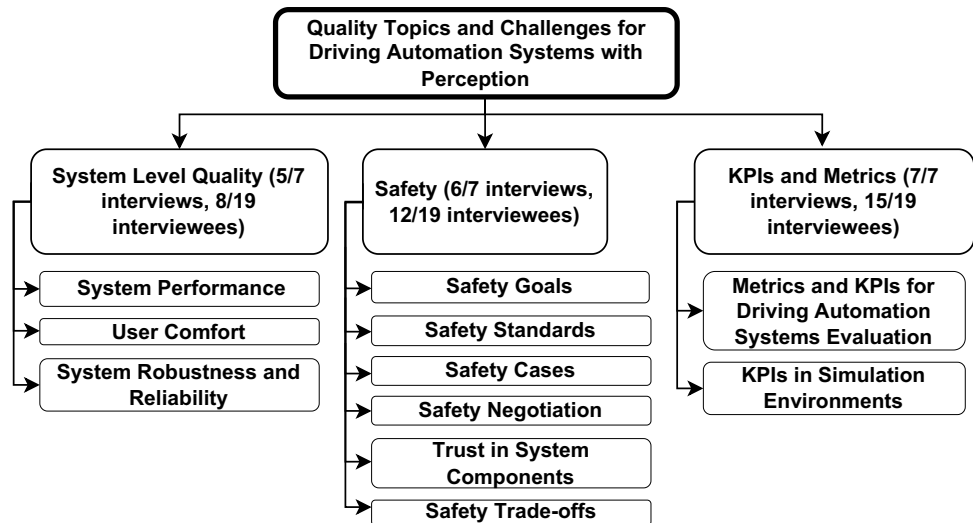
Unachievable requirements specifications: Two interviewees mentioned that sometimes clients provide unachievable requirements, even though requirements specifications are clear and precise.

“Sometimes clients come to us with a very well written set of requirements, like we want this annotator and want this precision or accuracy ... Then they send us data. But when we start looking at the data, it turns out that, given this data, these requirements are basically impossible to meet.” P18

Difficulties in specifying quantitative requirements: Due to confidentiality, interviewees were not able to elaborate on specific target levels for quantitative requirements. However, they did reflect generally about the difficulty in determining quantitative quality targets. For example, an interviewee questioned whether recognizing 99% of bounding boxes was enough and wondered how to make trade-offs between precision and recall.

Specification process: One interviewee emphasized that documentation of the rationale and goals of the project can serve as a form of requirement specification. They emphasized documenting underlying principles and the problem to be solved as part of requirements.

Fig. 4 Mind map illustrating relevant quality topics and challenges (RQ2) for driving automation systems with perception



Specification changes: The uncertain and highly iterative nature of perception systems and their development environment means that specifications are particularly prone to change.

“Requirements at any level are not something that is static. They should reflect your current best interpretation. These things can change because your understanding or your development process changes or the environment changes because there are suddenly new demands on how something is supposed to perform or you learn something new about the system or its environment.” P15

Difficulties in data and annotation specification: One interviewee said that specifying data requirements is difficult and different from functional specification, as it is hard to identify features and ensure data quality upfront.

“It’s very different how you write a data specification ... it’s hard to know what the future expects and what type of classes we want and how we want to combine certain objects ... we future proof our datasets quite well by specifying. We do specify a lot of classes.” P5

Another interviewee reported that it is difficult to specify quality (non-functional) requirements on data and annotation and to understand how qualities affect model performance.

“I work a lot with image quality before any ML is involved. Even that is very difficult to quantify. We can have very much right objectively measurable requirements on image quality, sharpness. Then how those translate to the actual performance of an ML algorithm is not at all linear.” P16

Another participant described challenges in specifying requirements for data annotation when dealing with external

partners. It is difficult to have an upfront, detailed specification of data classes and accuracy levels. Instead, data specification needs to be developed iteratively and experimentally with suppliers.

5 Results: quality (RQ2)

Based on the thematic analysis, we divide the Quality theme into the following sub-themes—“System-level Quality,” “Safety,” and “KPI and Metrics”—and important topics within each sub-theme. The sub-themes and topics are summarized in Fig. 4. We also note how many interviewees discussed the sub-theme. These topics and challenges are used to address RQ2.

5.1 System-level quality

This first sub-theme focuses on quality at the system level. This sub-theme came up in five interviews and was discussed by eight of the 19 participants.

System performance: As in other safety-critical domains, practitioners are required to satisfy performance and accuracy requirements for the entire system. They are often pressured to come up with precise accuracy numbers for ML models.

These requirements are typically quantitative in nature and often are prescribed as a bounded range of values rather than a single specific value to account for non-determinism. As an example, one interviewee discussed a parking situation where the exact performance depends on associated sensors and other factors, but the resulting behavior is deemed correct as long as it falls within the specified range, e.g., “... no farther than Xcm from the vehicle behind and no farther than Ycm from the wall (P11)”

However, the interviewees emphasized that performance measurement baselines are not always easy to define and that the current standards do not provide specific information on how to establish statistical expectations. Instead, the standards expect deterministic, specific behavior. Even without the use of ML, many factors result in non-determinism in embedded systems. The inclusion of ML leads to even further potential for non-determinism.

“That way of thinking doesn’t work with that specific standard (ISO26262), because that standard doesn’t have these kind of numbers. You can’t even write and say that things should be correct or the product should be correct with this kind of statistical numbers, but rather should be correct always. And I think that’s a bad way of thinking, let’s say I design a classical piece of electronics, even that one doesn’t work binary: yes or no always.” P1

Our interviewees stated that redundancy (e.g., redundant algorithms) can help improve performance. However, effective redundancy should be carefully planned—ideally not just with different approaches, but with different methods (e.g., parallel algorithms) and ideally different datasets. If there are multiple ways to gather data (e.g., radar, cameras, infrared), the possibility exists for sensor fusion. Various sensors have different abilities, e.g., radars are good at detecting orientation and speed while cameras are good at detecting objects. Sensor fusion enhances redundancy and, subsequently, performance in various scenarios. Another interviewee pointed out, however, that redundancy comes with a high cost and may lower usability:

“One pretty high level trade off would be cost or usability to the user. You could pack dozens of compute units and redundant sensors, which would drain of course the money of the user, or the customer, and the battery as well. You might just be able to drive around for a few minutes. But, you would shift the trade off towards more safety or availability just by piling up more redundancy. Of course, at some point, that it just is not feasible to use in an actual product that you could provide to customers.” P9

User comfort: The topic of user comfort in a vehicle also came up in our interviews. For example, in addition to allowing for redundancy, multiple sensors may improve a user’s comfort, e.g., by detecting speed and allowing for quicker adaptation.

System robustness and reliability: Several interviewees brought up topics surrounding robustness and reliability. Evaluating robustness and reliability requires consideration of a complex system made up of many sub-components—many of which come from suppliers. Driving automation systems consists of several components, and their robustness

and reliability are taken into consideration for the allocated ASIL (Automotive Safety Integrity Level) levels on the system.

“We have a total goal of robustness, reliability that includes the perception point of view, the components from a hardware point of view. We are discussing with the system suppliers which levels we are on right now from an ASIL point of view, and also from a reliability point of view, and confidence point of view.” P10

One interviewee pointed out the importance of data quality for robustness. If data quality is low, estimates of robustness may be inaccurate. In such situations, safety measures must be put in place to account for this uncertainty (P15).

5.2 Safety

As might be expected, given the criticality of driving automation systems, safety was one of the most popular discussion points in interviews. This sub-theme came up in six of the seven interviews, and was discussed by 12 of the 19 participants. Interviewees brought up “Safety Goals,” “Safety Standards,” “Safety Cases,” “Safety Negotiation,” “Trust in System Components,” and “Safety Trade-offs” as important topics in this sub-theme.

Safety goals: Multiple interviewees mentioned the importance of establishing safety goals. Such goals are the starting point of safety-related requirements specification. They are established at a high level and then connected in a hierarchy to lower-level individual system functions, where measures are put in place to ensure that the safety goal is realized throughout the system. Note that existing work has pointed out that safety goals for AD are entirely different from those defined for an ADAS system [57].

“Safety goals will basically be the starting point of the safety requirement specification... Then this will be the first parent requirement in the safety hierarchy and, then, in the next level—by some analysis, like fault analysis—you will try to understand what in the function level can violate that safety goal and introduce safety mechanisms.” P19

Safety standards: Safety is one of the most important quality attributes of driving automation systems. To ensure safety one needs to make sure that the nominal function is safe as well as in presence of faults according to the faults assumptions. However, the guidance on applying AI is limited and best practices are not fully established for safety argumentation.

Safety aspects of driving automation systems are normally structured according to ISO26262. However, it is commonly discussed if ISO26262 can actually address driving automation systems in an efficient way when they include

machine learning, e.g., [58]. At least one interviewee commented that the ISO26262 is not sufficient for systems using ML.

However, beyond ISO 26262 and SOTIF, there are emerging ISO standards (e.g., TS5083⁴), and guides (e.g., AMLAS [59]) focusing on AD, as well as those focusing on AI/ML (e.g., ISO PAS8800,⁵ TR5469⁶). One interviewee pointed out that the complexity of the topic leads to the emergence multiple standards....

Several interviewees mentioned SOTIF (Safety of the Intended Functionality⁷) as a further standard for functional safety of driving automation systems—used in conjunction with ISO26262. However, neither standard is fully adequate for ensuring both system and function safety, especially given that ML blurs the boundaries between functional and system safety.

One interviewee also mentioned multiple emerging safety standards for AI, including the TR5469 standard for functional safety in AI,⁸ PAS8800 for road vehicle safety for AI,⁹ and TR24029 for assessment of the robustness of neural networks (AI)¹⁰. Another interviewee referenced UL4600, which offers guidance on, among other topics, data safety [60].

Two interviewees also stated that they follow PAS-1883, an emerging standard for defining ODDs.¹¹

“You need to define ODDs and specify them... so far, there has not been a standard for how you define and how you specify ODDs. But we are very well aware about ... BSI standard PAS-1883. PAS is becoming a standard or instruction to follow. It’s an initiative to spell out how you should define an ODD, and how you cascade down, and how to specify an ODD.” P10

For these new standards, some challenges can be expected. One interviewee noted the challenge of applying the standards.

“I think it is also a little bit difficult for many of the engineers who are not experts in functional safety according to the standard to really understand what it means to them in their daily life...” P6

Furthermore, one interviewee pointed out that although standards can help developers avoid costly software failures, conforming to safety standards can also be quite costly. Another interviewee added that the top-down approach of safety standards does not match the current working procedures (e.g., agile way of working).

Safety cases: Safety cases are structured arguments used to show documented evidence that the system is sufficiently safe. Safety cases are an important aspect of driving automation systems development. However, it is important to ensure that the safety case matches the reality experiences by the developers. Furthermore, another interviewee brought up that data requirements are part of the evidence used in a safety case. Thus, it is important to have data management and data quality, for both training and validation data, as part of the safety case argumentation.

When asked about the impact of ML on safety cases, interviewees noted that safety cases can be defined in a modular manner over components of the system. ML-based components, naturally, must be part of such safety argumentation. However, one participant pointed out that safety cases are not yet well defined for driving automation systems and that establishing a methodology for the creation of safety cases in this context is difficult, as the new standards are still upcoming.

Another interviewee describes safety case argumentation as a joint process conducted with OEMs and suppliers:

“... We will definitely be interested in how our supplier has solved that problem and what safety argumentation they give us because we need to integrate it in in our safety case for the vehicle. We always have joint reviews, and we go into the detail on that.” P11

Safety negotiation: Driving automation systems are often built, at least in part, from existing components offered by external suppliers. Therefore, just as safety case argumentation is built jointly with suppliers, safety requirements must also be developed in conversation with suppliers.

Interviewees described a process where suppliers are assumed to have already assessed the basic safety of their already extant components outside of the context of a particular driving automation system before being contacted. This can be considered a safety element out of context, as discussed in [61]. Then, the driving automation system developers present safety requirements for a particular driving automation system as part of contract negotiation.

“When we meet them [suppliers] for the first time and we start talking we expect them to have done their homework over safety elements out of context—that they can present their assumptions and the safety holes that they have for their systems based on those assumptions. That means that they

⁴ <https://www.iso.org/standard/81920.html>

⁵ <https://www.iso.org/standard/83303.html>

⁶ <https://www.iso.org/standard/81283.html>

⁷ <https://www.iso.org/standard/77490.html>

⁸ <https://www.iso.org/standard/81283.html>

⁹ <https://www.iso.org/standard/83303.html>

¹⁰ <https://www.iso.org/standard/77609.html>

¹¹ <https://www.en-standard.eu/pas-1883-2020-operational-design-domain-odd-taxonomy-for-an-automated-driving-system-ads-specification/>

know what they're talking about, and they're good as a supplier. ... and then what we do is take care of the responsibility for the complete vehicle. Once we know our functions, we do a HARA (Hazard Analysis and Risk Assessment, from ISO26262), and from those we derive functional safety requirements that we give to them. You will need to fulfill this [requirement] to get this project." P11

Trust in system components: Because a driving automation system is constructed using components developed externally, interviewees noted that trust is initially a potential issue. Developing driving automation systems requires trusting that externally developed components are safe and reliable. GPS, in particular, was brought up by multiple interviewees due to the potential safety hazards of inaccurate GPS readings and map data. Obviously, they still will have a safety case at the end.

"I have always faced a challenge, which is how to trust the GPS, and even the map. Safety-wise, what we have always been told is, that a GPS should we treat it as a [ASIL level] QM (quality management, i.e., all assessed risks are tolerable), and the same with maps. But whenever you pinpoint to some specific target, or a supplier, or a sensor provider, the story is very different. They claim that they started talking about the accuracy and things like that, but functional safety-wise we will be fine. ... but you never can overcome that act of faith. I would say, show me some study that you have done to say that what you provide is up to an integrity level; that has never happened in my opinion. For me, it is still a challenge that we need to make a leap of faith when we select a specific supplier." P11

Safety trade-offs: Attainment of certain quality attributes, such as safety or performance, tends to require trade-offs with other quality attributes. Safety was a high priority for interviewees, even if it came at the cost of other quality attributes or driving automation system functionality. For example, one participant discussed limiting function availability to improve safety, choosing to only consider automation under certain situations (i.e., limiting the ODD) rather than widely allowing the car to be driven autonomously. Another interviewee discussed limiting vehicle functionality to improve safety, in this case, compromising the speed of the vehicle for improved safety.

"Speed is an obvious trade-off, so we also trade high speeds for safety. We go at lower speeds. Speed is such a dimensioning parameter in all the safety work since all the risks are high when you increase speed." P7

5.3 KPIs and metrics

Assessments of driving automation systems quality are made by tracking certain performance metrics, often called "key performance indicators" (KPIs), and comparing the attained value to selected thresholds. This sub-theme was discussed in all seven interviews, by 15 of the 19 participants. Interviewees brought up "KPIs in Simulation Environments" and "KPIs and Metrics for Driving Automation Systems Evaluation" as topics in this sub-theme.

Metrics and KPIs for driving automation systems evaluation: For evaluation of the driving automation system, interviewees explained that—rather than specifying deterministic properties—they track a set of high-level metrics related to safety, performance, functionality, comfort, and other factors. These metrics can be tracked over the execution of many different scenarios, either in simulation or in a real vehicle; then, statistical analysis of the collected observations can be used to make an assessment of the driving automation system.

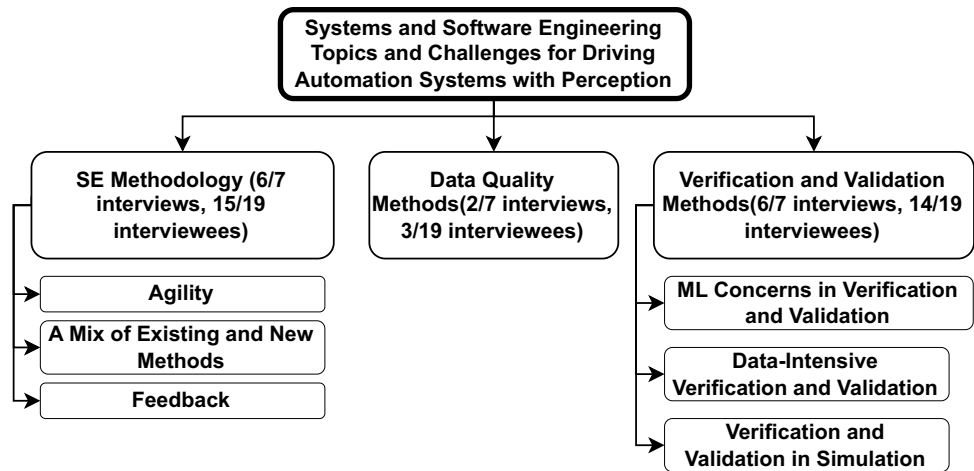
"If we say the requirements were a specification for the entire driving automation system stack, I think definitely it's quite hard to have very precise or detailed specification for all the functions, but actually we have some high level metrics like safety metrics or performance metrics, functionality or traffic comfort metrics, those metrics are on a very high level, which means we can use those metrics combining with the scenario database and then we run millions of the scenarios and get the statistical analysis report. We can't say we don't have anything for testing or validation. We have something but they are very different from the traditional understanding of the specification." P4

For evaluation of the low-level ML components—such as specific classification models—interviewees stated that they employ standard metrics used in other ML domains. They note, however, that they pay close attention to the specific data in the dataset (e.g., values for specific classes like pedestrians) to ensure that careful evaluation is performed.

When interacting with ML, assessments of safety—and the design of the functionality under assessment—must take into account uncertainty in the ML output. Uncertainty estimation is employed at all levels, from algorithm design, to data selection, to evaluation. One interviewee described their use of uncertainty estimations to inform data selection and to reveal gaps in their training data.

The KPIs of systems dependent on ML will ultimately be determined by the data used for training and validation of the ML components. Therefore, to ensure that KPIs are informative and realistic, high-quality training and validation data must be used, but that gathering such high-quality data is expensive. It is important that validation metrics

Fig. 5 Mind map illustrating relevant systems and software engineering topics and challenges (RQ3) for driving automation systems with perception



have good results, as performance on the road is unlikely to surpass these lab results. However, another interviewee emphasized that KPIs or quality measures are difficult to identify without a lot of background knowledge, e.g., the average size of cars.

Often, data are reused, especially for rare scenarios. However, reusing data from one scenario to another may decrease realism and can raise issues in the accuracy of KPI values during validation.

One interviewee also stressed the importance of communicating KPIs and metrics for data quality and variance to users who use annotated data.

“You should produce such KPIs (on data variance) and communicate this to the users. In this case, the users say alright, or if we need to adjust anything from ... the data collection perspective. And the same goes with the annotation as well, if they are of good quality or not, essentially it’s KPI numbers to the users as input.” P8

KPIs in simulation environments: Because simulations may not accurately reflect the real world, KPI observations gathered from simulations may also not match the observations that would be made in reality. One interviewee was skeptical of using values from simulations as evidence of safety.

“During the verification and validation, how you can verify that the KPIs that you get during the validation actually show the reality? Especially if you do some sort of synthetic simulation, ... then it’s not very clear to me how one can argue the safety aspect, that we can reach the same KPI in real world.” P16

Another interviewee noted that KPIs should be compared between simulations and real-world testing to help show that the simulation is realistic; otherwise, the value of such simulations is in doubt.

6 Results: systems and software engineering (RQ3)

As with the previous results, we divide the Systems and Software Engineering theme into sub-themes—“SE Methodology,” “Verification and Validation Methods,” and “Data Quality Methods”—and important topics within each sub-themes. The sub-themes and underlying topics are summarized in Fig. 5. We also note how many interviewees discussed the sub-theme. These sub-themes and topics address RQ3.

6.1 SE methodology

Software engineering methodologies refer to the frameworks or approaches that guide the processes, activities, and tasks involved in software development [62]. Such methodologies provide a structured and systematic way to plan, design, develop, test, evaluate, and deliver software products to customers. Different methodologies (e.g., waterfall, agile) offer a distinct set of principles, practices, and techniques to manage the software development lifecycle. The interviewees described several aspects of their software development methodologies and how they have changed in the face of ML use. This sub-theme came up in six of the seven interviews and was discussed by 15 of the 19 participants. Interviewees brought up agility, feedback, and a mix of existing and new methods as topics in this sub-theme.

Agility: As context, many of our interviewees have transitioned or are transitioning to a more agile way of systems development for driving automation system development. For example, participants are moving from a process based on the traditional V-model to a specific agile framework, e.g., SAFe.

“We are using SAFe as a formal agile framework in our software development process for autonomous vehicles.” P2

Such transitions and the use of agile methods at a large scale bring their own challenges unrelated to ML, e.g., [63]. These issues form a background to our exploration of ML use in driving automation systems.

Previous work focusing on methodological transitions has pointed out culture clashes between more agile and traditional ways of working [64, 65]. One interviewee points out similar clashes between control engineers—who develop and test software and hardware—and software developers, who are focused on data and simulations.

“I think the biggest clash is the way of working. Control engineers want to approach things in their way. They want to synthesize things and test things for real in the vehicle. Whereas the software companies have a more data-driven background and they start off with the data and they just work in different ways.” P2

A mix of existing and new methods: Many traditional methods and processes where, e.g., requirements are defined and broken down to different parts of the system will still apply. This includes large-scale agile methods, such as SAFe. However, several interviewees reported that it is difficult to implement existing software engineering methods when the system includes machine learning because of non-deterministic behavior. In addition, since the use of machine learning is rather opportunistic or technology-driven, the existing processes can be infeasible because of non-deterministic behavior of such systems. The participants focus on integrating and adapting methods while using machine learning.

“If we distinguish those two parts, we just say for the black box or part or autonomous driving business part, it’s hard to follow [existing methods], but for the rest, we still can leverage the classic (development) knowledge.” P3

Some interviewees described this mix of methods by describing both a top-down and a bottom-up start. Traditionally, development would start with up-front requirements elicitation, with a break-down facilitating system development (similar to Sec. 4). However, in the case of driving automation systems, the development process is less top-down but rather composed of many different components that must fulfill the product definition on top level.

“It’s like bottom-to-top and then back to the bottom. We have just like a rough starting point, like we want to drive this route, and that’s basically the scenario or the high-level requirements, and then we build some algorithms and try it. We start in that sense in the bot-

tom and then, when we have something in place, we can start testing in a structured way with replay of logs and with the structured scenario database and all that and then we learn.” P2

Several participants describe using a highly iterative development method for driving automation systems. Although most companies have already implemented some degree of agile development—which already prescribes working in an iterative manner—the level of agility in practice may vary. However, iterations appear to be key for machine learning development. Thus, this need for frequent iteration must be satisfied either through agile or other iterative ways of working.

6.2 Data quality methods

Data quality in machine learning systems is critical and has an effect on development methods and ways of working. This sub-theme came up in two of the seven interviews and was discussed by three of the 19 participants. Interviewees emphasized the importance of placing a high priority on data quality, particularly with safety in mind.

“First, you focus on data itself. You can try to identify the purpose of using this data in your system and identify other possible issues. There are some guidelines like data safety guideline from the safety-Critical System Club, and then they have a very structured guideline about how to identify them. I think they are hundreds of identified possible issues about the data in safety critical system. So, first you can probably try to explore the data itself.” P4

To ensure the quality of the data, tools, and requirements, one interviewee indicated that traditional quality assurance methods, for example, fault analysis for preventing and solving discovered issues, are still valid, even with new technologies. Another participant mentioned the importance of feedback as part of the data annotation process, which involves creating documents to clarify annotation uncertainties which can be continually used by annotators.

6.3 Verification and validation methods

Verification and validation methods refer to the systematic techniques and activities employed to confirm that the system meets the specified requirements and intended purposes, and the system has been designed, developed, and implemented correctly. This sub-theme came up in six of the seven interviews and was discussed by 14 of the 19 participants.

Interviewees brought up machine learning concerns in verification and validation, data-intensive verification and validation, and verification and validation in simulation

as topics in this sub-theme. Note that although there is a difference between validation and verification, interviewees are using both terms (verification and validation) as a high-level definition of general verification and validation activities.

Machine learning concerns in verification and validation: Practitioners must select appropriate verification methods for the specific machine learning algorithms being deployed. Verification of machine learning is a relatively new field of research, and practitioners may need to monitor ongoing research regarding particular types of machine learning.

“There is some research about how to verify the neural network like abstract interpretation and then you can use the safety engineering method like to design some measures or set a scope for the neural network and try to have like redundant pipeline and to monitor this algorithm.” P3

A common challenge when performing verification is that the model acts like a black box—it is difficult to infer how a model makes its decisions. From the engineering perspective, to ensure that the whole system works well, it can be important to focus on the explainability of the model. Practitioners start with something that they feel might work, then check the results, and then adjust accordingly. Once the system achieves basic functionality, if any limitations (e.g., poor performance) of the artificial intelligence algorithms or the neural network are observed, then the engineers attempt to understand the model to identify the reason.

“If we see the limitations of the artificial intelligence algorithms or the neural network, then probably we want to dig into this black box and try to figure out what’s the reason behind that.” P4

Although explainability affects the way in which systems are developed, our interviewees reported that it is typically not their first concern. They often focus on getting the system working first and then consider the qualities such as explainability.

“I think we’re in the stage where we’re just trying to get the whole thing working first and then we would go into more understanding why it works the way it is.” P2

Data-intensive verification and validation: With the incorporation of machine learning and artificial intelligence, the verification and validation of the perception systems become more data intensive. For example, an interviewee pointed out that effective verification and validation require representative data selection, i.e., they need the right data for data-driven validation. Furthermore, the interviewee describes performing statistical analysis on the collected data as part of verification. The participant also emphasized

the need for an efficient sampling method to cover the most common ODDs, as extensive data collection is not practical.

Driving automation system practitioners often meet the need for verification data with synthetic data, emulating data collected from sensors or cameras. Due to the complexity of collecting real data for scenarios that appear very rarely (e.g., edge cases), it would be beneficial to use such synthetic data, but the effort of validating synthetic data could be higher than actually collecting this kind of data.

“I don’t believe that synthetic camera data or LIDAR data has sufficient fidelity to be used for validation yet. It’s lacking in many aspects. And to prove that it is actually useful for validation would probably require more data than not even using it to begin. [Note: validating the usefulness of synthetic data would require so much data that synthetic data is no longer needed]” P7

Traditional systems are often verified by comparing observations to specific expected output. Given the quantity of data needed and the difficulty of specifying a deterministic outcome, driving automation systems with machine learning may need to be verified based on statistical analyses. In other words, they need to statistically verify that the created and used datasets capture the whole ODD.

Verification and validation in simulation: Interviewees reported a shift toward verification and validation in simulated environments from using actual hardware because of the safety-critical nature of driving automation systems, which require different steps in the validation and verification process to rule out as many issues as possible before going on the roads. Practitioners collect data and conduct simulations to test the driving automation system in a virtual environment. One interviewee stated that they are focusing more on scenario-based approaches for verification and validation.

“From the simulation team, we are trying just to shift the direction from the classic embedded system world ... that’s the reason why we have the scenario-based approach for validation. So if we see the data from some other companies like Waymo or Uber, they spend 99.95% of the test cases in their virtual environment, just 0.05% on the real vehicle. Because some of the scenarios are really dangerous, we can’t just to ask a driver to drive on the road and do some to test some edge cases.” P4

On the other hand, it can be challenging to perform verification and validation of driving automation system in a simulation. An interviewee points out that there might be limitations when using synthetic simulations for safety argumentation; it may be difficult to achieve the same KPI results in the real world compared to simulated validation

and verification. Furthermore, when machine learning is incorporated into the system, it can be difficult or very expensive to create realistic data for the verification and validation process.

Interviewees consider the balance between the percentage of verification that should take place in the simulation and in the real environment. The challenge of realism indicates that some verification should take place on real hardware. However, it is more time-consuming to test in the real world than using simulations, and therefore, simulations are an efficient method to be used before doing the final tests on the roads. Regardless, interviewees stressed that it is important to ensure that the test cases are representative of the driving scenarios in both the real and simulated environments.

The degree of DAS automation for simulations¹² also has an impact on how validation and verification is conducted. For lower levels of automation (i.e., very little driving automation), one interviewee points out that simulation can be outsourced to suppliers. However, for higher-levels of driving automation, they will need to develop in-house simulation. V & V for DAS is also affected by standardization, including new and upcoming standards like TS5083.

7 Summary and discussion

In this section, we summarize our results, answer our RQs, list future directions for research and practice, and discuss threats to validity.

7.1 Requirements engineering topics and challenges (RQ1)

We have identified a number of RE topics and challenges in Sect. 4, as summarized in Fig. 3. These topics and challenges can be seen as a checklist when working with machine learning-based perception systems—a list of issues that should be considered.

Our interviewees emphasize that the definition and limits of ODDs are an integral part of perception systems, and these ODDs have important impacts on data requirements and collection, confirming findings in Heyn et al. [35]. Similarly, perception systems development relies heavily on the use of scenarios and associated edge cases. Such scenarios play a key role in dictating annotation, data collection, and simulation.

In terms of challenges, our results indicate that **ODD detection** and **ODD exit detection** are challenging, as these require information not only about what to detect in the environment, but also how to detect and the accuracy

of the detection confirming findings in [37]. In addition, **data requirements** are highly influenced by the content of an ODD; therefore, ODDs can be used to evaluate whether a data distribution is sufficient for good machine learning model performance. However, it is not always easy to **collect the data** specified by ODDs. Heyn et al. also emphasized the importance of ODDs in driving automation systems and noted the lack of a common definition for ODDs [24]. Our participants go further and mention the need for ODD standardization (and efforts in that regard).

One major challenge is that simulations should reflect **realistic scenarios**, echoed by Acuna et al. [66]. To ensure safe perception, the collected data and scenarios must be thorough, and the perception system must avoid failure in all scenarios. In addition to covering normal scenarios, it is important to **specify edge cases** among scenarios, which are then used to determine data distributions. However, edge cases introduce challenges as they create **confusion among annotators** and are challenging to **test in reality** due to safety concerns.

Breaking down requirements for data and annotations can be very difficult, and additional challenges are introduced due to requirements dependencies and the need for multiple teams to collaborate. In general, we believe that the **gap between standard RE methods and machine learning components** is both a technical gap and a gap in training and backgrounds, as the machine learning components are often engineered by data scientists without a software engineering background confirming the results in [9, 20].

Difficulties in breakdown, machine learning opaqueness, as well as the introduction of more elements to trace (e.g., ODDs, scenarios, training data), make it difficult to establish **traceability**. These challenges add to the known challenges with motivating and using traceability in practice [67].

Creating **specifications for data and annotations** is challenging, as it is difficult to have an upfront specification for data classes, e.g., pedestrians and crosswalks. Furthermore, sometimes machine learning components are assigned unrealistic and **unachievable requirements**. Although requirements change is a frequently acknowledged RE problem [68], with perception systems, the **level of uncertainty and change** is particularly high due to uncertainty about the system, including machine learning, and the environmental targets. **Quantifying quality requirements** (e.g., accuracy) is also particularly challenging in perception systems, echoing the results of Vogelsang and Borg [9].

7.2 Quality topics and challenges (RQ2)

As part of the Quality theme, our interviewees have identified a number of topics and challenges. Practitioners consider **performance, reliability, robustness, safety, and user comfort** as important quality attributes. It is

¹² <https://www.sae.org/blog/sae-j3016-update>

interesting to note that the space of qualities that interviewees focused on is generally small, compared to the space of NFR qualities explored in past academic work [18].

Interviewees found that for driving automation systems, **performance was difficult to measure accurately**. As a means of ensuring both performance and safety, redundancy of algorithms and sensors was important, but interviewees noted that **redundancy must be carefully designed**, particularly in terms of data and algorithms, and redundancy as a principle for system design has both limitations and trade-offs.

In the context of driving automation systems, safety is particularly critical, as also noted by [45]. Safety assurance is already challenging for conventional driving automation systems software and becomes even more challenging with the inclusion of machine learning. Practitioners set safety goals, often in negotiations with component suppliers. **Safety negotiation** with suppliers has been a challenge, including issues of **trust in components and suppliers**. In addition, ensuring driving automation systems safety requires collaboration and effort from different parties is challenging, confirming findings in [53]. Interviewees also acknowledged that safety does not operate in a vacuum, recognizing **safety trade-offs** with, for example, security, availability, and functionality. Safety cases are a critical element of assuring that the safety goals are met, but are also more complex given the uncertainty of machine learning.

To ensure safety, practitioners must comply with **evolving safety and AI standards**, including established standards such as ISO26262 [69]—which are not sufficient to account for the incorporation of machine learning. This confirms the findings of [46, 47] and underlines the need of newer machine learning-specific standards such as TS5083 [70]. In general, given the wide range of standards, there have been challenges in understanding, managing, and conforming with the relevant standards. Although there are significant costs associated with safety incidents, conforming to standards is also costly.

The large input space and non-determinism of machine learning complicate quality assurance. Instead of specifying concrete expected behaviors in key scenarios, quality assurance is performed by **tracking critical KPIs** during the execution of catalogs of scenarios and performing statistical analysis on captured observations. KPIs can be **defined on multiple levels**, including **driving automation systems-specific KPIs** and **standard KPIs for machine learning components** (e.g., **precision and recall**). KPI assessment is affected by **training and validation data, uncertainty**, and **simulation realism** confirming the results of [45].

7.3 Systems and software engineering topics and challenges (RQ3)

Our interviews also revealed topics and challenges related to systems and software engineering development methodologies for perception systems. Our findings show that the presence of **machine learning adds further complexity to agile ways of working**. Existing traditional and agile methodologies are not sufficient to meet the needs of large-scale machine learning, echoing the findings of [51]. Our interviewees apply a **mix of methods** using more traditional, top-down engineering in some areas and more iterative, bottom-up development in others. In this way, from a methodological point of view, a **continuous feedback cycle** is key to successful delivery.

The focus on safety further complicates development methods, confirming the findings in [54], as many **safety methodologies do not adequately address machine learning**. The importance of data to perception systems requires changes in development practices. In particular, practitioners need **data quality methods**.

Our findings show that **verification and validation is more challenging and data-intensive in the presence of machine learning**. Data selection and consideration of data quality are required to ensure effective verification. The use and acquisition of **synthetic data** are an important topic, but raises data quality issues. Rather than comparing observed behavior with specific, expected outcomes, **V & V is based on statistical analyses** of quantitative metrics. To gather sufficient observations and limit the risk to vehicle operators, verification and validation use simulation. However, it is a challenge to have **realistic simulation** and to **determine in which situations simulation can replace real verification**.

7.4 Future directions in research and practice

Some of the identified challenges in RQ1-3 are relatively new from an RE and SE perspective (e.g., ODD detection, missing edge case, the proliferation of machine learning-related safety standards, machine learning verification and validation), while others have been long recognized (e.g., traceability [67], specification changes [68], and quality trade-offs [71]). Our findings point to a number of new research topics. We outline these areas, highlighting examples of existing state-of-the-art work on these topics.

Although the focus of our work has been on perception systems, we believe that many of the topics and challenges found apply more generally to other domains reliant on machine learning. For example, challenges breaking down specifications would hold due to the volatility and opaqueness of machine learning. Future work should contrast RE challenges and practices in other machine learning-enabled domains.

ODD methods: Our findings illustrated the importance and challenges associated with ODD development as part of complex machine learning systems. Existing work has focused on various aspects of ODD development as part of autonomous driving, e.g., [72], including a consideration of safety [35, 73], but our industrial partners still find this topic a challenge. ODD can be linked to broader work on context in RE, e.g., [74, 75], but future work should explore what aspects of ODD context are domain-specific or general.

Data requirements: How to capture and define requirements over data is an issue that should be a focus of future work. Although previous work has looked at data-driven RE, e.g., [76], this focuses more on gathering standard requirements from sources such as social media, rather than requirements for the data needed for machine learning. Other work has looked at data quality, but from an era before the rise of machine learning, e.g., ¹³

Requirements traceability with machine learning: Tracing requirements to system elements is essential for safety argumentation and change management, but becomes challenging when traces must include machine learning components like models and data. Although traceability has been heavily investigated from a requirements perspective [77], traceability for machine learning is only just starting to be explored, e.g., [78]. Other works look specifically at traceability from the perspective of data provenance or data lineage as part of machine learning [79].

Scenarios as specifications: Given the challenges of defining up-front, complete requirements for driving automation systems, practitioners have turned to scenarios both for specifying data and verification, including simulation. Scenarios have been an active topic in RE, e.g., [80], but mainly for improving the quality of gathered requirements, including completeness, and not as a stand-in for traditional requirements. Using scenarios as specifications data-driven development requires further attention.

Quantifying machine learning requirements: Interviewees expressed difficulties with placing specific targets on quality requirements for driving automation systems, e.g., performance requirements, echoing some of the challenges found in recent work [18]. Although quantifying quality requirements, or metrics, has been an active area of investigation for many years, e.g., [81], most work on metrics for machine learning is specific to particular qualities, e.g., uncertainty [82]. Setting targets for such metrics is particularly difficult and context specific.

Redundancy: Redundancy as a means to improve performance and safety arose as a prominent issue. Redundancy has been studied in general for systems engineering, e.g., [83], and has been studied from the perspective of multiple

machine learning models, e.g., ensemble learning [84]. However, more consideration can be made at combining these perspectives, considering machine learning redundancy at the system level.

Safety standards and machine learning: Safety standards and machine learning is an active topic, e.g., [58], but we see in our findings a proliferation of many possible standards or frameworks which can be applicable for driving automation systems, but the selection or integration of these multiple standards has not been extensively explored.

Large-scale agile with machine learning: Much work has been dedicated to reporting and making recommendations concerning transitions to agile methods for large-scale systems, e.g., [63]. However, most work does not consider the challenges introduced with machine learning. In terms of machine learning development, methods like CRISP-ML [85] and a more recent focus on MLOps [86] attempt to guide development. But the combination of these machine learning and data-driven methods with established, large-scale agile methods like SAFe¹⁴ is still mainly unexplored.

Large-scale agile with machine learning and safety: Adding to the complexity of the previous direction, such large-scale, agile, machine learning-enabling methods should also be usable in safety-critical contexts. Safety challenges as part of large-scale agile have been investigated [87], but not in an explicit machine learning context.

Verification and validation for machine learning: Testing and related activities for machine learning are already an active area of investigation [88]; however, we feel we should highlight this direction, in particular the areas of synthetic data curation and simulation, as it was raised by several interviewees. Others have begun to investigate the utility of synthetic data, e.g., [89], but further investigations in a driving automation systems or safety-critical contexts are needed.

7.5 Threats to validity

Internal validity: We internally peer-reviewed the interview guide and conducted a pilot interview to improve the guide and process. We sent a preparation email to all the interview participants with the details and purpose of the interview study. To maintain consistency in the interview process, at least three authors conducted each interview, with two authors present in all interviews.

All interviews were conducted in English, and the auto-generated transcripts were ‘fixed’ by authors by listening to audio recordings and correcting any transcription errors. Note that the working language of each company was English, so the language should not have created barriers.

¹³ <https://iso25000.com/index.php/en/iso-25000-standards/iso-25012>

¹⁴ <https://scaledagileframework.com/>

Although qualitative coding always comes with some bias, we mitigated this threat by following established literature [55], coding in multiple rounds, using inductive and deductive codes, and having multiple authors participate in each round of coding, with in-depth discussion on code meanings and assignments. Many of our provided quotes in the results section reflect individual opinions; however, we were able to synthesize the results from different opinions into themes. Moreover, the interviews were group interviews; quite often, an individual opinion was reflective of the opinion of the group; otherwise, it would have produced contrary discussion. We note that participants in our group interviews had rather harmonized opinions.

External validity: We used a mixture of purposive and snowball sampling. As our study needed a certain set of expertise to answer our questions, we could not conduct random sampling, using our networks and their contacts. Still, due to the size of the study, with participants covering a wide variety of roles with varying experience levels, covering differing company roles and sizes in the perception system ecosystem, we believe we have a relatively representative sample.

Furthermore, we argue that we reached a sufficient point of saturation with our interview data, as we noticed a sharp decline in emerging codes after analyzing the fifth group interview.

Note that one cannot link participants to interviews and companies this is done deliberately to protect the anonymity of our participants. Although this may affect transferability of our results, we feel this level of anonymity does not greatly hurt our results. Though our study results are limited to perception systems in DAS, we argue that some findings can apply to other safety-critical or perceptions systems. This applicability should be explored in future studies.

8 Conclusion

Our study investigated requirements engineering, quality, and systems and software engineering topics and challenges during the development of DAS. We interviewed 19 participants from five companies and identified a number of topics and challenges that have a major impact on the specification, development, and quality of DAS. The results of this study offer guidance to practitioners and suggest future research directions in the intersection of requirements engineering, software quality, development methodologies, and machine learning to help mitigate the challenges practitioners are facing.

Acknowledgements Support for this project was provided by Vinnova pre-study 2021-02572. We also thank all interview participants.

Funding Open access funding provided by University of Gothenburg.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Mallozzi P, Pelliccione P, Knauss A, Berger C, Mohammadiha N (2019) Autonomous vehicles: state of the art, future trends, and challenges. In: Dajsuren Y, van den Brand M (eds) *Automotive systems and software engineering*. Springer, Cham, pp 347–367
2. Borg M, Englund C, Wnuk K, Duran B, Levandowski C, Gao S, Tan Y, Kaijser H, Lönn H, Törnqvist J (2018) Safely entering the deep: a review of verification and validation for machine learning and a challenge elicitation in the automotive industry. arXiv preprint [arXiv:1812.05389](https://arxiv.org/abs/1812.05389)
3. Cooling J (2019) *The complete edition—software engineering for real-time systems: a software engineering perspective toward designing real-time systems*. Packt Publishing Ltd, 35 Livery Street Birmingham B3 2PB
4. Highsmith J (2013) *Adaptive software development: a collaborative approach to managing complex systems*. Addison-Wesley, New York
5. Knight J (2012) *Fundamentals of dependable computing for software engineers*. CRC Press, Boca Raton, FL
6. Salas E, Goodwin GF, Burke CS (2008) *Team effectiveness in complex organizations: cross-disciplinary perspectives and approaches*. Routledge, London
7. Lurie Y, Mark S (2016) Professional ethics of software engineers: an ethical framework. *Sci Eng Ethics* 22:417–434
8. Belani H, Vukovic M, Car Ž (2019) Requirements engineering challenges in building AI-based complex systems. In: 2019 IEEE 27th international re conference workshops (REW), IEEE, pp 252–255
9. Vogelsang A, Borg M (2019) Requirements engineering for machine learning: Perspectives from data scientists. In: 2019 IEEE 27th international requirements engineering conference workshops (REW), IEEE, pp 245–251
10. Lwakatare LE, Raj A, Bosch J, Olsson HH, Crnkovic I (2019) A taxonomy of software engineering challenges for machine learning systems: An empirical investigation. In: *International conference on agile software development*, Springer, Cham, pp 227–243
11. Arpteg A, Brinne B, Crnkovic-Friis L, Bosch J (2018) Software engineering challenges of deep learning. In: 2018 44th Euromicro conference on software engineering and advanced applications (SEAA), IEEE, pp 50–59
12. Amershi S, Begel A, Bird C, DeLine R, Gall H, Kamar E, Nagappan N, Nushi B, Zimmermann T (2019) Software engineering for machine learning: A case study. In: 2019 IEEE/ACM 41st international conference on software engineering: software engineering in practice (ICSE-SEIP), IEEE, pp 291–300

13. Habibullah KM, Heyn H-M, Gay G, Horkoff J, Knauss E, Borg M, Knauss A, Sivencrona H, Li J (2023) Requirements engineering for automotive perception systems: An interview study. In: International working conference on requirements engineering: foundation for software quality, Springer, pp 189–205
14. Heyn H, Habibullah K, Knauss E, Horkoff J, Borg M, Knauss A, Jing Li P (2023) Automotive perception software development: Data, annotation, and ecosystem challenges. In: 2nd international conference on ai engineering–software engineering for AI, IEEE
15. Pradhan SK, Heyn H-M, Knauss E (2023) Identifying and managing data quality requirements: a design science study in the field of automated driving. *Softw Qual J*. <https://doi.org/10.1007/s11219-023-09622-8>
16. Villamizar H, Escovedo T, Kalinowski M (2021) Requirements engineering for machine learning: A systematic mapping study. In: 2021 47th Euromicro conference on SE and advanced applications (SEAA), IEEE, pp. 29–36
17. Ali MA, Yap NK, Ghani AAA, Zulzalil H, Admodisastro NI, Najafabadi AA (2022) A systematic mapping of quality models for AI systems, software and components. *Appl Sci* 12(17):8700
18. Habibullah KM, Gay G, Horkoff J (2022) Non-functional requirements for machine learning: An exploration of system scope and interest. In: 2022 IEEE/ACM 1st international workshop on software engineering for responsible artificial intelligence (SE4RAI), IEEE, pp 29–36
19. Pei Z, Liu L, Wang C, Wang J (2022) Requirements engineering for machine learning: A review and reflection. In: 2022 IEEE 30th international requirements engineering conference workshops (REW), IEEE, pp 166–175
20. Ahmad K, Abdelrazek M, Arora C, Bano M, Grundy J (2023) Requirements engineering for artificial intelligence systems: a systematic mapping study. *Inform Softw Technol* 158:107176
21. Maalej W, Pham YD, Chazette L (2023) Tailoring requirements engineering for responsible Ai. *Computer* 56(4):18–27
22. Gjorgjevikj A, Mishev K, Antovski L, Trajanov D (2023) Requirements engineering in machine learning projects. *IEEE Access* 11:72186–72208
23. Ahmad K, Bano M, Abdelrazek M, Arora C, Grundy J (2021) What's up with requirements engineering for artificial intelligence systems? In: 2021 IEEE 29th international re conference (RE), IEEE, pp 1–12
24. Heyn H-M, Knauss E, Muhammad AP, Eriksson O, Linder J, Subbiah P, Pradhan SK, Tungal S (2021) Requirement engineering challenges for AI-intense systems development. In: 2021 IEEE/ACM 1st workshop on AI engineering-SE for AI (WAIN), IEEE, pp 89–96
25. Heyn H-M, Knauss E, Malleswaran I, Dinakaran S (2023) An investigation of challenges encountered when specifying training data and runtime monitors for safety critical ml applications. In: International working conference on requirements engineering: foundation for software quality, Springer, pp 206–222
26. Farrell M, Mavridou A, Schumann J (2023) Exploring requirements for software that learns: A research preview. In: International working conference on requirements engineering: foundation for software quality, Springer, pp 179–188
27. Villamizar, H., Kalinowski, M., et al. (2022) A catalogue of concerns for specifying machine learning-enabled systems. *arXiv preprint arXiv:2204.07662*
28. Al Islam MN, Ma Y, Alarcon P, Chawla N, Cleland-Huang J (2022) Resam: Requirements elicitation and specification for deep-learning anomaly models with applications to uav flight controllers. In: 2022 IEEE 30th international requirements engineering conference (RE), IEEE, pp 153–165
29. Liebel G, Tichy M, Knauss E, Ljungkrantz O, Stieglbauer G (2018) Organisation and communication problems in automotive requirements engineering. *Requirements Eng* 23(1):145–167
30. Pernstål J, Gorschek T, Feldt R, Florén D (2013) Software process improvement in inter-departmental development of software-intensive automotive systems—a case study. In: International conference on product focused software process improvement, Springer, pp 93–107
31. Allmann C, Winkler L, Kölzow T et al (2006) The requirements engineering gap in the oem-supplier relationship. *J Univ Knowl Manag* 1(2):103–111
32. M. Mahally M, Staron M, Bosch J (2015) Barriers and enablers for shortening software development lead-time in mechatronics organizations: A case study. In: Proc. of the 2015 10th joint meeting on foundations of SE, pp 1006–1009
33. Staron M (2019) Requirements engineering for automotive embedded systems. In: Dajsuren Y, van den Brand M (eds) *Automotive systems and software engineering*. Springer, Cham, pp 11–28
34. Ribeiro QA, Ribeiro M, Castro J (2022) Requirements engineering for autonomous vehicles: a systematic literature review. In: Proc. of the 37th ACM/SIGAPP symposium on applied computing, pp 1299–1308
35. Heyn H-M, Subbiah P, Linder J, Knauss E, Eriksson O (2022) Setting AI in context: A case study on defining the context and operational design domain for automated driving. In: International working conference on re: foundation for software quality, Springer, pp 199–215
36. Ågren SM, Knauss E, Heldal R, Pelliccione P, Malmqvist G, Bodén J (2019) The impact of requirements on systems development speed: a multiple-case study in automotive. *Requir Eng* 24(3):315–340
37. Zhang X, Tao J, Tan K, Törngren M, Sanchez JMG, Ramli MR, Tao X, Gyllenhammar M, Wotawa F, Mohan N et al (2022) Finding critical scenarios for automated driving systems: a systematic mapping study. *IEEE Trans Softw Eng* 49(3):991–1026
38. Luo Y, Zhang X-Y, Arcaini P, Jin Z, Zhao H, Zhang L, Ishikawa F (2022) Hierarchical assessment of safety requirements for configurations of autonomous driving systems. In: 2022 IEEE 30th international requirements engineering conference (RE), IEEE, pp 88–100
39. Zhang R, Albrecht A, Kausch J, Putzer HJ, Geipel T, Halady P (2021) Dde process: A requirements engineering approach for machine learning in automated driving. In: 2021 IEEE 29th international requirements engineering conference (RE), IEEE, pp 269–279
40. Felderer M, Ramler R (2021) Quality assurance for ai-based systems: Overview and challenges (introduction to interactive session). In: *Software quality: future perspectives on software engineering quality: 13th international conference, SWQD 2021, Vienna, Austria, January 19–21, 2021, Proceedings vol 13*, Springer, pp 33–42
41. Horkoff J (2019) Non-functional requirements for machine learning: challenges and new directions. In: 2019 IEEE 27th international requirements engineering conference (RE), IEEE, pp 386–391
42. Habibullah KM, Horkoff J Non-functional requirements for machine learning: understanding current use and challenges in industry. In: 2021 IEEE 29th International RE Conference (RE), pp. 13–23 (2021). IEEE
43. Habibullah KM, Gay G, Horkoff J (2023) Non-functional requirements for machine learning: understanding current use and challenges among practitioners. *Requir Eng* 28:1–34
44. Fujii G, Hamada K, Ishikawa F, Masuda S, Matsuya M, Myojin T, Nishi Y, Ogawa H, Toku T, Tokumoto S et al (2020) Guidelines for quality assurance of machine learning-based artificial intelligence. *Int J Softw Eng Knowl Eng* 30(11n12):1589–1606

45. Winner H, Lemmer K, Form T, Mazzega J (2019) Pegasus-first steps for the safe introduction of automated driving. In: Meyer G, Beiker S (eds) Road vehicle automation. Springer, Berlin, pp 185–195
46. Falcini F, Lami G (2017) Challenges in certification of autonomous driving systems. In: 2017 IEEE international symposium on software reliability engineering workshops (ISSREW), IEEE, pp. 286–293
47. Jenn E, Albore A, Mamalet F, Flandin G, Gabreau C, Delseny H, Gauffriau A, Bonnin H, Alecu L, Pirard J, et al. (2020) Identifying challenges to the certification of machine learning for safety critical systems. In: European congress on embedded real time systems (ERTS 2020)
48. Fisher, M., Collins, E., Dennis, L., Luckcuck, M., Webster, M., Jump, M., Page, V., Patchett, C., Dinmohammadi, F., Flynn, D., et al. (2018) Verifiable self-certifying autonomous systems. In: 2018 IEEE international symposium on software reliability engineering workshops (ISSREW), IEEE, pp 341–348
49. Jain A, Patel H, Nagalapatti L, Gupta N, Mehta S, Guttula SC, Mujumdar S, Afzal S, Mittal RS, Munigala V (2020) Overview and importance of data quality for machine learning tasks. In: Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining
50. Poth A, Meyer B, Schlicht P, Riel A (2020) Quality assurance for machine learning—an approach to function and system safeguarding. In: 2020 IEEE 20th international conference on software quality, reliability and security (QRS), IEEE, pp 22–29
51. Giray G (2021) A software engineering perspective on engineering machine learning systems: state of the art and challenges. *J Syst Softw* 180:111031
52. Hesenius M, Schwenzfeier N, Meyer O, Koop W, Gruhn V (2019) Towards a software engineering process for developing data-driven applications. In: 2019 IEEE/ACM 7th international workshop on realizing artificial intelligence synergies in software engineering (RAISE), IEEE, pp 35–41
53. Nahar N, Zhou S, Lewis G, Kästner C (2022) Collaboration challenges in building ml-enabled systems: Communication, documentation, engineering, and process. In: Proceedings of the 44th international conference on software engineering, pp 413–425
54. Adler R, Feth P, Schneider D (2016) Safety engineering for autonomous vehicles. In: 2016 46th Annual IEEE/IFIP International conference on dependable systems and networks workshop (DSN-W), IEEE, pp 200–205
55. Saldaña J (2021) The coding manual for qualitative researchers. The coding manual for qualitative researchers pp 1–440
56. Sutcliffe A (2003) Scenario-based requirements engineering. In: Proc.. 11th IEEE international RE conference, 2003., IEEE Computer Society, pp 320–320
57. Warg F, Skoglund M, Thorsén A, Johansson R, Brännström M, Gyllenhammar M, Sanfridson M (2020) The quantitative risk norm—a proposed tailoring of hara for ads. In: 2020 50th annual IEEE/IFIP international conference on dependable systems and networks workshops (DSN-W), IEEE, pp 86–93
58. Henriksson, J., Borg, M., Englund, C. (2018) Automotive safety and machine learning: Initial results from a study on how to adapt the iso 26262 safety standard. In: Proceedings of the 1st international workshop on software engineering for AI in autonomous systems, pp 47–49
59. Hawkins R, Paterson C, Picardi C, Jia Y, Calinescu R, Habli I (2021) Guidance on the assurance of machine learning in autonomous systems (amlas). [arXiv:2102.01564](https://arxiv.org/abs/2102.01564)
60. Koopman P (2023) UI 4600: what to include in an autonomous vehicle safety case. *Computer* 56(05):101–104
61. Johansson R, Sivencrona H (2021) Developing seooc-original concepts and implications when extending to ads. In: CARS 2021 6th international workshop on critical automotive applications: robustness & safety
62. Sommerville I (2011) Software engineering, 9/E. Pearson Education India, Bangalore
63. Kasauli R, Knauss E, Horkoff J, Liebel G, Oliveira Neto FG (2021) Requirements engineering challenges and practices in large-scale agile system development. *J Syst Softw* 172:110851
64. Nerur S, Mahapatra R, Mangalaraj G (2005) Challenges of migrating to agile methodologies. *Commun ACM* 48(5):72–78
65. Chan FK, Thong JY (2009) Acceptance of agile methodologies: a critical review and conceptual framework. *Decis Support Syst* 46(4):803–814
66. Acuna D, Philion J, Fidler S (2021) Towards optimal strategies for training self-driving perception models in simulation. *Adv Neural Inf Process Syst* 34:1686–1699
67. Wohlrab R, Steghöfer J-P, Knauss E, Maro S, Anjorin A (2016) Collaborative traceability management: challenges and opportunities. In: 2016 IEEE 24th international RE conference (RE), IEEE, pp 216–225
68. Jayatilke S, Lai R (2018) A systematic review of requirements change management. *Inf Softw Technol* 93:163–185
69. ISO: ISO 26262:2018 (2018) Road vehicles—functional safety. International Organization for Standardization, Geneva. www.iso.org
70. ISO: ISO/CD TS 5083 (2023) Safety for automated driving systems—design, verification and validation, under development. International Organization for Standardization, Geneva. www.iso.org
71. Chung L, Nixon BA, Yu E, Mylopoulos J (2012) Non-functional requirements in software engineering, vol 5. Springer, Berlin
72. Czarnecki K (2018) Operational design domain for automated driving systems. Taxonomy of Basic Terms. Waterloo Intelligent Systems Engineering (WISE) Lab, University of Waterloo, Canada
73. Gyllenhammar M, Johansson R, Warg F, Chen D, Heyn H-M, Sanfridson M, Söderberg J, Thorsén A, Ursing S (2020) Towards an operational design domain that supports the safety argumentation of an automated driving system. In: 10th European congress on embedded real time systems (ERTS 2020)
74. Broens, T., Quartel, D., Van Sinderen, M. (2007) Capturing context requirements. In: Smart sensing and context: second european conference, EuroSSC 2007, Kendal, England, October 23–25, 2007. Proceedings vol 2. Springer, pp 223–238
75. Ali R, Dalpiaz F, Giorgini P (2010) A goal-based framework for contextual requirements modeling and analysis. *Requir Eng* 15:439–458
76. Maalej W, Nayebi M, Ruhe G (2019) Data-driven requirements engineering—an update. In: 2019 IEEE/ACM 41st international conference on software engineering: software engineering in practice (ICSE-SEIP), IEEE, pp 289–290
77. Torkar R, Gorschek T, Feldt R, Svahnberg M, Raja UA, Kamran K (2012) Requirements traceability: a systematic review and industry case study. *Int J Softw Eng Knowl Eng* 22(03):385–433
78. Njomou AT, Africa AJB, Adams B, Fokaefs M (2021) Msr4ml: Reconstructing artifact traceability in machine learning repositories. In: 2021 IEEE International conference on software analysis, evolution and reengineering (SANER), IEEE, pp 536–540
79. Kästner C (2022) Machine learning in production: from models to products
80. Seyff N, Maiden N, Karlsen K, Lockerbie J, Grünbacher P, Graf F, Ncube C (2009) Exploring how to use scenarios to discover requirements. *Requir Eng* 14:91–111
81. Boehm BW, Brown JR, Lipow M (1976) Quantitative evaluation of software quality. In: Proceedings of the 2nd international conference on software engineering, pp 592–605

82. Psaros AF, Meng X, Zou Z, Guo L, Karniadakis GE (2023) Uncertainty quantification in scientific machine learning: methods, metrics, and comparisons. *J Comput Phys* 477:111902
83. Coit DW (2003) Maximization of system reliability with a choice of redundancy strategies. *IIE Trans* 35(6):535–543
84. Sagi O, Rokach L (2018) Ensemble learning: a survey. *Wiley Interdiscip Rev Data Min Knowl Discov* 8(4):1249
85. Studer S, Bui TB, Drescher C, Hanuschkin A, Winkler L, Peters S, Müller K-R (2021) Towards crisp-ml (q): a machine learning process model with quality assurance methodology. *Mach Learn Knowl Extr* 3(2):392–413
86. Kreuzberger D, Kühl N, Hirschl S (2023) Machine learning operations (mlops): overview, definition, and architecture. *IEEE Access* 11:31866–31879
87. Steghöfer J-P, Knauss E, Horkoff J, Wohlrab R (2019) Challenges of scaled agile for safety-critical systems. In: Product-focused software process improvement: 20th international conference, PROFES 2019, Barcelona, Spain, November 27–29, 2019, Proceedings vol 20, Springer, pp 350–366
88. Zhang JM, Harman M, Ma L, Liu Y (2020) Machine learning testing: survey, landscapes and horizons. *IEEE Trans Softw Eng* 48(1):1–36
89. Hittmeir M, Ekelhart A, Mayer R (2019) On the utility of synthetic data: an empirical evaluation on machine learning tasks. In: Proceedings of the 14th international conference on availability, reliability and security, pp 1–6

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.