

Automotive Perception Software Development: An Empirical Investigation into Data, Annotation, and Ecosystem Challenges

Hans-Martin Heyn^{*§}, Khan Mohammad Habibullah^{*}, Eric Knauss^{*}, Jennifer Horkoff^{*}

Markus Borg[†], Alessia Knauss[‡], Polly Jing Li[¶]

^{*}Chalmers | University of Gothenburg, Sweden

[†]RISE Research Institutes of Sweden

[‡]Zenseact AB, Sweden, [¶]Kognic AB, Sweden

[§]Corresponding author, Hans-Martin.Heyn@gu.se

Abstract—Software that contains machine learning algorithms is an integral part of automotive perception, for example, in driving automation systems. The development of such software, specifically the training and validation of the machine learning components, require large annotated datasets. An industry of data and annotation services has emerged to serve the development of such data-intensive automotive software components. Wide-spread difficulties to specify data and annotation needs challenge collaborations between OEMs (Original Equipment Manufacturers) and their suppliers of software components, data, and annotations.

This paper investigates the reasons for these difficulties for practitioners in the Swedish automotive industry to arrive at clear specifications for data and annotations. The results from an interview study show that a lack of effective metrics for data quality aspects, ambiguities in the way of working, unclear definitions of annotation quality, and deficits in the business ecosystems are causes for the difficulty in deriving the specifications. We provide a list of recommendations that can mitigate challenges when deriving specifications and we propose future research opportunities to overcome these challenges. Our work contributes towards the on-going research on accountability of machine learning as applied to complex software systems, especially for high-stake applications such as automated driving.

Index Terms—accountability, annotations, data, ecosystems, machine learning, requirements specification

I. INTRODUCTION

Driving automation refers to system that can automatically intervene in the driving task [1]. This includes advanced driver assistance systems (ADAS) which can be seen as a pre-stage to conditional or even full autonomous driving [2]. The aim of ADAS is to provide comfort and especially additional safety to manual driving tasks because the majority of accidents are still caused by human error [3]. Prominent features of ADAS include, among others, *collision avoidance*, *lane detection*, *traffic sign recognition*, *pedestrian detection*, *parking assistance* and *driver monitoring* [4]. All these features rely on the availability of data from a great variety of sensors, e.g., cameras, radar, lidar, ultrasonic sensors, fused and processed in real-time in what is referred to as the *perception system*. Besides rule-based software components, the automotive industry relies on machine learning (ML) to enable the perception

system of ADAS to be fast, robust, and precise enough in processing the incoming data [5]. The advent of ML in the automotive industry, however, has caused a paradigm shift, because software engineers no longer express all logic in source code. Instead, they train ML models with large, often pre-annotated datasets. Traditional processes for specifying, developing, and testing automotive software can no longer apply. Correctly working software that incorporates ML algorithms require not only the avoidance of systematic mistakes during software development but also sets expectations on the datasets available at design- and run-time [6], including the annotation of such data.

The datasets must fulfil the desired expectations in order to fulfil the desired performance of the software. In our pre-study *Precog* we explored how practitioners specify expectations on ML models, data, and data annotations as well as which trade-offs are taken for different quality aspects in software serving automotive perception systems, see [7] for more detailed information. Through an interview-based study with industrial practitioners from the Swedish automotive industry

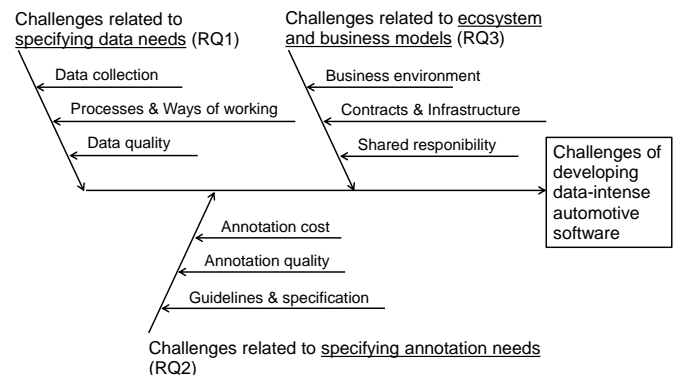


Fig. 1: Cause-Effect diagram showing the major themes regarding the ability to specify data and annotation needs, and the ability of ecosystems and business models to handle shared data-intensive software development.

we identified challenges in eight themes: *AI and ML models, annotation, data, ecosystem and business, quality, requirement engineering, perception, and system and software engineering*. In this article, we investigate closer the challenges identified in the data, annotation, and ecosystem and business themes. The reason we combine these themes together is the current advent of a *data industry* in the automotive supplier sector [8]–[10]. Original Equipment Manufacturers (OEMs) rely heavily on the collaboration with suppliers for the development of vehicle systems and software. The success of such collaborations depend on the ability of OEMs to specify expectations towards its suppliers. For the shared development of data-intensive software, this implies that both the OEM and the supplier must be able to specify and understand expectations on data, expectations on the annotation of data, and to maintain suitable interactions and collaborations for handling data-intensive software development.

The study is guided by three research questions:

- RQ1** What challenges do practitioners experience when specifying data needed for the development of automotive perception software that include machine learning components?
- RQ2** What challenges do the same practitioners experience when specifying annotation needs for data as part of the software development process?
- RQ3** What are implications towards industry ecosystems and business models for handling shared development of data-intensive software?

Figure 1 gives an overview of the main themes we found in relation to the research questions. Concerning RQ1, we found that the ability to specify data needs is negatively impacted by nontransparent data selection as part of the data collection process, missing process guidelines and a lack of common metrics describing data variation as a means of representing data quality. In regards to RQ2, the most critical challenges we found are inconsistent manual annotations and missing specifications and guidelines for the annotation processes. In answering RQ3, we found that, in relation to the business environment, conventional value chains and sourcing policies impede shared data-intensive software developments. We saw a trend towards sharing development tools and utilising open source policies. Furthermore, we saw that new forms of collaborated development and contracts are required to facilitate transparency and shared responsibility in data driven developments.

II. RELATED WORK

Rahimi et al. called for more attention from the requirement engineering community towards the ability of specifying of, what they referred to as, *Machine-Learned Components* (MLC) [11]. They explicitly mentioned datasets as an own aspect of MLCs that need to be properly specified for achieving a desired outcome of the MLC. Especially in safety-critical perception systems, requirements need to be specified towards robustness of the MLC [12]. Robustness is achieved if *small* changes in the input images do not lead to undesired behaviour.

However, often it is not clearly specified what *small* changes in the input space entail [13]. In an interview study with data scientists, Vogelsang and Borg identified gaps in mutual understanding of technical concepts and measures between customers and data scientists who prepare data for ML models as a consequence of the lack of proper data specifications [14]. The problem of missing data specifications becomes apparent if software application-specific requirements are to be incorporate as prior domain knowledge into the training dataset [15]. Another consequence of a lack of context-based specification for datasets is that the datasets tend to become intractably large which makes it impossible to scrutinise their content [16].

Paullada et al. elaborate on several challenges encountered in data handling in ML research that also seem to apply towards commercial software application of ML [17]. Jo and Gebru, for example, argue that unspecified data collection resembles a *wild west* mentality resulting in the risk of bias in the datasets. Instead, they propose the use of documentation methods from archiving, such as *mission statements* and *process records* [18]. Similar to model cards for ML [19], datasheets for data have been proposed as a first step towards data specifications [20].

Besides a lack of data specifications, there is also a lack of specifications for the annotations of the data. Most commonly, data for ML training is annotated manually, sometimes even with the help of *crowd-working* platforms. The use of crowd-working however can obscure the annotation process [21]. A clear task design is required to avoid human-induced errors which again calls for a proper requirements specification [22]. Even the definition of success metrics impact the result of the annotation process: for example a high annotation *accuracy* not necessarily results in high correctness of model predictions [23]. In automotive use cases of ML, the under-specification of data and annotations causes ambiguities in the requirements towards data which has negative implications towards the verifiability of safety-critical software that uses ML [6], [24]. For example, internal audits that should ensure correct behaviour of a ML model in relation to a company’s ethical values cannot be conducted without proper specification of the data collection and processing [25]. Accountability can only be established through transparency and ownership of the dataset development lifecycle which requires rigorous documentation of each stage in it [26]. Initial attempts to create large publicly available datasets are for example nuScenes for autonomous driving [27].

Knowledge sharing, quality definition, and effective communication are some necessary aspect of running a shared software development ecosystems [28]. However, in cross-company software development projects, a lack in documented domain assumptions and missing cross-organisational documentation, e.g., data specifications, during data collection, labelling, and cleaning have been identified as a major cause for failures of the resulting software that contains ML models [29].

III. METHOD

The study bases on a series of group interviews and a workshop conducted with different stakeholders involved in the development of perception systems for driving automation in the Swedish automotive industry. The reasons for choosing the Swedish automotive industry was on the one hand convenience of having local access to the stakeholders, and on the other hand the high competitiveness and international interdependence of the Swedish automotive industry.

With some Swedish OEMs being subsidiaries of larger international companies, such as Geely or Volkswagen AG, or being themselves large international players with several subsidiaries, such as Volvo Group AB, the Swedish automotive industry provides in our opinion a representative sample of the worldwide automotive industry.

Figure 2 depicts the research process followed during the study.

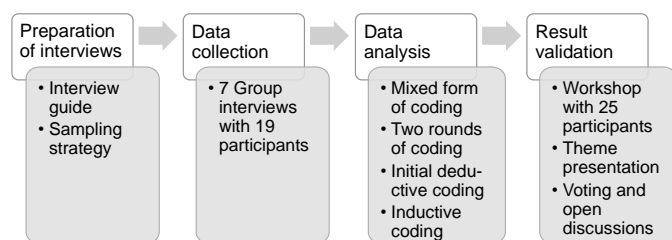


Fig. 2: Overview of the interview study

A. Preparation of interviews

The guiding theme of the interviews was the exploration of requirements and software engineering practises for automotive perception software that incorporates ML. Based on a-priori formulated research questions, the involved researchers created an interview guide¹ with nine parts, whereof two parts were optional depending on available time. The first part aimed at collecting information about the background of the interviewees. In the second part, we used a diagram⁴ of a typical architecture for automotive perception systems and asked the interviewees to position themselves in terms of the architecture. The aim was to find common ground with the interviewees and to understand their typical field of work. In the next part, we concentrated on discussing the processes used to ensure correctness of the perception system. Figure 3 was used to discuss with the interviewees how requirements on the quality of the function relate to requirements on data and annotation needs. Part four of the interview guide contains questions regarding the safety case, and how ML models can become a part of the safety case for software components. The next part contained questions regarding the context description in which a safety case of the software component is valid. Part six contained questions about the ecosystem and business environment in which the perception software is developed. Here, we were curious in how safety critical software that

¹accessible in the replication package <https://doi.org/10.7910/DVN/HCMVLI>

includes ML models are developed together with partners and suppliers and how expectations in the software, the ML model and data are communicated. The next two parts were optional and contained questions regarding quality trade-offs. The final part closed the interview and allowed the interviewees to suggest improvements and additional interviewees.

Sampling strategy: Our sampling strategy was a mix of purposeful and snowball sampling. For the latter, we sent open calls to contacts in the Swedish automotive industry, including the partners involved in the *Precog* study, and explicitly asked both before and after the interview for additional contacts we could interview. With respect to purposeful sampling we aimed to interview practitioners from industry who had experience with developing software for automotive perception systems, data science, machine learning, requirement engineering, and safety engineering. Because we knew that we cannot find all these qualifications in a single person, we decided to conduct group interviews to cover all desired competencies. We also hoped to enable better discussions in a small group settings.

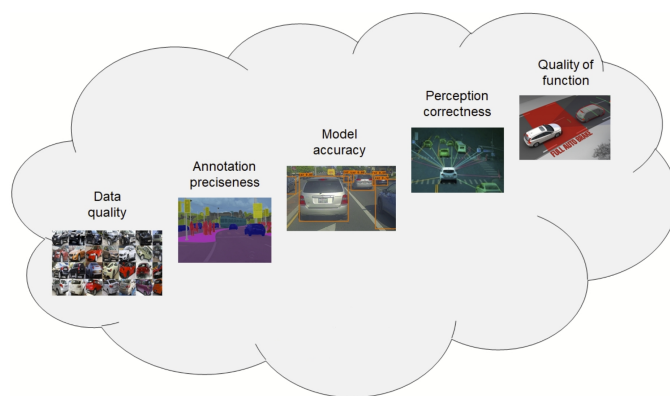


Fig. 3: Building a safety case for automotive perception software

B. Data collection

We conducted two campaigns for data collections: Group interviews and a validation workshop.

Interviews: In seven interview session we interviewed 19 participants from five companies. Three international automotive OEMs, one supplier, and one annotation company participated in the study. In each session, at least two scientists were present. The interview sessions were conducted between December 2021 and April 2022 via Microsoft Teams and took 90 to 120 minutes. A list of the participants for each interview is given in Table I. All interviews were recorded, automatically transcribed, anonymised, mistakes in the automatic transcripts manually corrected, and the final transcripts stored as spreadsheets for further analysis.

In the invitation e-mail to each group of interviewees we informed about the research project's goal, financial support, and duration. In the beginning of the interview, we provided some formal information such as data privacy compliance and permission for recording and data processing. Then, the interview guide was followed. For each section of the

TABLE I: Overview of conducted interviews

Inter-view	Field of work	Participants
A	Object detection	Product owner Product owner Test engineer
B	Autonomous Driving	AI engineer Software developer System architect
C	Vision systems	Product owner Requirement engineer Deep learning engineer
D	AD and ADAS	System engineer Manager AD ¹
E	Testing and validation AD ¹	System architect Product owner Product owner Compliance officer Data Scientist
F	Data annotations	AI engineer Data scientist
G	Autonomous Driving	System safety engineer

¹The participants' IDs are random and not shown to ensure confidentiality.

interview guide one interviewer asked the questions, while the remaining interviewers observed and took notes.

The interviews were semi-structured using the set of pre-determined open-ended questions formulated in the interview guide. However, we allowed deviations from the questions and the order of questions to facilitate discussions, also among the group of interviewees.

C. Data analysis

We applied a mixed coding strategy for analysing the qualitative data obtained through the interviews. Mixed coding strategies can be suitable for settings in which coding of in-depth interviews is conducted by teams using a shared platform, such as in our case Microsoft Office 365 [30]. Each coding team consisted of at least three researchers who conducted the first round of coding together. The team started with a number of high-level deductive codes which were based on the interview questions and researchers' experience. Then, while applying the deductive codes, new codes emerged as part of an inductive coding scheme. These emerging codes were added to a shared list used by all coding teams. After five interviews, we observed saturation by noticing that not many new inductive codes emerged in the following interviews. In a second round of coding, a new group of at least two researchers first revisited each interview and then applied the final list of emerged codes. Afterwards, pattern coding was used to identify emerging themes and sub-categories [31]. Finally, each statement of the interviewees was assigned to the identified sub-categories. The final codes of each interview were reviewed by an additional independent researcher.

D. Result validation

Towards the end of the *Precog* study in April 2022, we conducted a 2.5-hour workshop with 25 participants from industry, where five participants were also interviewees. The aim of the workshop was to validate and discuss the preliminary findings of the interview study. The workshop was conducted on-site

in Göteborg, Sweden, but with the possibility of connecting remotely using the Zoom conference software. During the workshop we first presented the identified sub-categories for each theme on the online whiteboard platform Miro. One theme at a time we let the participants vote on which categories are most important or relevant to them. The participants could give a maximum of one vote to each category and they were allowed to vote for more than one category in each theme. We use these results to gauge the relative importance of our discovered sub-themes according to participating practitioners.

IV. RESULTS

This section presents the identified themes around challenges when specifying data and annotations of data for automotive perception software that incorporates ML. Furthermore, this section presents themes regarding the ability of automotive industry's ecosystem and business models to handle data-intensive developments, such as the design, developing and deployment of software that incorporates ML models.

A. RQ1: The ability to specify data for the development of automotive perception software

Unlike conventional software systems in which rules specify the desired behaviour, a software component that incorporates ML infers these rules from data. Therefore, data plays a prominent role when trying to ensure correct behaviour of ML-based software for perception systems. If the data provided to an ML algorithm is biased, the resulting ML model will learn that bias, and consequently the decisions of the perception software will be biased.

During the interviews we wanted to learn about the interviewees' understanding of "data quality". We furthermore asked about the processes and ways of working used in relation to specifying and collecting data for the development of data-intensive software.

Figure 1 illustrates the three major themes that we identified in the interviews regarding the ability to specify data used for the development of software for perception systems: 1) Data collection, 2) Processes & Way of Working, and 3) Data quality. For each themes, a set of sub-categories was identified. Table II lists these categories. The indicated score is the percentage of votes given for a sub-category out of all votes given for a theme.

1) *Data collection*: Ensuring the correct behaviour of software that contains ML models requires a highly data-driven development. The participants of the validation workshop ranked the sub-category *data selection* as a key aspect of data collection to ensure safe and correct behaviour. Often, the right dataset as input to ML algorithms is found through a set of iterations between the ML experts and the data scientists. For example, uncertainty measures can be used to decide which additional data needs to be selected to reduce uncertainty:

"To ensure some form of safety measures from the model, we produce uncertainty estimations from outputs. Those are used in the data

TABLE II: Overview of all themes and sub-categories for RQ1: Challenges affecting the ability to specify data used for data-intensive software development (n=46), n is the number of submitted votes

ID	Description	Score
D1	Data collection	
-I	data selection	13%
-II	simulation	9%
-III	data collection	7%
-IV	metadata	4%
-V	experimentation	2%
-VI	synthetic data	2%
D2	Processes and Ways of Working	
-I	data specification	7%
-II	data requirements	4%
-III	data verification	4%
-IV	data cleaning	2%
-V	data collaboration	2%
-VI	data storage	0%
D3	Data quality	
-I	data variation	11%
-II	bias	9%
-III	future-proof dataset	9%
-IV	data quality	7%
-V	data correctness	4%
-VI	data re-usability	4%
-VII	data maintainability	0%

selection of course to look for what type of data are we uncertain. What do we need to learn more from?" - Interviewee B-I

Because finding the right data can be expensive, development teams try to use simulations for training and validation purposes. The participants even consider data from simulation more important than data originating from planned experiments, such as test drives.

"Furthermore, simulations are very often an integral part of test strategies for machine learning based systems." - Interviewee C-IV

There can be two reasons for preferring data from simulation: It is often significantly cheaper to obtain data through simulations. And, it is possible to obtain data for rare case scenarios that could be impossible to obtain in real world experiments:

"Especially like for the rare case scenario that is not really easy to replicate in real world, so we cannot." - Interviewee E-II

2) *Processes and Way of working*: We asked during the interviews which processes and ways of working in regard to data used in the development of safety-critical data-intensive software is being used. Most safety standards, such as ISO 26262, rely on the correctness of processes to build up a safety case for a product. Therefore, it is important to define processes that ensure the safety of software with ML components.

The most important capability, defined formally through a process or informally through a way of working, is the creation of data specifications. We saw earlier that data selection is the most important activity in a data driven development process. Data specifications are a logical prerequisite to data selection.

Ideally, data specifications precise the requirements set on the data in relation to for example physical data properties, data quality, and quantitative targets [14]. But it is often quite unclear what a data specification should entail:

"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 do want to combine certain objects." - Interviewee B-I

An iterative processes can be used to find the final data specification of the system:

"It's more of a sort of a data driven and then statistical analysis of the data in a continuous way. So we start with logging data, annotating it, training or models, and so on. And then we can also draw some statistics. OK, how is the class balance in this dataset? How does that affect the per class accuracy? Do we need to look for more? And then we can of course feedback that to the data selection team and they can start looking for certain classes, for example." - Interviewee B-I

The statement shows that a specification typically consists of a set of requirements such as accuracy, balance, etc. We found however that the exact scope of a data requirement is not entirely clear. Data requirements can for example describe desired probability distributions and quantity of the data:

"We write documents, word documents, basically where we describe the distribution of the data and the quantity of the data that needs to be collected." - Interviewee C-II

Data requirements can also entail specific data quality aspects, such as pixel density, brightness, size of bounding boxes, etc. In both understandings of data requirements, they allow for data verification. Data verification means checking that the data is representative for the desired "real-world" scenario.

Data specifications and data requirements entailed within the specifications are key enablers for many companies to collaborate in data collection and processing for example with supplier companies:

"We have a 3rd party company driving around all this mileage and collecting data. They want you to send over that data to them for doing the simulations. And then they will put their requirements on what sort of data we are collecting." - Interviewee D-I

A final data specification that describes the utilised data can be a key input towards a safety case of a ML-based perception system and allows for verification of design decision.

3) *Data quality*: Data requirements often entail some desired data quality aspect. Interestingly, the most important data quality aspects mentioned by the interviewees do not describe physical properties of data, such as pixel density, contrast, resolution, brightness, etc., but instead focus on the represented information in the data.

Because data selection was identified as the most important aspect for data collection, it is not surprising that data variation has been chosen by the study participants as the most impor-

tant data quality characteristic, even before data correctness. Data variation is directly causally related to bias in data; a lack of data variation will result in bias, which can propagate into the ML model.

A challenge regarding data variation is the definition of KPIs, or in general measures of variety. How do you measure variety, and when do you know that your data has enough variety?

“How would you divide that space and define it in a way that allows a measure of have I covered not only enough children, but also **enough variety** of children [as vulnerable road users]?” - Interviewee F-I

Both data variation and data correctness require however an a-priori understanding of the environment in which the software will be deployed. If the operational domain is unknown, it will be difficult to describe what variety entails:

“You need to **understand the distribution of where to collect data and that requires an understanding of where the function in the end will be used.**” - Interviewee C-I

Collecting and processing data is often a costly part of the development process. Therefore, the re-usability, and future-proofness of data are considered important data quality aspects.

“What do we need to ensure to make use of the data we’ve collected up to now? I mean, **how do we make sure we don’t have to start from scratch?**” - Interviewee E-I

B. RQ2: The ability to specify annotations for data used in automotive perception software

The development of ML models often relies on supervised learning which requires annotated datasets. There are approaches towards automatised annotation of data² but these approaches regularly do not succeed in replacing human-in-the-loop annotators [35]. Because annotation plays a major role in the development of perception software that incorporates ML, both from a performance and cost point-of-view, we investigated which challenges practitioners encountered when specifying annotations for data.

Figure 1 shows the three major themes we identified within annotation challenges. Refined sub-categories for each theme together with the score they received in the validation workshop are listed in Table III.

1) **Annotation costs:** The final cost of annotation often is described as a trade-off between annotation quality and quantity. All interview partners agree that the cost of annotation rises exponentially with the level of annotation quality and linearly with the quantity of annotated data. Both a higher quality of annotations and a higher quantity of annotated data can result in performance increases of the trained ML model:

²automated annotation for example has been attempted for images [32], for textual data [33], or videos [34]

TABLE III: Overview of all themes and sub-categories for RQ2: Challenges affecting the ability to specify annotations for data used for data-intensive software development (n=31)

ID	Description	Score
A1	Annotation costs	
-I	annotation cost	3%
A2	Annotation quality	
-I	annotation consistency	16%
-II	annotation correctness	13%
-III	annotation quality	13%
-IV	annotation validation	10%
-V	annotation re-usability	6%
-VI	annotation precision	3%
-VII	pixel precision	0%
A3	Guidelines & Specification	
-I	ground truth	13%
-II	annotation specification	10%
-III	annotation guidelines	6%
-IV	labelling	6%

“And there’s some kind of scaling. So if you don’t have a lot of frames [annotated], your model performance will be way worse. And if you would have a higher quality then you would of course get a higher performance. But at some point there’s a diminishing of returns.” - Interviewee F-II

Given a fixed budget for annotation, a trade-off is described as choosing between either high quantity of data with low quality annotations, or vice versa:

“[...] OK if you have higher quality then maybe you can do with less [annotated] data instead. And then you save money from one perspective. But yeah, it takes longer time to do one precise annotation. So that **you have to balance** the two a bit.” - Interviewee C-II

2) **Annotation quality:** Annotations are regarded a *quality aspect of data*, but there is no clear distinction in what this quality aspect entails. It can refer to qualitative aspects of annotations such as the precision of annotation boxes or the correctness of annotations. But it can also refer to quantitative aspects, such as the amount of annotated data or how many features are labelled within a single frame. Furthermore, there is no distinct definition of *annotation quality*. The uncertainty in the definition of annotation quality can stem from the inability to define clear quality metrics which can be used for setting requirements on the annotation (similar to the inability to define metrics for data variety as a data quality criteria):

“And just to add, it’s not really clear **how we can measure the quality of annotations itself** as well. Like how to make sure that even like if you put a requirement on the annotations and they have reached their level of quality that you asked for.” - Interviewee D-I

The uncertainty in the specifications of annotation can result in uncertainties and even consistency problems in the annotated dataset:

*“I think it’s a pitfall. Maybe that it’s easy to look at, you know, like pixel precision: ah, you are two pixels out of the actual border of the object here, but I think maybe we’ve seen a bigger problem that **one type of object has one class in one label and another class in the next image** because hundreds of annotation people have interpreted the specifications differently. I think for us that’s a bigger problem in annotation quality than the pixel precision.”* - Interviewee B-II

Being able to provide consistent annotation has been ranked as the most pressing challenge by the validation workshop participants. There are approaches described in literature that in theory allow for testing consistency between annotators (for example [36]), but it seems not to be used in practise yet. At the same time, pixel precision is not considered an equally pressing challenge:

*“It’s **much easier to solve the preciseness problem than the consistency problem**. It’s our experience.”* - Interviewee B-I

Furthermore, as a consequence of high annotation costs there is a desire to re-use annotated data:

*“[...] we still would like to use that data because **I paid a lot of money to annotate it** and then we can do different things [...] that actually contributes to robustness.”* - Interviewee A-I

3) *Guidelines & Specification*: In theory, specifications serve two purposes: A *requirements specification* documents the requirements that need to be fulfilled by an item and a *technical specification* can document the features of an item that fulfil the desired requirements. They are key communication artefacts between OEMs and suppliers or other external companies [37]. More and more, annotations are conducted by external companies. Yet, for the process of annotating data, clear requirements specifications often are not formulated. This results in ambiguous expectations on the resulting annotations:

*“[A]s soon as something is not explicitly stated, it’s they [the annotators] don’t know what to do because we experience [that] they can extrapolate what they know, but **there’s still a lot of these either ambiguous or hard to tell scenarios that are usually quite unprecise**.”* - Interviewee F-I

As a consequence of missing requirements specification for annotation, there is ambiguity in the guidelines that describe the annotation process. For example, how should uncertainties during annotation be handled?

*“But which are still sort of not properly described or it’s not properly defined how to act in that situation. [...] What is annotated should it be [marked as] correct, or is it better to mark [it] as unclear? **Maybe it’s worse to be wrong than to be not sure in this case**.”* - Interviewee D-I

The question about certainty in the annotation is relevant, because many annotation processes apply time constraints towards the annotators. With a fixed time budget, for example a maximum of 20 seconds per image, the quality of the anno-

tations might be significant worse than if no time constraint is applied in the annotation process because the annotator has more time to concentrate on the details of each frame:

*“In one way it’s a **trade off of the time taken to peruse annotations versus the actual quality**. So I mean, for example, if we say you have one minute to do each task, or if you have, how however much time you need. We would naturally assume that the later approach will enable us to create more detailed annotations.”* - Interviewee F-I

Finally, without *annotation specifications*, and therewithin information about the annotation process, the company receiving the annotated data from a subcontractor cannot judge on the reliability of the annotated data:

*“I mean it depends on what, what if they have it in house, or if they outsource it and so on. But when it comes down to the labelling, that means they have to tell us **what their actions or processes [are] and then how they guarantee the quality or efficiency of the labelling**. And then we have to take that judgement.”* - Interviewee D-I

The ability to provide annotation specifications is also of significance towards safety evidence. Because the annotation process has influence on the final performance and correct behaviour of a ML model that is part of a perception software, it needs to be clearly documented as part of a safety case.

C. RQ3: Automotive industry’s ecosystems and business models for data-intensive software developments

The previous section highlighted the importance of data and annotation specifications for the development of software that incorporates ML components in the automotive industry. Requirement documents, for example requirements specifications, play a major role in steering the process flow between OEMs and their suppliers [38]. Typically, major parts of a vehicle’s software are developed externally through suppliers in agreement with the requirement specifications they receive from the OEMs. However the typical specification driven approach fails in scenarios of complex system development, for example components with significant linkage between different systems [39], and highly data driven developments such as deep learning [40].

In this section we investigate causes of the challenges that affect the ability of the automotive industry’s ecosystems and business models to handle data intensive development. Based on the conducted interviews and the validation workshop we identified three major themes as shown in Figure 1 and twelve sub-categories listed in Table IV.

1) *Business environment*: The “conventional” value chain in the automotive industry is based on sourcing suppliers which provide the OEMs with the technology needed for their products. Because of the emergence of data-intensive software in vehicles and the agile transformation that embraced the automotive industry³, this type of partnership in the value chain cannot work anymore:

³where the agile transformation can be an effect of the introduction of data-intensive software [41]

TABLE IV: Overview of all themes and sub-categories for RQ3: Challenges affecting the ability of the automotive industry’s ecosystems and business models to handle data-intensive software development (n=46)

ID	Description	Score
B1	Business environment	
-I	value chain	7%
-II	ecosystem	4%
-III	feedback	4%
B2	Contracts & Infrastructure	
-I	tools	22%
-II	contracts	7%
-III	negotiations	4%
B3	Shared responsibility	
-I	transparency	20%
-II	collaboration	15%
-III	legal aspects	7%
-IV	risk of litigation	4%
-V	business model	4%
-VI	crowd sourcing	2%

“Yeah, **partnerships is a very hot word instead of suppliers**. So and I think that captures it, it cannot be a classical, what’s it called, purchaser supplier relation, where you just write what you want and then you get it. This is like the sourcing new managers dream world. **They get the blueprint, the drawing and then they go to five different vendors and then they pick the cheapest one. But it doesn’t work like that at all.** It was hard already with classical software. It’s like impossible here, so it needs to be more of a partnership where it’s possible to have a dialogue and to iterate without, you know, having to go through the whole commercial process once more, so it needs to be a bit more loose ended.” - Interviewee B-II

Instead of traditional sourcing, OEMs seek project partnerships for example with major technology companies. Alternatively, they partner up (and eventually buy) technology start-up companies, especially around data collection, annotation services, or automatic driving in general. The reason is that OEMs aim at integrating suppliers closer into their own development:

“[...] we did have a pretty extensive thorough sourcing project when we choose the suppliers that we’re working with now. And that was a major factor. We didn’t just pick the cheapest one or the one that we thought had their absolute best, you know, maybe accuracy. **We picked the one that we feel that we can work with these people and you know, we can have a dialogue and it’s possible to make adaption without getting into a commercial discussion.** You know, from the first minute. It’s more like a partnership spirit or setup.” - Interviewee B-II

The comment that an OEM is willing to pay a premium for joint teams development stems from the need of regular feedback in complex system development. Traditionally, suppliers only allow limited access to the development details in *joint reviews*. However, these discrete feedback points are not sufficient for an agile mindset required for data-intensive software projects:

“Normally with the suppliers, if you assume that they are a normal supplier, you cannot see those stuff unless you go to some kind of on site joint reviews because **it’s IP basically.**” - Interviewee G-I

We identified two major causes for the need of continuous feedback loops: First, the operational design domain (ODD), i.e., the context in which the system can operate as expected, is initially not entirely known. Instead, the true ODD is jointly “discovered” during the development.

“[...] it’s basically giving some feedback to the system design, function design and so on to modify the function and introduce the limitations to the function which we name it ODD. **An iterative loop is going on and on and on until all of these triggering events are acceptable basically.**” - Interviewee G-I

Second, the desired quality level of the software, especially in relation to safety and security, can only be achieved by continuously improving for example the data selection or annotations, even after the software has been deployed.

2) *Contracts & Infrastructure*: New forms of cooperation require new ways of *negotiating with* and *contracting of* suppliers, as well as tools that support continuous collaborations. A typical contract can contain function requirements that a supplier needs to fulfil. However, for data driven and probabilistic systems, function requirements are often very abstract and provide little guidance for system development:

“What I get [from customers] are usually those very abstract function requirements that I should never miss a single pedestrian, ever.[...] So what I try to do is that I break down and do decomposition [until] I **end up with feasible requirements on each component.**” - Interviewee C-III

The key issue with contracting of suppliers for data driven development is the definition of success. It seems difficult to define clear development goals due to the iterative nature of data driven development:

“Like **what’s the definition of success** when you’re really building a [machine learning] model and going back from that, then what’s the annotation preciseness you need” - Interviewee F-II

This inability to define success relates back to the inability of defining clear *metrics* and KPIs for data and annotation quality.

“And I mean it’s I think I can say they[, the OEMs,] push us [suppliers]. They want to have **more defined quality metrics** and so on, but it’s also very hard for us to come up with them.” - Interviewee C-II

An interesting problems arises around the *tools* used for the joint development of systems between suppliers and OEMs. Typically, in-house tools are proprietary and companies might be reluctant to share this intellectual property. A mentioned solution is the use and active development of open source tools:

“Yeah, so far we have not found any off the shelf product that can solve [what we need for collaboration], so **it’s more or less a combination of in-house tools and open source tools.**” - Interviewee C-I

The ease of collaboration through open source tools might explain the success of such tools in the field of ML and data driven development, even in a competitive environment such as the automotive industry.

3) *Shared responsibility*: New forms of partnerships between suppliers and OEMs causes shared responsibility for quality assurance, new *business models* for data and annotation services, while still keeping due diligence of the OEM. This has some consequences on how responsibility is organised between the partners. *Transparency* in the development processes becomes more important because building a safety case for perception software requires traceability and documentation of all design decisions, including traceability to the data and annotations used for the ML components:

“[...] talking from a safety case argumentation point of view; they will come, **they will ask for all the documentation and traceability** and they want to know what sort of process had you followed when it comes to machine learning as such. So we see a joint venture between us. To help out in the total vehicle certification point of view and when it comes to quality of the machine learning.” - Interviewee D-I

Building safety cases, and fulfilling the necessary quality criteria in both data and annotations is a *collaborative* effort. Furthermore, close collaboration also enables OEMs better to reduce *risk of litigation*, because in a collaborative environment, they better understand how the system was developed and what data has been used:

“We need to prove that we have done due diligence. So **it’s not good enough to just believe what our supplier says that it will work.** [...] we need to evaluate, how is that even possible and what data has been used and collect enough data in case there is a litigation and then we will also need to work with data, perhaps for a different business models.” - Interviewee B-III

V. DISCUSSION

The results of the study gave an impression of the major change that is currently affecting the automotive sector: Data intensive developments, such as the development of software that includes ML models, cannot be conducted in the same way as “traditional” software components. There is a lack of knowledge on how to properly specify aspects of data-intensive software. A notable finding of our study is that requirements specifications seem to play a major role in the sourcing and certification processes of automotive products, yet there is no common approach for specifying data or annotations of data. This is problematic in the sense that there are businesses emerging that specialise in the data procurement and data annotation for ML development, but automotive OEMs are not yet routinely able to collaborate with “data” companies in the same manner as they can with system or

components suppliers through established procurement processes and management.

A. Recommendations

The interviews and the results from the validation workshop provided insight into the current state of practise when specifying data and annotations with the aim of achieving acceptable performance for automotive perception software that incorporates ML. From our observations, based on the lessons learned in the interviews, and the indicative scoring from the workshop, Table V provides several recommendations for practitioners in the industry.

VI. THREATS TO VALIDITY

Threats to validity can arise from the interviews, the workshop, and the data analysis process. In this section, we discuss possible threats to validity, and how we implemented mechanisms to reduce them.

A. Threats to internal validity

Threats to internal validity arise when confounding variables cause bias in the result. This can occur through a lack of rigour (i.e., degree of control) in the study design [47]. We established several mechanisms to reduce potential confounding: The interview guide was internally peer-reviewed and a test session of the interview was conducted before starting data collection. Furthermore, to avoid personal bias, at least two authors conducted each interviews. One of the authors was present at all interviews, while the others authors took turns in joining the interviews. After each interview, the authors aligned their interviewing experience in group meetings. The workshop was lead by all researchers, and a briefing after the workshop was conducted to share and discuss impressions obtained during the workshop. Another potential bias can arise from the sampling process. We deployed a mixture of purposeful and snowball sampling for both the interviews and the workshop. We needed a certain set of expertise to answer our questions, yet we also allowed companies to suggest additional interview partners. The companies were contacted through an open call. Additionally, we actively approached all OEMs in Sweden and received participants suggested by them. Furthermore, the workshop reduced potential selection bias, because participants outside of the companies of the interview study were included. Another threat to validity arises when saturation is not reached in the collected data. We can argue that we reached a point of saturation because we noticed a sharp decline in emerging codes after analysing the fifth group interview.

B. Threats to external validity

Threats to external validity arise when generalisability of the research results cannot be guaranteed. To support generalisability of the results, a sampling strategy was chosen that included different roles on different levels and at a number of different companies of different size. However, our study-results and conclusions are limited to the automotive sector,

TABLE V: Recommendations for practitioners and towards researchers based on lessons-learnt in the Precog study

ID (sub-category)	Recommendation
Data-I (D1-I, D2-II, D3-I, B3-I)	Establish clear traceability of data selection decisions as input towards a safety case of software with ML. Comment: Data selection can have major implication on the correct behaviour of software that contains ML, because data selection strongly can influence bias and therefore should be traceable. Related work: [12], [20]
Data-II (D1-I, D2-I, B3-II)	Accept that a data specification can only be created iteratively. Comment: Data preparations and selection for data-intensive software developments are obviously highly data driven activities. Practitioners need to analyse the available data before being able to understand which additional data might be needed. Therefore, conventional OEM-supplier sourcing processes might be unsuitable for data intensive developments and need to be changed. Related work: [11], [14]
Data-III (D2-II, D2-III, D3-I, D3-IV, D3-V)	Establish common metrics on data variation and other relevant data quality aspects to facilitate clearer communication between companies. Comment: Data variation has been voted as the most important aspect of data quality, yet there is a lack of clear metrics that allow for specifying data variation. Related work: [25]
Annotation-I (A1-I, A2-I, A2-III)	Evaluate if an increase in annotation quality in lieu of an increase of the annotation quantity, i.e., the amount of annotated data, can result in better ML model performance. Comment: An increase in annotation quality seems to have stronger positive effects at the same cost compared to an increase in the amount of annotated data. Related work: [42], [43]
Annotation-II (A2-I, A2-II, A2-VI, A2-VII)	Concentrate on annotation consistency rather than pixel precision to increase annotation quality. Comment: According to our interviewees, inconsistent annotations have worse effects on ML model performance than variations in the precision of the bounding boxes. Related work: [23], [44]
Annotation-III (A2-IV, A3-I, A3-II, B3-I, B3-IV)	Clearly specify annotations and the annotation process. Comment: An annotation specification allows for judging the reliability of the annotated data, which can provide safety evidence and accountability towards the data-intensive software component. Related work: [6], [22], [24]
Ecosystems-I (B1-I, B2-II, B2-III, B3-II)	Avoid conventional automotive OEM-supplier sourcing processes in data-intensive developments. Comment: An upfront requirements specification is often not feasible for data-intensive developments, because “discovering” the right data and training a desired ML model is a highly iterative process. Instead, development partnerships with less bureaucracy are a trend mentioned by several interviewees. Related work: [45]
Ecosystems-II (B1-II, B2-I, B3-I, B3-II, B3-III, B3-VI)	Increase the use of open source tools in shared data intensive developments. Comment: Open source tools are from a legal perspective easier to share with new collaborators, they make it easier for smaller companies such as start-ups to participate in the development, and they establish transparency in the development process which can be positive when building a safety case or arguing for security. Related work: [46]

and specifically to the development of software for perception systems. However, we argue that perception system represent a typical situation in which a highly data-intensive software development is needed. Therefore, our results might also be valid for other data intensive development environments conducted in a more conservative business area like the automotive sector. An example can be the medical sector where ML plays more and more an important role in software for image based diagnostics.

VII. CONCLUSION AND OUTLOOK

This interview-based study investigated challenges that arise in the automotive industry when specifying data-intensive software components, such as software for perception systems. In seven group interviews with a total of 19 participants and through a validation workshop with 25 participants, we identified challenges that impact the ability to specify data and annotations of data. The inability to coherently measure data variation, unclear data collection processes, and the need of iterative development methodologies for data selection are examples of challenges that compromise the ability to specify data effectively for data depending software products in an automotive application. Unclear definition of annotation quality, a misleading focus on preciseness and quantity instead of consistency, and a lack of transparency in the annotation processes are examples of impediments that hinder proper annotation specifications. Furthermore, the study investigates current practises in the business environment and ecosystems deployed in the automotive industry, especially concerning a new trend towards emphasising joint development

projects over the traditional OEM-supplier relationship in data-intensive developments. We concluded this study by providing a number of recommendations based on our observations.

We expect a major change in how the automotive industry is going to collaborate with suppliers and other partners in the development of data-intensive systems. The results of our study suggest a number of further research topics: The problem of defining clear metrics for data quality aspects and annotation aspects, and how partners can agree on proper metrics is not solved. There is research needed in understanding how different quality aspects of annotations should be specified for achieving a desired ML model performance. Furthermore, the development of an “annotation industry” is in progress [48], and the success of these companies and “crowd-sourcing” approaches depend on the ability to collaborate fruitfully with established companies such as OEMs in the automotive sector. Currently, an emphasis is set on quantity and precision over consistency in annotations. We need to learn how a suitable trade-off between different annotation aspects can be achieved, such that the cost for developing the software components is minimised and the resulting performance maximised.

Based on the findings of this study, we propose further research 1) on how incremental collaborative specifications of data selection can be achieved, 2) on how such a specification for data selection can be validated, 3) on how the annotation process can be specified and eventually integrated in a safety assurance life-cycle, and 4) on how the information and knowledge sharing between OEMs and suppliers can be improved towards more joint responsibility in the development of machine learning models.

ACKNOWLEDGEMENTS

This project has received funding from Vinnova Sweden under the FFI program with grant agreement No 2021-02572 (precog), from the EU’s Horizon 2020 research and innovation program under grant agreement No 957197 (vedliot), and from a Swedish Research Council (VR) Project: Non-Functional Requirements for Machine Learning: Facilitating Continuous Quality Awareness (iNFORM). We are thankful to all interviewees and companies who supported us in this research.

REFERENCES

- [1] On-Road Automated Driving (ORAD) Committee, *Taxonomy and Definitions for Terms Related to Driving Automation Systems for On-Road Motor Vehicles*. Warrendale, U.S.: SAE International, 2021.
- [2] V. K. Kukkala, J. Tunnell, S. Pasricha, and T. Bradley, “Advanced driver-assistance systems: A path toward autonomous vehicles,” *IEEE Consumer Electronics Magazine*, vol. 7, no. 5, pp. 18–25, 2018.
- [3] A. J. Khattak, N. Ahmad, B. Wali, and E. Dumbaugh, “A taxonomy of driving errors and violations: Evidence from the naturalistic driving study,” *Accident Analysis & Prevention*, vol. 151, p. 105873, 2021.
- [4] M. I. Chacon-Murguía and C. Prieto-Resendiz, “Detecting driver drowsiness: A survey of system designs and technology,” *IEEE Consumer Electronics Magazine*, vol. 4, no. 4, pp. 107–119, 2015.
- [5] A. Moujahid, M. E. Tantaoui, M. D. Hina, A. Soukane, A. Ortalda, A. ElKhadimi, and A. Ramdane-Cherif, “Machine learning techniques in adas: a review,” in *2018 International Conference on Advances in Computing and Communication Engineering (ICACCE)*. IEEE, 2018, pp. 235–242.
- [6] M. Borg, C. Englund, K. Wnuk, B. Duran, C. Levandowski, S. Gao, Y. Tan, H. Kaijser, H. Lönn, and J. Törnqvist, “Safely entering the deep: A review of verification and validation for machine learning and a challenge elicitation in the automotive industry,” *Journal of Automotive Software Engineering*, vol. 1, no. 1, pp. 1–19, 2019.
- [7] Vinnova. Precog: Requirements engineering toward safe machine learning-based perception systems for autonomous mobility. [Online]. Available: <https://bit.ly/3SSGLaQ>
- [8] C. Salinesi, I. Kusumah, and C. Rohleder, “New approach for supporting future collaborative business in automotive industry,” in *2018 IEEE International Conference on Engineering, Technology and Innovation (ICE/ITMC)*. IEEE, 2018, pp. 1–9.
- [9] K. E. Martin, “Ethical issues in the big data industry,” in *Strategic Information Management*. Routledge, 2020, pp. 450–471.
- [10] D. Wang, S. Prabhat, and N. Sambasivan, “Whose ai dream? in search of the aspiration in data annotation,” in *CHI Conference on Human Factors in Computing Systems*, 2022, pp. 1–16.
- [11] M. Rahimi, J. L. Guo, S. Kokaly, and M. Chechik, “Toward requirements specification for machine-learned components,” in *2019 IEEE 27th International Requirements Engineering Conference Workshops (REW)*. IEEE, 2019, pp. 241–244.
- [12] B. C. Hu, R. Salay, K. Czarnecki, M. Rahimi, G. Selim, and M. Chechik, “Towards requirements specification for machine-learned perception based on human performance,” in *2020 IEEE Seventh International Workshop on Artificial Intelligence for Requirements Engineering (AIRE)*. IEEE, 2020, pp. 48–51.
- [13] B. C. Hu, L. Marssó, K. Czarnecki, R. Salay, S. Huakun, and M. Chechik, “If a human can see it, so should your system: Reliability requirements for machine vision components,” in *2022 IEEE 44th International Conference on Software Engineering (ICSE)*. IEEE, 2022, pp. 1145–1156.
- [14] A. Vogelsang and M. Borg, “Requirements engineering for machine learning: Perspectives from data scientists,” in *2019 IEEE 27th International Requirements Engineering Conference Workshops (REW)*. IEEE, 2019, pp. 245–251.
- [15] L. Zhou, S. Pan, J. Wang, and A. V. Vasilakos, “Machine learning on big data: Opportunities and challenges,” *Neurocomputing*, vol. 237, pp. 350–361, 2017.
- [16] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, “On the dangers of stochastic parrots: Can language models be too big?” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021, pp. 610–623.
- [17] A. Paullada, I. D. Raji, E. M. Bender, E. Denton, and A. Hanna, “Data and its (dis) contents: A survey of dataset development and use in machine learning research,” *Patterns*, vol. 2, no. 11, p. 100336, 2021.
- [18] E. S. Jo and T. Gebru, “Lessons from archives: Strategies for collecting sociocultural data in machine learning,” in *Proceedings of the 2020 conference on fairness, accountability, and transparency*, 2020, pp. 306–316.
- [19] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru, “Model cards for model reporting,” in *Proceedings of the conference on fairness, accountability, and transparency*, 2019, pp. 220–229.
- [20] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford, “Datashets for datasets,” *Communications of the ACM*, vol. 64, no. 12, pp. 86–92, 2021.
- [21] C. Hube, B. Fetahu, and U. Gadiraju, “Understanding and mitigating worker biases in the crowdsourced collection of subjective judgments,” in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp. 1–12.
- [22] J. W. Vaughan, “Making better use of the crowd: How crowdsourcing can advance machine learning research,” *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 7026–7071, 2017.
- [23] D. Tsipras, S. Santurkar, L. Engstrom, A. Ilyas, and A. Madry, “From imagenet to image classification: Contextualizing progress on benchmarks,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 9625–9635.
- [24] R. Salay and K. Czarnecki, “Using machine learning safely in automotive software: An assessment and adaptation of software process requirements in iso 26262,” *arXiv preprint arXiv:1808.01614*, 2018.
- [25] I. D. Raji, A. Smart, R. N. White, M. Mitchell, T. Gebru, B. Hutchinson, J. Smith-Loud, D. Theron, and P. Barnes, “Closing the ai accountability gap: Defining an end-to-end framework for internal algorithmic auditing,” in *Proceedings of the 2020 conference on fairness, accountability, and transparency*, 2020, pp. 33–44.
- [26] B. Hutchinson, A. Smart, A. Hanna, E. Denton, C. Greer, O. Kjartansson, P. Barnes, and M. Mitchell, “Towards accountability for machine learning datasets: Practices from software engineering and infrastructure,” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021, pp. 560–575.
- [27] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liang, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nusenes: A multimodal dataset for autonomous driving,” in *CVPR*, 2020.
- [28] C. Alves, J. A. P. de Oliveira, and S. Jansen, “Software ecosystems governance—a systematic literature review and research agenda,” *ICEIS (3)*, pp. 215–226, 2017.
- [29] N. Sambasivan, S. Kapania, H. Highfill, D. Akrong, P. Paritosh, and L. M. Aroyo, ““everyone wants to do the model work, not the data work”: Data cascades in high-stakes ai,” in *proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021, pp. 1–15.
- [30] N. M. Deterding and M. C. Waters, “Flexible coding of in-depth interviews: A twenty-first-century approach,” *Sociological methods & research*, vol. 50, no. 2, pp. 708–739, 2021.
- [31] J. Saldaña, *The coding manual for qualitative researchers*, 2nd ed., J. Seaman, Ed. SAGE Publishing, 2013.
- [32] M. M. Adnan, M. S. M. Rahim, A. Rehman, Z. Mehmood, T. Saba, and R. A. Naqvi, “Automatic image annotation based on deep learning models: a systematic review and future challenges,” *IEEE Access*, vol. 9, pp. 50 253–50 264, 2021.
- [33] C. Ding, M. Utiyama, and E. Sumita, “Nova: A feasible and flexible annotation system for joint tokenization and part-of-speech tagging,” *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)*, vol. 18, no. 2, pp. 1–18, 2018.
- [34] A. Berg, J. Johnander, F. Durand de Gevigney, J. Ahlberg, and M. Felsberg, “Semi-automatic annotation of objects in visual-thermal video,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [35] X. Wu, L. Xiao, Y. Sun, J. Zhang, T. Ma, and L. He, “A survey of human-in-the-loop for machine learning,” *Future Generation Computer Systems*, 2022.
- [36] J. Wang, Y. Yang, and B. Xia, “A simplified cohen’s kappa for use in binary classification data annotation tasks,” *IEEE Access*, vol. 7, pp. 164 386–164 397, 2019.

- [37] A. Shishodia, P. Verma, and V. Dixit, "Supplier evaluation for resilient project driven supply chain," *Computers & Industrial Engineering*, vol. 129, pp. 465–478, 2019.
- [38] C. Allmann, L. Winkler, and T. Kölzow, "The requirements engineering gap in the oem-supplier relationship," *Journal of Universal Knowledge Management*, vol. 1, no. 2, pp. 112–122, 2006.
- [39] J. Bach, J. Langner, S. Otten, M. Holzäpfel, and E. Sax, "Data-driven development, a complementing approach for automotive systems engineering," in *2017 IEEE International Systems Engineering Symposium (ISSE)*. IEEE, 2017, pp. 1–6.
- [40] C. Kaiser, A. Stocker, G. Viscusi, M. Fellmann, and A. Richter, "Conceptualising value creation in data-driven services: The case of vehicle data," *International Journal of Information Management*, vol. 59, p. 102335, 2021.
- [41] R. Hoda, N. Salleh, and J. Grundy, "The rise and evolution of agile software development," *IEEE software*, vol. 35, no. 5, pp. 58–63, 2018.
- [42] Y. Marton and A. Sayeed, "Thematic fit bits: Annotation quality and quantity interplay for event participant representation," in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, 2022, pp. 5188–5197.
- [43] L. Schmarje, V. Grossmann, C. Zelenka, S. Dippel, R. Kiko, M. Oszust, M. Pastell, J. Stracke, A. Valros, N. Volkmann *et al.*, "Is one annotation enough? a data-centric image classification benchmark for noisy and ambiguous label estimation," *arXiv preprint arXiv:2207.06214*, 2022.
- [44] V. Taran, Y. Gordienko, A. Rokovyi, O. Alienin, and S. Stirenko, "Impact of ground truth annotation quality on performance of semantic image segmentation of traffic conditions," in *International Conference on Computer Science, Engineering and Education Applications*. Springer, 2019, pp. 183–193.
- [45] M. Lempp and P. Siegfried, "Characterization of the automotive industry," in *Automotive Disruption and the Urban Mobility Revolution*. Springer, 2022, pp. 7–24.
- [46] S. Kochanthara, Y. Dajsuren, L. Cleophas, and M. van den Brand, "Painting the landscape of automotive software in github," in *2022 IEEE/ACM 19th International Conference on Mining Software Repositories (MSR)*. IEEE, 2022, pp. 215–226.
- [47] M. K. Slack and J. R. Draugalis Jr, "Establishing the internal and external validity of experimental studies," *American journal of health-system pharmacy*, vol. 58, no. 22, pp. 2173–2181, 2001.
- [48] A. Sorokin and D. Forsyth, "Utility data annotation with amazon mechanical turk," in *2008 IEEE computer society conference on computer vision and pattern recognition workshops*. IEEE, 2008, pp. 1–8.