

Data Annotation Quality Problems in AI-Enabled Perception System Development

Hina Saeeda*

hinasa@chalmers.se

Chalmers University of Technology and University of
Gothenburg
Sweden

Mazen Mohamad

RISE Research Institutes

Gothenburg, Sweden

mazen.mohamad@ri.se

Tommy Johansson

Kognic AB

Gothenburg, Sweden

tommy.johansson@kognic.com

Eric Knauss

Chalmers University of Technology and University of

Gothenburg

Gothenburg, Sweden

eric.knauss@cse.gu.se

Abstract

Data annotation is a critical yet error-prone activity in developing AI-enabled perception systems (AIEPS) for automated driving. Annotation quality directly affects AI model development performance, safety and reliability. Yet, there is little empirically grounded understanding of how annotation errors arise and propagate across the multi-organisational automotive supply chain. This study investigates the types, causes, and effects of data annotation errors through a multi-organisation case study involving six companies (e.g., technology and AI companies that develop software for advanced driver assistance systems and self-driving cars) and four research institutes in Europe and the UK. We conducted 19 semi-structured interviews with 20 experts (≈ 50 hours of transcripts) and applied a six-phase thematic analysis. The resulting data annotation errors taxonomy identifies 18 recurring error types across three major data-quality dimensions: *completeness* (attribute omission, missing feedback loop, privacy/compliance omission, edge-case omission, selection bias, sensor synchronisation issues), *accuracy* (wrong class label, bounding-box errors, granularity mismatch, insufficient guidance, bias-driven errors), and *consistency* (inter-annotator disagreement, ambiguous instructions, lack of purpose knowledge, misaligned hand-offs, limited review and logging, lack of frameworks and standards, cross-modality misalignment). For research rigour and triangulation, we further validated the proposed taxonomy of data annotation errors with industry. Practitioners confirmed the usefulness of the taxonomy for root-cause analysis of data annotation errors, supplier quality reviews, new project onboarding, and optimising data annotation guidelines. Practitioners described it as a “failure-mode catalogue” comparable to FMEA (Failure Mode and Effects Analysis). By framing annotation quality as an AI-enabled system development lifecycle and supply-chain concern, this work advances SE4AI by providing a shared vocabulary, diagnostic checklist, and actionable guidance for trustworthy AIEPS development.

CCS Concepts

• **Software and its engineering** → **Empirical software engineering**; • **Information systems** → *Data quality*; • **Computing methodologies** → AI-enabled perception system development problems; • **Applied computing** → Automotive Domain.

Keywords

AI-enabled Perception Systems, Autonomous Driving, Data Annotation Errors, Taxonomy, Expert-Based Evaluation, Multi-Organisation Case Study, Data Quality Assurance, Automotive Domain, Automotive Supply Chain, AI-Enabled System Development Life Cycle

1 Introduction

AI-enabled automotive perception systems (AIEPS) are central to automated driving, supporting object detection, tracking, and classification for enhanced safety and efficiency [1]. These systems underpin advanced driver assistance systems (ADAS), which enhance driving safety, cost efficiency, and convenience [16, 20]. Core perception capabilities such as *pedestrian detection*, *traffic sign recognition*, and *obstacle avoidance* serve as essential inputs to ADAS functions by interpreting sensor data from cameras, radar, LiDAR, and ultrasonics [6, 15, 33]. Within the AIEPS development lifecycle, data annotation, the process of labelling raw sensor data and converting it into structured datasets, remains one of the most expensive and error-prone tasks. Ultimately, the performance of these systems depends on the quality of annotated data used for training and validation [15, 23].

Motivation. While annotation quality is widely recognised as critical for AIEPS, the systematic impact of data annotation errors on system performance remains insufficiently understood. In safety-critical domains such as autonomous driving, even minor annotation errors such as missed or inconsistently labelled objects can cascade into unsafe behaviours and degraded model reliability [7, 12, 40, 41]. However, how such errors originate, propagate across the AIEPS development lifecycle, and ultimately impact system performance remains underexplored. [10, 22, 31]. This gap is amplified in multi-organisational automotive supply chains, where heterogeneous tools, inconsistent annotation practices, and fragmented quality assurance processes increase error propagation. Addressing this issue is essential to shift from reactive correction toward proactive, standardised, and trustworthy annotation processes [11].

Research Gap. Prior work has focused on isolated error types or technical solutions such as weak supervision and automated

labelling [25, 26, 32]. Crucially, there are no comprehensive, empirically grounded studies that systematically capture the full spectrum of annotation errors across multi-organisational automotive supply chain contexts. In this domain, different tiers produce annotated datasets, develop AIePS, test them, and deliver integrated products to original equipment manufacturers (OEMs). The lack of an industry validated, comprehensive knowledge base (e.g. a detailed taxonomy) of annotation errors hinders alignment of processes, knowledge sharing, and quality assurance across organisational boundaries.

Study Aim. Annotation quality is not merely a data issue but an AIePS development lifecycle concern, as errors propagate from annotated data through the AI model to final system decisions. Accordingly, this study aims to provide insights and evaluations of the data annotation errors to support proactive, lifecycle-oriented quality assurance. This study provides a taxonomy of data annotation errors and a classification of their root causes and impacts. This study also provides insights into how companies can use this information to their advantage after validating the usefulness of the proposed taxonomy with industrial practitioners.

Research Questions. Guided by these objectives, this multiorganisational case study investigates:

- **RQ1:** What are the different types of data annotation errors, their causes, and how do they affect AIePS development and performance?
- **RQ2:** How do practitioners perceive the usefulness of the proposed taxonomy of data annotation errors?

Methodology Overview. We conducted 19 semi-structured interviews with 20 participants from six international companies and four research institutes. Using multi-level triangulation, integrating case study data, thematic analysis, and expert validation, the study ensures both empirical depth and practical relevance.

Novelty. Our study is novel in (a) it presents the first systematic multi-organisational empirical study in the automotive domain covering the full supply chain, (b) it introduces an industry validated taxonomy of data annotation errors with practical relevance, and (c) it bridges data quality assurance (QA) with SE4AI practice by framing annotation errors as AIePS development lifecycle bottlenecks rather than isolated data issues.

2 Background and Related Work

The AIePS development Lifecycle and Role of Data Annotation. The development of AI-enabled systems involves multiple stages, data collection, annotation, model training, evaluation, deployment, and monitoring [9, 23]. Among these, data annotation—the transformation of raw sensor data into labelled datasets—is a crucial and resource-intensive phase [11]. By labelling images, point clouds, or sensor sequences with bounding boxes, segmentation masks, or object classes, annotators create the ground truth essential for model learning, where annotation quality directly impacts performance, generalisation, and decision accuracy in safety-critical domains [12, 15].

Data Annotation Error. The annotation errors are incorrect, incomplete, or inconsistent labels in training or validation

data [17, 19], arising from human variability, ambiguous guidelines, limited tools, or weak quality assurance [13, 22]. Examples include misclassifying pedestrians on scooters, inconsistently labelling occlusions, or omitting small objects, all of which introduce noise that weakens model learning [31, 40]. These errors can propagate throughout the AI lifecycle, leading to false detections, sensor misinterpretations, and reduced safety in real-world driving [7, 41].

Data Quality Dimensions. Research identifies three core dimensions of data quality: completeness, accuracy, and consistency as key indicators of reliable data representation [2, 3, 27, 37]. Originally established in information systems [27, 37], this framework now underpins modern AI models and annotation quality assessment. Our taxonomy builds upon these classical dimensions. Completeness reflects whether all required data are present and relevant real-world entities or attributes are captured [2, 3]. In annotation, missing objects or scenarios (e.g., pedestrians in poor lighting) represent completeness errors that reduce dataset representativeness. Accuracy measures how closely annotations (e.g., bounding boxes, segmentation masks) match true real-world conditions [2, 27]; misaligned boxes or incorrect labels exemplify accuracy errors. Consistency concerns uniformity across annotators, scenes, or systems [2, 3]. Inconsistent labelling of similar objects or frames constitutes a consistency error that undermines reproducibility, fairness, and reliability.

Annotation Quality Across the AI Lifecycle in Industrial Supply Chains. AIePS development spans a complex supply chain of OEMs, Tier-1, and Tier-2 suppliers [11, 15]. Each uses distinct annotation tools, standards, and quality controls, resulting in heterogeneous datasets where early-stage errors can propagate through model training, validation, and integration [22]. Consequently, annotation quality represents a dynamic, lifecycle-spanning assurance challenge that directly influences the dependability, maintainability, and trustworthiness of AI-enabled systems [10, 31].

3 Research Methodology

We followed the qualitative case study guidelines established by Runeson and Höst et al. [30], and aligned our study design and reporting with the empirical standards proposed by Ralph et al. [28] to ensure transparency, triangulation, and validity (see Figure 1).

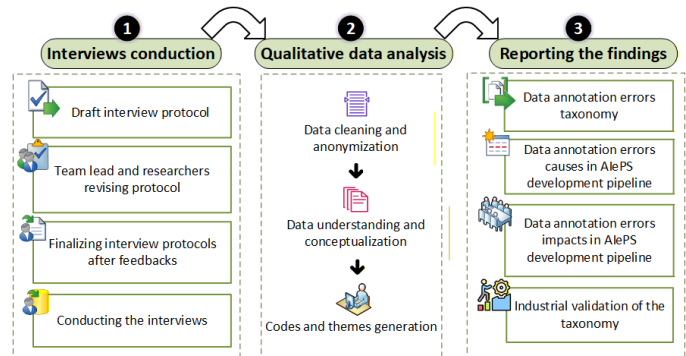


Figure 1: Research method followed.

Research Context This study is part of Project X, which aims to enhance academia-industry collaboration by developing concepts, models, and techniques for effective and safe AlePS productions. We conducted 19 semi-structured interviews [30] with 20 participants from industry and academia. These interviews were performed online between October 2024 and April 2025, each lasting ≥ 120 minutes. ID16 participated in two separate interviews (labelled 16A and 16B), each lasting approximately four hours. In contrast, participants ID18 and ID19 took part in a single focus group session, counted as one interview (see Table 1). In total, this amounts to over **50 hours of rich transcripts**. Using purposeful sampling [28], we selected participants based on their roles, organisational context, and experience with data annotation or AI model development in autonomous driving. Interviews were recorded, transcribed, manually corrected, and concluded with participant feedback. This approach ensured diverse and in-depth insights into annotation quality, associated errors, and their impact on the AlePS development life cycle.

Companies Context To capture diverse perspectives, we conducted fourteen interviews with representatives from six companies engaged in AlePS development. We included the entire supply chain: one OEM, three Tier-1 suppliers, two Tier-2 suppliers, a perception training company, and a code analysis firm, all based in two European countries and the UK.

Additionally, five interviewees from four research institutes, including a state-owned European institute and leading universities from Europe and the UK, helped bridge theory and practice across AI, data annotation, quality, and regulation. Thematic Saturation was reached with Interview 15, with no new information or candidate themes emerging in Interviews 16, 17, 18, and 19.

Participants As shown in Table 1, targeted experts included leading contributors to AlePS, such as safety standard developers, founders of relevant companies, and the chief investigator of top-tier research labs, ensuring deep domain expertise. This is a niche domain where only a few globally recognised organisations and experts are active. We deliberately reached out to these leading actors and succeeded in including the most relevant and established experts in this emerging field.

Interview Preparation: A semi-structured interview guide [28] was developed iteratively by four researchers in collaboration with three industrial experts experienced in data annotation and AlePS, thereby establishing content validity. The guide featured open-ended, neutral, and clear questions aligned with the research questions and relevant literature presented in the background and related work section 2. To enhance clarity, structure, and duration, the guide was piloted with two independent industry partners (not part of the final study), representing annotation requirements (Tier 1) and production (Tier 2), whose feedback informed refinements. (For details, see: [→ Interview Guide](#))

Data Analysis: We applied a six-phase thematic analysis process following Braun and Clarke *et al.* [5]. The analysis was guided by a combined deductive-inductive codebook, and the units of analysis were 19 interview transcripts, each coded in text chunks of approximately three sentences. Across all transcripts, about 1050 text segments (coded units) were analysed.

(1) Familiarisation: Two researchers transcribed, anonymised, and validated all interviews to ensure ethical integrity and data

Table 1: Categorisation of interviewees based on company type and specialisation, including their roles and years of experience

Company ID: Type	Focus Area	Expert (Years of Experience)
A: Tier 2	Data Annotation	ID1: Annotation Expert (6) ID2: Perception Expert (10) ID3: Quality Expert (9)
B: Tier 1	Safety Software	ID4: Machine Learning Expert (7) ID5: Data Scientist (11)
C: University	Research	ID6: Senior Researcher (9)
D: University	Research	ID7: Researcher (8)
E: University	Research	ID8: Researcher (5)
F: Research Institute	Research	ID9: Researcher (10) ID10: Researcher (20)
G: OEM	Automotive	ID11: Machine Learning Expert (5) ID12: V&V Expert (5) ID13: Data Engineer (4) ID14: Researcher (3)
H: Tier 1	Safety Software	ID15: Researcher (10)
I: Tier 2 I: Tier 2	Quality Assurance Quality Assurance	ID16 A: Quality Expert (18) ID16 B: Quality Expert (18)
J: Tier 1	Digital Solution	ID17: Head of Research (17) ID18,19: Research Engineer (1,2)

accuracy. Through multiple readings, they familiarised themselves with the data, developing an initial understanding of emerging patterns and potential themes.

(2) Initial Coding: A deductive codebook was first developed based on the research questions and interview guide (see: [→ Codebook](#) and [→ Codebook-Aligned Initial Codes](#)). During open coding, inductive codes were added as new insights emerged. Each transcript (ID1–ID19) was coded systematically, resulting in a comprehensive set of detailed codes extracted from all interviews (see: [→ Complete List of Emerging Codes from Transcripts](#)).

(3) Code Validation: Two researchers independently coded all transcripts and cross-verified results, achieving a strong inter-coder reliability (Cohen’s $\kappa = 0.8$) corresponding to 80% agreement. The remaining minor differences (about 20%) were resolved through discussion and consensus, ensuring the consistency and credibility of the coding process.

(4) Theme Development: All codes were clustered into potential themes, related to the three key data quality dimensions *Completeness*, *Consistency*, and *Accuracy* and aligned with the pre-defined research questions (see: [→ Detailed Themes Development](#)). A consolidated analysis produced an initial *Master List of Data Annotation Errors*, capturing all identified errors across the AlePS development lifecycle. The master list compiles identified errors organised by their nature, frequency, and impact, with representative examples forming the basis for the final taxonomy of annotation errors.

(5) Theme Refinement: The five researchers collaboratively reviewed, merged, and refined overlapping codes through a three-stage iterative process to ensure analytical robustness. In the first round, all codes within each theme were compiled, and their frequencies (*i.e.*, the number of mentions across interviews) were recorded. Across the three quality dimensions, we identified 23 codes under *Completeness*, 23 under *Consistency*, and 18 under *Accuracy*. In the second round, codes mentioned three times or fewer

were classified as low-frequency. In the final round, these were removed, retaining only the most recurrent data annotation errors consistently mentioned by multiple experts. This resulted in six high-frequency codes for *Completeness*, seven for *Consistency*, and five for *Accuracy* (see: → [Final Refine Data Annotation Errors List](#)).

(6) Reporting: The final themes and associated codes were verified, documented, and linked with representative quotes to maintain transparency and traceability. Five researchers collaboratively reviewed the complete analysis using Excel, OneDrive, and Miro, ensuring structured documentation and analytic rigour.

Replication Package: We make the interview guide, codebook, thematic analysis themes, and validation survey available to support future researchers in replicating the study. The replication package is available on HARVARD Dataverse ¹.

4 Findings

The Findings section presents a taxonomy of data annotation errors, highlighting their causes and impacts (Figure 2). It also reports expert validation results confirming the taxonomy’s usefulness, industrial relevance, and practical applicability (Section 4.2).

4.1 RQ1: Data Annotation Errors, their Types, Causes and Impacts in AlePS Development

(1) Completeness Errors. This section addresses completeness errors related to both data coverage and the data annotation processes that sustain and verify that coverage over time.

(1.1) Attribute Omission: Missing or incomplete attributes (e.g., colour, state, occlusion, or diversity attributes) leave datasets insufficiently descriptive. As one participant noted, “Another issue is adding proper metadata [attribute]... otherwise, older datasets become outdated or inconsistent.” (ID10) Causes include unclear requirements, annotator fatigue, and lack of validation or schema linkages (ID5–ID10). Such omissions degrade reliability and performance in edge case scenarios.

(1.2) Missing Feedback Loop: This error refers to the absence of a systematic review and refinement process that allows recurring annotation errors to persist across iterations. As highlighted, “The most striking aspect missing from this process flow is the feedback loop... teams cannot adjust their approach during annotation, and mistakes get repeated.” (ID3) Weak communication, rare quality reviews, and tools lacking traceability hinder continuous improvement (ID9–ID18), slowing iterations and degrading perception quality. The “Missing feedback loop” is categorised under completeness because it represents a process through which completeness fails to be achieved or maintained, rather than an isolated data issue.

(1.3) Privacy and Compliance Omission: This error refers to the insufficient enforcement or documentation of privacy and compliance measures during data annotation. Inadequate anonymisation or weak enforcement of GDPR-like policies introduces ethical and legal risks. Participants noted, “Ensuring data privacy... is critical when dealing with sensitive data,” (ID15) and “GDPR requirements remain somewhat unclear and there is an inherent risk in misinterpreting these guidelines” (ID9). Missing anonymisation tools or audit

logs (ID13–ID18) reduce visual uniformity and trustworthiness. Although “Privacy and Compliance Omission” primarily relates to ethical and legal governance, it is included under the completeness theme as it directly influences dataset coverage. Weak enforcement or unclear interpretation of privacy requirements often leads to the omission of sensitive samples or attributes, thereby reducing the representativeness and completeness of annotated data.

(1.4) Edge cases /Unforeseen Scenarios Omission: This error refers to the omission of rare, unexpected, or safety-critical cases that fall outside predefined annotation guidelines. Rare or unforeseen scenarios are often ignored. One participant explained, “At the same time, we should acknowledge the existence of unknown scenarios, particularly edge cases. The challenge with unknown scenarios is that we cannot write explicit guidelines for them because they are unknown.” (ID12) Another added, “Ambiguities arose—whether a pedestrian on a scooter should be annotated as relevant or not.” (ID6) Missing mechanisms to flag unusual cases (ID8, ID13) limit robustness in safety-critical settings.

(1.5) Selection Bias: Overrepresentation of common conditions (e.g., daylight, clear weather) creates unbalanced datasets. “Some projects focus on specific events... this inherently introduces bias.” (ID5) and “Most datasets capture sunny conditions... models struggle in rain or at night.” (ID18) Lack of planning and diversity requirements (ID14–ID17) reduces model adaptability and fairness.

(1.6) Synchronisation \ Calibration Issues: This error refers to the misalignment or incomplete calibration of multi-sensor systems such as cameras, radar, and LiDAR used during data collection and annotation. Misaligned or unsynchronised sensors produce annotation drift. As noted, “You need to synchronise sensors so that data points align in time and space... otherwise, they show up in different places.” (ID10) Missing timestamp checks and calibration metadata (ID16–ID18) reduce precision in multi-sensor fusion.

(2) Accuracy Errors. Under this theme, we address both procedural and data related factors contributing to accuracy errors. We are adopting a holistic life cycle perspective rather than a purely data-centric view.

(2.1) Wrong Class Label: This error refers to the incorrect assignment of semantic categories to objects or scenes, where visually or contextually similar instances are labelled under the wrong class. Misclassifications occur when annotators confuse similar categories. “In autonomous driving, you have different classes of vehicles...OK, this is a car, this is a truck, this is a van ... different people might annotate the same object differently.” (ID18) Ambiguous class boundaries, annotator fatigue, and poor validation (ID2–ID7) lead to noisy supervision and degraded model reliability.

(2.2) Bounding Box Errors: This error refers to inaccurate, inconsistent, or corrupted bounding boxes used to localise objects in images or sequences. Bounding-box errors arose from weak quality assurance tolerances, ambiguous guidelines for occlusions, subjective boundary interpretations, annotator fatigue, and tool instability, causing inconsistent rendering or corrupted coordinate data. One participant shared, “Some clients insist on pixel-perfect accuracy like a 10-pixel margin.” (ID2) Another added, “When annotating objects... avoid optimistic annotations; underestimating size causes issues in collision avoidance.” (ID12) Missing QA tolerances and unstable rendering tools (ID8–ID14) amplify spatial drift and tracking instability.

¹The replication data associated with this study is publicly accessible and available online here.



Figure 2: Taxonomy of data annotation errors in AI-enabled perception systems, organised into three main categories: Completeness, Accuracy, and Consistency.

(2.3) Granularity Mismatch: This error refers to inaccuracies in the level of annotation detail, either overly fine-grained or overly coarse, relative to the intended schema or perception task. Misalignment between annotation detail and schema depth leads to inconsistency. *“If the requirement specifies capturing the entire head but annotators only mark facial features, this introduces bias.”*

(ID8) and *“We used a taxonomy... it didn’t include trams... the taxonomy was from the US.”* (ID17). Outdated taxonomies (ID9, ID18) reduce dataset interoperability and model transferability. Inconsistent levels of detail, such as incomplete object coverage or missing 3D dimensions (ID8, ID18), impair model calibration and feature

learning, ultimately degrading recognition accuracy and the ability of AlePS to adapt across diverse operational domains.

(2.4) Insufficient Guidance: This error refers to vague, incomplete, or ambiguous annotation guidelines that fail to provide clear operational instructions for annotators. Unclear and incomplete annotation guidelines lead to subjective interpretation. *“If annotators are given vague guidelines... they may not annotate consistently.”* (ID2) and *“Initially, it wasn’t clear how to annotate in cases where visibility was poor.”* (ID5) Insufficient guidance resulted from incomplete documentation of quality criteria, ambiguous or evolving rules, limited annotator expertise and feedback, and tool limitations such as missing guideline prompts, version control, and traceability to design goals. Ambiguous or incomplete annotation guidelines increase inter-annotator variability, introduce systematic noise, and weaken label consistency, ultimately reducing model accuracy, reliability, and trust in AlePS (ID5, ID6, ID7, ID16). Although “Insufficient Guidance” stems from process or documentation shortcomings, it is classified as an accuracy error because it directly leads to incorrect or suboptimal labelling. Ambiguous or evolving guidelines cause misinterpretations, class confusion, and boundary errors, systemic sources of inaccuracy rather than mere procedural flaws. While “Insufficient guidance” can also lead to inter-annotator variability (e.g., consistency errors), it is primarily categorised under accuracy because the root cause lies in inadequate or ambiguous annotation instructions that lead annotators to produce incorrect or imprecise labels relative to the intended ground truth.

(2.5) Bias-Driven Errors: This error refers to systematic deviations in annotation outcomes caused by human, contextual, or automation-induced biases that distort the accurate representation of real-world phenomena. Systematic deviations arise from cognitive or automation-induced bias. *“Subsets systematically mislabeled (e.g., red cars)... supervisors might have interpreted differently.”* (ID2) and *“Automation tools rubber-stamped labels... annotators over-trusted tool outputs.”* (ID6) Bias-driven errors arose from unbalanced datasets, misaligned taxonomies, and overreliance on automated labelling, compounded by human cognitive biases, automation overtrust, and inconsistent quality assurance, which reinforced systemic data skew and unfairness (ID9, ID11, ID15). Overreliance on automated pre-labelling and unbalanced samples (ID8–ID15) propagates bias into perception models, harming fairness and trustworthiness.

(3) Consistency Errors. Under this theme, we cover both procedural and data oriented consistency errors. While accuracy concerns whether a label is correct with respect to a known ground truth, consistency concerns whether multiple annotators or tools would label the same instance in the same way.

(3.1) Inter-Annotator Disagreement: This error refers to the variation in labels assigned by different annotators to the same data instance, caused by subjective judgments, ambiguous guidelines, or inconsistent interpretation, resulting in annotation noise and reduced dataset reliability. As one participant emphasised, *“Yes, consensus is essential. If multiple annotators work on the same data, their outputs should align to ensure accuracy.”* (ID1) Another added, *“Different annotators ... they will subjectively annotate the object depending on distance or shape ... then we need a second person to correct this.”* (ID18) Human judgment variability, vague definitions of relevance or environmental context (ID3, ID6–ID7), and

cost-driven limits on extensive review cycles (ID9, ID15) increase annotation noise and reduce dataset reliability, requiring additional reconciliation efforts before model training.

(3.2) Ambiguous Instructions: This error refers to unclear, complex, or contradictory annotation guidelines that cause annotators to interpret labelling rules differently, leading to inconsistent annotations and reduced data quality. One expert explained, *“Clients often provide guidelines, but we frequently revise them to improve efficiency. I advocate for simplification annotators should not have to read 140 pages to understand the guidelines.”* (ID2) Similarly, another observed, *“Determining which pedestrians are ‘relevant’ was unclear; guidelines did not define it strictly, leading to divergent [inconsistent] interpretations.”* (ID6) Ambiguous instructions, such as overcomplicated or overly descriptive guidelines (ID1–ID3, ID5), missing contextual criteria (e.g., occlusions, lighting, traffic density) (ID7–ID9), lack of communication channels for clarifications (ID13), and unrealistic attempts to eliminate all ambiguity without iterative updates (ID16) amplify disagreement and rework, as annotators interpret rules differently under time pressure or limited supervision.

(3.3) Lack of Purpose Knowledge: This error refers to annotators’ limited understanding of the dataset’s intended use or downstream model goals, leading to misaligned labelling decisions and inconsistent relevance judgments. Annotators often lack awareness of the dataset’s end-use, leading to divergent judgments and inconsistent relevance assessments. As one participant stated, *“The annotators sometimes don’t understand why they are annotating certain things or how it ties into the overall purpose. They often say, ‘If we knew what this was supposed to train for, we could do it better.’”* (ID4) Another noted, *“Requirements pass through multiple stakeholders, and annotators often don’t know the true purpose of labelling. This loss of context leads to inconsistent quality.”* (ID6) Poor communication of downstream objectives (ID8, ID11), insufficient domain expertise (ID12), and limited onboarding for new annotators result in inconsistent boundary choices, missing edge-case handling, and reduced overall alignment with model requirements.

(3.4) Misaligned Hand-offs: This error refers to inconsistencies and information loss that occur when annotation tasks, requirements, or guidelines are transferred between teams or stakeholders without proper documentation or communication, resulting in fragmented, inconsistent datasets. Inconsistencies often emerge during transitions between client, QA, and annotation teams. *“OEMs define requirements, but these may not be directly shared with suppliers. By the time they reach annotators, the instructions are incomplete or altered, creating inconsistencies.”* (ID6) Another participant added, *“If a team does annotation internally instead of outsourcing, they might not document it properly. Instead, they’ll convey details in meetings, and that information can be very temporary, easily forgotten.”* (ID11) Under-documented communication, missing shared tracking tools, and inconsistent rule interpretation across multiple vendors (ID4–ID6, ID8, ID11, ID14) lead to fragmented datasets and require post-hoc harmonisation to ensure uniformity before model deployment.

(3.5) Limited Review & Logging: This error refers to the absence of consistent quality checks and detailed record keeping during the annotation process, which prevents error detection, traceability, and reproducibility of annotation decisions. Scarce QA

resources and weak process traceability leave inconsistencies undetected. One participant explained, “Given the vast data, manual review is not economically viable. With too few QA resources, errors remain undetected.” (ID6) Another added, “Annotation companies rely on manual reviewing ... they said they want 95% accuracy, but it was all done manually ... automated QA was lacking.” (ID18) Insufficient automated validation (ID10, ID13), limited sampling or inter-annotator agreement checks (ID14–ID17), and high costs of independent verification (ID15–ID16) result, corrupted or inconsistent labels (e.g., bounding-box drift or missing objects). These errors propagate into training data, lowering model stability and generalisation.

(3.6) Lack of Frameworks and Standards: This error refers to the absence of unified annotation protocols, ontologies, or benchmarking criteria across projects or organisations, leading to inconsistent practices, poor comparability, and reduced reproducibility of annotated datasets. The absence of shared annotation standards or tool benchmarks leads to inconsistent practices across organisations. One participant summarised, “One of the biggest challenges is the lack of standardisation in the annotation market... It is difficult to compare different annotation providers because there is no universal benchmark.” (ID3) Similarly, another noted, “There is no common standard ... every company follows their own guideline ... sometimes datasets lack information about how annotations were done.” (ID18) Missing international standards (ID6, ID9, ID13), inconsistent adherence to ISO/UNECE guidelines (ID10, ID15), and incomplete documentation in datasets limit comparability, benchmarking, and reproducibility across AIePS.

(3.7) Cross-Modality Misalignment: This error refers to spatial or temporal inconsistencies between annotations across different sensor modalities (e.g., camera, LiDAR, radar), caused by calibration errors or synchronisation gaps, leading to inaccurate data fusion and degraded perception performance. Misalignment between sensor modalities (e.g., camera–LiDAR–radar) results in inconsistent spatial annotations. As one participant highlighted, “When multiple sensors record the same scene, they need to be registered together... You need calibration to make sure they align.” (ID10) Another added, “We projected annotations from point cloud onto the image plane and checked alignment... sometimes the projection didn’t capture the object correctly.” (ID18) Temporal sampling differences (ID13, ID17), inadequate calibration protocols (ID9–ID10), and missing 3D-to-2D projection validation distort object positioning and tracking accuracy across modalities, impairing fusion-based perception and sensor reliability.

4.2 RQ2: Experts’ Validation of the Data Annotation Errors Taxonomy

Validation Purpose and Scope. The validation focused on confirming the accuracy, relevance, and industrial applicability of the identified error categories rather than generating new data. That’s the reason the validation involved four experts (1 OEM, 1 Tier-1, 2 Tier-2); these participants were senior practitioners directly involved in data annotation governance and quality assurance, selected for their domain expertise and cross-organisational perspective. While all experts were based in European automotive contexts, their organisations operate globally, ensuring exposure

to international standards and practices. This confirmatory and triangulated approach, grounded in design science and empirical software engineering principles [14, 30], ensured the taxonomy’s practical relevance and completeness within real-world annotation workflows. The aim of this validation was depth and expert triangulation rather than statistical generalisation.

Validation Approach. Experts reviewed the data annotation error taxonomy using a structured questionnaire covering four dimensions: (1) their background and experience in annotation workflows, (2) the taxonomy’s usefulness and relevance, (3) its clarity, categorisation, and improvement areas, and (4) its practical applicability within supply chain processes. This systematic design ensured a balanced evaluation of both conceptual soundness and industrial usability (see → [Structured Validation Questions](#)). The (Table 2) shows details of the complete validation responses.

(1) Expert Profiles: Starting Questions (Q1–Q4). The validation involved experts, including a System Safety Expert (OEM), Machine Learning Engineer (Tier-1), Director of Customer Services (Tier-2), and Perception Expert (Tier-2). Each possessed 4–5 years of experience in annotation workflows and AIePS development, with direct roles in data collection, annotation, and AI system integration.

(2) Usefulness, Relevance, and Applicability (Q5–Q16). Across all experts, the taxonomy was perceived as highly relevant, practical, and directly applicable within industrial settings for structuring, diagnosing, and improving data annotation workflows.

- The **OEM system safety expert** viewed the taxonomy as a “failure-mode catalogue” analogous to FMEA (Failure Mode and Effects Analysis), helping identify and trace data-related weaknesses in perception-system development. It was considered particularly valuable for root-cause analysis, supplier quality reviews, and training of perception validation teams. The expert also noted that it supports the diagnosis of perception failures (e.g., inconsistent object recognition or missed detections) and facilitates cross-team communication among the safety, perception, and data quality units.

- The **Tier-1 ML engineer** described the taxonomy as a practical “checklist” for setting up new annotation pipelines, reviewing guidelines, and validating existing annotation processes. It was also seen as a “cheat sheet” for systematically verifying completeness, identifying recurring issues, and ensuring that all major error sources are considered during the design or revision of data annotation workflows.

- The **Tier-2 director of customer services and perception expert** emphasised that the taxonomy supports “systematic analytics of annotation errors” and improves “onboarding for new projects.” They found it particularly useful for structuring internal quality assurance (QA), facilitating client discussions, and aligning understanding between annotation providers and customers. Perception expert further highlighted its potential as onboarding material for new employees and clients, and as a reference during guideline optimization and annotation tool configuration.

(3) Reflection on Clarity, and Improvements (Q9–Q12). All experts agreed that the three overarching quality dimensions, completeness, accuracy, and consistency, effectively capture the main data quality challenges observed in industrial contexts. The OEM noted that these dimensions align with established quality

Table 2: Expert validation outcomes, industrial utility, and implications for the taxonomy of data annotation errors

Theme	Validation Outcomes	Representative Comments (OEM / Tier-1 / Tier-2)	Industrial Utility	Implications
Usefulness	High relevance for structuring and diagnosing annotation failures; recognised as a cross-industry quality reference.	"OEM: Can be used as a failure-mode catalogue." "Tier-1: Helpful when setting up new annotation pipelines." "Tier-2: Useful as a checklist and for systematic analytics."	Strengthens QA traceability. Enables structured failure diagnosis and cross-team quality review.	Validated taxonomy's generalisability; retained as a framework for audits and supplier quality reviews.
Practical Experience	The taxonomy is clear and comprehensive, suggesting the merging of overlapping categories.	"Tier-2: Human vs automation not enough—add tool and guideline causes."	Improves clarity of defect reporting. Supports consistent error documentation across annotation teams.	Expanded causes to include "tool-related" and "guideline-related" error types.
Relevance	Consensus on three dimensions: completeness, accuracy, and consistency, aligned with ISO-based quality standards.	"OEM: These are the same dimensions we face daily."	Ensures interpretability in audits. Aligns with existing industry QA frameworks.	Maintained three main categories; clarified overlaps through refined definitions.
Suggestions	Encouraged merging sub-items, simplifying structure, and adding quantitative risk indicators.	"Tier-2: Add a risk-assessment matrix linking severity and probability."	Enables risk-based prioritisation. Supports safety-critical assessment of annotation errors.	Planned extension of taxonomy with severity-frequency mapping (akin to FMEA scoring).
Use Case Scenarios	Realistic applications include training, QA reviews, and root-cause analysis.	"OEM: Use for checking perception-system root causes." "Tier-1: Cheat sheet for reviewing annotation processes."	Enhances training and collaboration. Standardises terminology across suppliers and OEMs.	Evidence of immediate applicability; supports integration into QA tools and training programs.
Errors and Integration	Systematic mislabelling, vague instructions, and missing edge cases are considered most damaging.	"Tier-2: Systematic errors lead to ML models learning incorrect patterns."	Supports model reliability and safety. Enables early identification of critical annotation faults.	Prioritised these error types in the final taxonomy; marked as high-severity recurrent issues.

and safety standards such as ISO 21448, while the Tier-1 expert emphasised simplifying and consolidating categories to facilitate adoption, a recommendation incorporated into the final taxonomy refinement. Tier-2 experts proposed distinguishing more clearly between human-, automation, and tool-induced errors, which informed the cause impact grouping in our analysis. Collectively, experts encouraged extending the taxonomy through studies of real-world failure cases and linking each error type with measurable severity indicators such as frequency or safety impact.

(4) Preliminary Risk Categorization of Errors. The experts consistently emphasised high-risk errors such as systematic mislabelling, incomplete or edge-case omissions, inconsistent annotations, and schema mismatches, all of which align with the taxonomy's 18 core errors. In contrast, several low-frequency but practically significant errors also emerged, including cross-modality misalignment, low-quality pre-annotations, limited feedback loops, and human fatigue. While less frequently mentioned, these emergent cases reveal evolving challenges arising from automation, organisational constraints, and human variability, indicating where future quality assurance efforts and process improvements should focus. (see: → [Table 3- Expert Classification of Annotation Errors by Risk Level.](#))

5 Threats to Validity

Following Runeson *et al.* [30] and Maxwell [21], potential validity threats and mitigation measures are outlined below.

Construct Validity. To ensure participant relevance, experts were purposefully selected from OEMs, Tier-1, Tier-2, and research institutes. The interview protocol underwent internal review and two pilot tests for clarity and completeness. Key terms (*e.g.*, data annotation, perception systems) were clearly defined and supported with visuals to avoid misinterpretation. For expert validation, participants received structured briefings and could seek clarification

before evaluation. Researcher bias was minimised through team-based reflexive discussions and peer debriefing.

Internal Validity. To mitigate personal or organisational bias, participants were drawn from six automotive companies and four research institutes across multiple supply chain tiers. Interviews were recorded, transcribed, and participant-verified for accuracy. Dual independent coding yielded $\kappa = 0.8$ (Cohen's Kappa), with discrepancies resolved during weekly meetings from May to September 2025. Three collaborative workshops with researchers and industry partners refined the themes, while triangulation between interview (RQ1) and validation (RQ2) results enhanced conceptual and practical credibility.

External Validity. Although not statistically generalisable, the study achieves analytical generalisation through diverse participants, 20 experts from 11 organisations across two European countries and the UK. Coverage across OEMs, Tier-1/Tier-2 suppliers, and research institutions, combined with alignment to international standards (ISO/IEC 5259, IEEE P2801, ISO 26262, SAE J3016), strengthens transferability and ensures lifecycle representativeness within AlePS.

Reliability. To ensure transparency and replicability, all instruments (interview guide, codebook, validation questionnaire) are openly available. Interviews were conducted with informed consent, anonymised in accordance with GDPR, and securely stored. All research steps, from data collection to coding logic, were systematically documented to support replication and independent verification.

6 Conclusions and Discussion

This study presents a structured taxonomy of 18 data annotation errors in AI-enabled automotive perception, integrating theoretical

data-quality concepts with industrial insights. Grounded in the dimensions of completeness, accuracy, and consistency.

Data Annotation Errors and the State of the Art. Completeness errors capture the underrepresentation of critical driving scenarios and incomplete information capture in perception datasets [1, 38]. Traditional issues such as *edge-case omissions*, *selection bias*, and *attribute omission* align with known “long-tail” challenges [3, 15, 17]. Yet, our taxonomy extends this view by introducing novel dimensions such as *privacy/compliance omissions*, *missing feedback loops*, and *sensor synchronisation issues* [7, 29, 36, 41]. These additions highlight that completeness is not merely a matter of dataset size but of regulatory, temporal, and multimodal adequacy across sensors and stakeholders. The Figure 3 highlights the most frequently mentioned completeness-related errors in the study. Missing feedback loop and selection bias were the most common issues, followed by edge/unforeseen scenarios. Less frequent but notable errors included attribute omission, privacy/compliance omissions, and calibration/synchronization issues. Overall, the findings suggest that systemic gaps, particularly in feedback and data selection, are the main contributors to completeness errors.

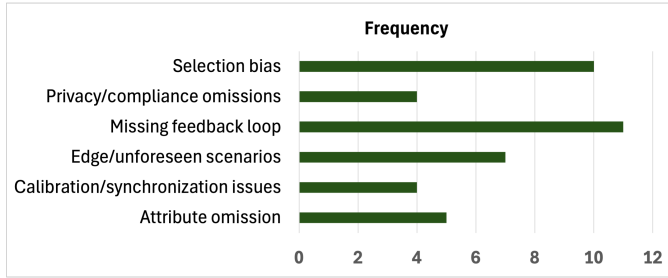


Figure 3: High mentioned completeness errors.

Accuracy errors represent distortions between the intended ground truth and the final annotation. Common cases such as *wrong class labels*, *bounding-box errors*, and *granularity mismatches* remain frequent causes of label noise [18, 25, 32]. At the same time, industrial evidence revealed human automation interaction issues *automation over trust*, *insufficient guidance*, and *lack of calibration feedback* that amplify inaccuracies through over-reliance on automated pre-labelling tools [3, 29, 39]. These findings bridge technical and socio-organisational sources of label errors, showing that accuracy degradation often stems from both data limitations and human decision making biases. Figure 4 presents the most frequently mentioned accuracy related errors in the study. Insufficient guidance was the most commonly reported issue, followed by bounding box errors and wrong class labels, each occurring frequently. Bias driven errors were also notable, while granularity mismatch appeared less often. Overall, the findings suggest that accuracy problems stem largely from unclear annotation instructions and technical labelling inconsistencies, underscoring the need for clearer guidelines and improved annotation precision.

Consistency errors denote deviations in labelling uniformity across annotators, projects, and modalities [3, 17, 35]. While inconsistent instructions and inadequate calibration remain established concerns, our taxonomy identifies process-level consistency failures

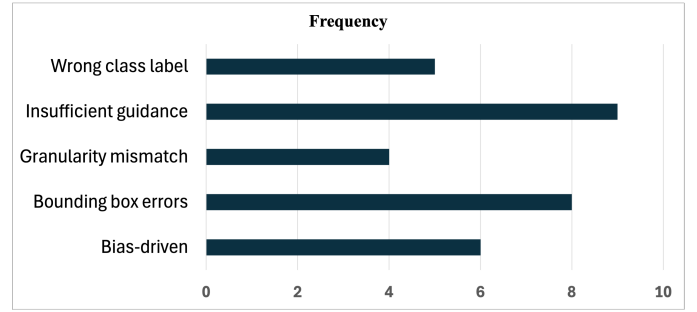


Figure 4: High mentioned accuracy errors.

unique to distributed industrial pipelines such as *misaligned hand-offs*, *cross-modality misalignment*, and lack of *standardisation frameworks* [4, 36, 41]. Many of these issues stem from socio-technical misalignments, including ambiguous guidelines, workload variability, and lack of motivational factors and annotator fatigue [3, 24]. Addressing such issues thus requires not only better annotation tools but also systematic processes such as feedback loops, clearer organisational responsibilities, and continuous cross-team calibration. As shown in the Figure 5, ambiguous instructions were the most frequently reported consistency-related error in our study, followed by inter-annotator disagreement and lack of frameworks or standards. Other notable issues, such as misaligned hand-offs and limited review and logging, point to process level errors, while lack of purpose knowledge and cross-modality alignment appeared less often but remain relevant. We observed that annotation inconsistency arises from both human variability and systemic gaps in process standardisation and guidance.

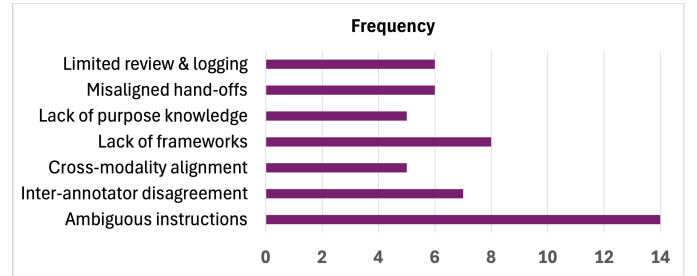


Figure 5: High mentioned consistency errors.

Data Annotation Errors and Quality Degradation. The study demonstrates that annotation errors are not random or isolated incidents but systemic drivers of downstream quality degradation in perception models. Completeness errors, such as missing feedback loops or edge case omissions, reduce dataset representativeness and weaken model generalisation [1, 38]. Accuracy-related errors, including mislabeling, bounding-box misplacement, and granularity mismatches, distort learning signals and reduce predictive reliability [3, 17, 25]. Likewise, consistency issues from inter-annotator disagreement to cross-modality misalignment introduce latent noise that propagates through training and deployment [15, 24]. Together, these defects can cascade throughout the AI

lifecycle, undermining performance in safety-critical metrics such as false positives, missed detections, and perception uncertainty. Grounded in the established principles of *completeness*, *accuracy*, and *consistency* [2, 27, 37], the taxonomy operationalises abstract data-quality constructs into concrete, auditable categories. This provides a data-centric assurance foundation aligned with the broader trustworthy-AI agenda, emphasising traceability, explainability, and governance in line with EU and ISO standards.

Emergent Low-Frequency Annotation Errors. Beyond the 18 high-frequency error types, our analysis also revealed several low-frequency but practically significant error types across the completeness, accuracy, and consistency dimensions. From a completeness perspective, these low-mention cases point to partial or uneven data coverage, including missing modalities, skipped frames, and restrictive dataset curation driven by technical, legal, or organisational constraints [29, 36, 38, 41]. We also found automation related omissions, for instance, when foundation models overlook rare or privacy sensitive instances and quality lapses such as incomplete pre-annotations or malfunction induced data loss [3, 39]. For accuracy, the low-frequency issues reflect small spatial or temporal misalignments, hybrid or automated labelling errors, and human factors such as fatigue, bias, and overreliance on automation [18, 24, 25, 32, 34, 38]. Even minor inaccuracies in annotation precision can propagate through training pipelines, degrading perception reliability. Under consistency, low-frequency cases reveal gradual or structural drifts in annotation uniformity across projects or over time [3, 17]. These arise from evolving guidelines, automation carryover errors, and mismatches between annotators’ understanding and supervisors’ expectations [34, 39]. Variations in workload, motivation, and team calibration, along with limited onboarding or expert oversight, further amplify such inconsistencies [8, 24]. Collectively, these low-frequency but conceptually important errors reveal emerging challenges in modern annotation pipelines, where automation, organisational constraints, and human variability intersect to shape data quality beyond the dominant error categories captured in the taxonomy.

Industrial Validation and Implications. Industry experts across OEMs, Tier-1, and Tier-2 suppliers validated the taxonomy’s clarity, relevance, and operational value. They characterised it as a “failure-mode catalogue” analogous to FMEA, an interpretation absent in prior literature. The taxonomy enables direct traceability from perception anomalies (*e.g.*, false positives, missed detections) to upstream data defects (*e.g.*, incomplete coverage, misaligned hand-offs). It provides a shared vocabulary for diagnosing whether an issue originates from human annotation, automation, or tooling, thereby improving cross-tier communication and quality control. Experts also recognised its educational value for onboarding, internal QA reviews, and continuous improvement, and suggested integrating it into annotation dashboards to institutionalise a data-quality culture.

Trustworthy AI Development and Data Governance Implications. The data annotation errors taxonomy contributes to the broader trustworthy AI agenda by linking data quality, process transparency, and accountability across the AI systems development lifecycle. It supports traceability of errors from perception anomalies to underlying annotation causes, an essential capability for compliance with emerging standards such as ISO/IEC 5259,

IEEE P2801, ISO 26262, SAE J3016 lifecycle assurance with these frameworks. The taxonomy provides a practical foundation for governance-by-design, where annotation quality becomes an auditable and traceable element of AI system certification. In doing so, it operationalises fairness, transparency, and accountability as measurable engineering properties rather than abstract ethical goals.

Implications for the AI-Enabled System Development Lifecycle. The taxonomy influences all phases of AI-enabled system development by embedding data quality considerations into engineering workflows. During *data collection and curation*, it improves dataset representativeness by identifying completeness-related gaps. In *annotation and tooling*, it guides workflow design to mitigate human-automation drift and reinforce auditability. For *model training and evaluation*, it provides traceability between annotation noise and model reliability, informing retraining priorities. Within *verification and validation (V&V)*, it supports safety traceability by linking perception errors to underlying data issues. In *requirements engineering*, it formalises data-specific criteria such as inter-annotator agreement thresholds, while in *quality assurance and maintenance*, it establishes a foundation for audits, supplier evaluations, and continuous improvement. These cross-cutting implications position the taxonomy as a unifying framework aligning data governance, model assurance, and system level quality management across the AI lifecycle.

Implications for Research and Practice. Beyond AIePS development life cycle integration, the taxonomy also offers broader implications for both research and industrial practice. For research purposes, this study establishes an empirically grounded taxonomy that systematises how data annotation errors manifest in the development of AIePS. It highlights the need for future research to move beyond isolated notions of “data quality” toward integrated frameworks linking annotation processes, organisational factors, and model performance. The taxonomy can serve as a foundation for quantitative studies measuring the impact of specific error types on model outcomes, as well as for developing automated tools that detect or predict annotation errors. For practice, the taxonomy provides practitioners with a structured lens for diagnosing and preventing quality issues throughout the annotation lifecycle. It can inform the design of quality assurance checklists, onboarding materials, and supplier evaluation criteria, fostering shared terminology across OEMs, Tier-1s, and annotation vendors. By identifying high-impact error types such as missing feedback loops, insufficient guidance, and weak quality control organisations can prioritise interventions that strengthen both annotation quality and governance. Moreover, the taxonomy encourages proactive error monitoring and the integration of feedback into existing annotation workflows, enabling continuous improvement. It also supports traceability by linking annotation errors to downstream perception failures, enhancing transparency in safety-critical development. Finally, it establishes a conceptual foundation for data governance standards and policy frameworks, ensuring that annotation quality is treated as a first-class element of system assurance rather than a peripheral task.

Future Work. Future research should expand validation across diverse geographical and industrial contexts to strengthen generalisability. Another critical direction is the development of semi-automated quality assurance systems that apply the taxonomy for

real-time monitoring, drift detection, and prioritised correction, thus enabling continuous data assurance throughout the AI lifecycle. Furthermore, quantitative studies correlating the frequency and severity of specific error types with model performance metrics would deepen understanding of their practical impact. Integrating severity probability matrices, similar to FMEA (Failure Mode and Effects Analysis) scoring, can support risk-based prioritisation in industrial annotation workflows. By embedding such automated mechanisms into annotation tools and MLOps pipelines, future work can transform the taxonomy from a diagnostic artefact into a proactive engineering instrument for lifecycle assurance in SE4AI.

ACKNOWLEDGMENTS This research is supported by Vinova, Program Fordons strategisk Forskning och Innovation (FFI), Project FAMER (2023-00771).

References

- [1] Mrinal R Bache and Javed M Subhedar. 2021. Autonomous driving architectures: insights of machine learning and deep learning algorithms. *Machine Learning with Applications* 6 (2021), 100164.
- [2] Carlo Batini, Cinzia Cappiello, Chiara Francalanci, and Andrea Maurino. 2009. Methodologies for data quality assessment and improvement. *ACM Computing Surveys (CSUR)* 41, 3 (2009), 1–52. doi:10.1145/1541880.1541883
- [3] Jacob Beck. 2023. Quality aspects of annotated data: A research synthesis. *AStA Wirtschafts- und Sozialstatistisches Archiv* 17, 3 (2023), 331–353.
- [4] Markus Borg, Jens Henriksson, Kasper Socha, Olof Lennartsson, Elias Sonnsjö Lönegren, Thanh Bui, Piotr Tomaszewski, Sankar Raman Sathyamoorthy, Sebastian Brink, and Mahshid Helali Moghadam. 2023. Ergo, SMIRK is safe: a safety case for a machine learning component in a pedestrian automatic emergency brake system. *Software quality journal* 31, 2 (2023), 335–403.
- [5] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in qualitative research. *Qualitative Research in Psychology* 3, 2 (2006), 77–101.
- [6] Mario I Chacon-Murguia and Claudia Prieto-Resendiz. 2015. Detecting driver drowsiness: A survey of system designs and technology. *IEEE Consumer Electronics Magazine* 4, 4 (2015), 107–119.
- [7] Rung-Ching Chen, Vani Suthamathi Saravananarajan, Long-Sheng Chen, and Hui Yu. 2022. Road segmentation and environment labeling for autonomous vehicles. *Applied Sciences* 12, 14 (2022), 7191.
- [8] Harshil Darji, Jelena Mitrović, and Michael Granitzer. 2024. Challenges and Considerations in Annotating Legal Data: A Comprehensive Overview. *arXiv preprint arXiv:2407.17503* (2024).
- [9] Florenc Demrozi, Marin Jereghi, and Graziano Pravadelli. 2021. Towards the automatic data annotation for human activity recognition based on wearables and BLE beacons. In *2021 IEEE International Symposium on Inertial Sensors and Systems (INERTIAL)*. IEEE, 1–4.
- [10] Sangeeta Dey and Seok-Won Lee. 2023. A Multi-layered collaborative framework for evidence-driven data requirements engineering for machine learning-based safety-critical systems. In *Proceedings of the 38th ACM/SIGAPP Symposium on Applied Computing*. 1404–1413.
- [11] Teodor Fredriksson, David Issa Mattos, Jan Bosch, and Helena Holmström Olsson. 2020. Data labeling: An empirical investigation into industrial challenges and mitigation strategies. In *International Conference on Product-Focused Software Process Improvement*. Springer, 202–216.
- [12] Luiz G Galvão and M Nazmul Huda. 2023. Pedestrian and vehicle behaviour prediction in autonomous vehicle system—A review. *Expert Systems with Applications* (2023), 121983.
- [13] Khan Mohammad Habibullah, Hans-Martin Heyn, Gregory Gay, Jennifer Horkoff, Eric Knauss, Markus Borg, Alessia Knauss, Håkan Sivencrona, and Jing Li. 2023. Requirements engineering for automotive perception systems: An interview study. In *International Working Conference on Requirements Engineering: Foundation for Software Quality*. Springer, 189–205.
- [14] Alan R. Hevner, Salvatore T. March, Jinsoo Park, and Sudha Ram. 2004. Design Science in Information Systems Research. *MIS Quarterly* 28, 1 (2004), 75–105. doi:10.2307/25148625
- [15] Hans-Martin Heyn, Khan Mohammad Habibullah, Eric Knauss, Jennifer Horkoff, Markus Borg, Alessia Knauss, and Polly Jing Li. 2023. Automotive perception software development: An empirical investigation into data, annotation, and ecosystem challenges. In *2023 IEEE/ACM 2nd International Conference on AI Engineering—Software Engineering for AI (CAIN)*. IEEE, 13–24.
- [16] Asad J Khattak, Numan Ahmad, Behram Wali, and Eric Dumbaugh. 2021. A taxonomy of driving errors and violations: Evidence from the naturalistic driving study. *Accident Analysis & Prevention* 151 (2021), 105873.
- [17] Jan-Christoph Klie, Richard Eckart de Castilho, and Iryna Gurevych. 2024. Analyzing dataset annotation quality management in the wild. *Computational Linguistics* (2024), 1–48.
- [18] Christopher Klugmann, Rafid Mahmood, Guruprasad Hegde, Amit Kale, and Daniel Kondermann. 2024. No Need to Sacrifice Data Quality for Quantity: Crowd-Informed Machine Annotation for Cost-Effective Understanding of Visual Data. *arXiv preprint arXiv:2409.00048* (2024).
- [19] Frank Krüger. 2022. Keynote: Adventures in Annotation: Providing High Quality Labels for Supervised Machine Learning. In *2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)*. IEEE, 254–254.
- [20] Vipin Kumar Kukkala, Jordan Tunnell, Sudeep Pasricha, and Thomas Bradley. 2018. Advanced driver-assistance systems: A path toward autonomous vehicles. *IEEE Consumer Electronics Magazine* 7, 5 (2018), 18–25.
- [21] Joseph Maxwell. 1992. Understanding and validity in qualitative research. *Harvard Educational Review* 62, 3 (1992), 279–301.
- [22] Maab Mohammedali and Muntasir Adam. 2023. The influence of data annotation process requirements on performance criteria of ML models. *Gothenburg University Library* (2023).
- [23] Alireza Najafi and Azzedine Boukerche. 2024. On The Performance of Perception Systems of Autonomous Vehicles. In *ICC 2024-IEEE International Conference on Communications*. IEEE, 5365–5370.
- [24] Emily Ohman. 2020. Challenges in Annotation: Annotator Experiences from a Crowdsourced Emotion Annotation Task. In *DHN*. 293–301.
- [25] Sarina Penquitt, Jonathan Klees, Rinor Cakaj, Daniel Kondermann, Matthias Rottmann, and Lars Schmarje. 2025. From Label Error Detection to Correction: A Modular Framework and Benchmark for Object Detection Datasets. *arXiv preprint arXiv:2508.06556* (2025).
- [26] Heinrich Peters, Alireza Hashemi, and James Rae. 2024. Generalizable Error Modeling for Human Data Annotation: Evidence From an Industry-Scale Search Data Annotation Program. *ACM Journal of Data and Information Quality* 16, 3 (2024), 1–15.
- [27] Leo L. Pipino, Yang W. Lee, and Richard Y. Wang. 2002. Data quality assessment. *Commun. ACM* 45, 4 (2002), 211–218. doi:10.1145/505248.506010
- [28] Paul Ralph, Nauman bin Ali, Sebastian Baltes, Domenico Bianculli, Jessica Diaz, Yvonne Dittrich, Neil Ernst, Michael Felderer, Robert Feldt, Antonio Filieri, Breno Bernard Nicolau de França, Carlo Alberto Furia, Greg Gay, Nicolas Gold, Daniel Graziotin, Pinjia He, Rashina Hoda, Natalia Juristo, Barbara Kitchenham, Valentina Lenarduzzi, Jorge Martinez, Jorge Melegati, Daniel Mendez, Tim Menzies, Jefferson Moller, Dietmar Pfahl, Romain Robbes, Daniel Russo, Nyyti Saarimäki, Federica Sarro, Davide Taibi, Janet Siegmund, Diomidis Spinellis, Mirosław Staron, Klaas Stol, Margaret-Anne Storey, Davide Taibi, Damian Tamburri, Marco Torchiano, Christoph Treude, Burak Turhan, Xiaofeng Wang, and Sira Vegas. 2020. Empirical Standards for Software Engineering Research. *arXiv:2010.03525 [cs.SE]* <https://arxiv.org/abs/2010.03525>
- [29] Federico Ruggeri, Eleonora Misino, Arianna Muti, Katerina Korre, Paolo Torroni, and Alberto Barrón-Cedeño. 2024. Let Guidelines Guide You: A Prescriptive Guideline-Centered Data Annotation Methodology. *arXiv preprint arXiv:2406.14099* (2024).
- [30] Per Runeson and Martin Höst. 2009. Guidelines for conducting and reporting case study research in software engineering. *Empirical software engineering* 14 (2009), 131–164.
- [31] V Samuktha, S Abhilash, Nitish Kumar, and P Rajalakshmi. 2024. A Framework for Object Classification via Camera-Radar Fusion with Automated Labeling. In *2024 IEEE Sensors Applications Symposium (SAS)*. IEEE, 1–6.
- [32] Marius Schubert, Tobias Riedlinger, Karsten Kahl, Daniel Kröll, Sebastian Schoenen, Siniša Šegvić, and Matthias Rottmann. 2024. Identifying label errors in object detection datasets by loss inspection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 4582–4591.
- [33] Devendra Sharma. 2020. Evaluation and Analysis of Perception Systems for Autonomous Driving. US Patent 10,872,251.
- [34] Anting Shen and Forrest Nelson Iandola. 2020. Automated annotation techniques. US Patent 10,872,251.
- [35] Gijs van Dijk, Carlos Aguilera, and Shashank M Chakravarthy. 2024. Deciphering disagreement in the annotation of EU legislation. *Artificial Intelligence and Law* (2024), 1–36.
- [36] Marceli Wac, Raul Santos-Rodriguez, Chris McWilliams, and Christopher Bourdeaux. 2023. Capturing requirements for a data annotation tool for intensive care: Experimental user-centered design study. *arXiv preprint arXiv:2309.16500* (2023).
- [37] Yair Wand and Richard Y. Wang. 1996. Anchoring data quality dimensions in ontological foundations. *Commun. ACM* 39, 11 (1996), 86–95. doi:10.1145/240455.240479
- [38] Yuning Wang, Zeyu Han, Yining Xing, Shaobing Xu, and Jianqiang Wang. 2024. A survey on datasets for the decision making of autonomous vehicles. *IEEE Intelligent Transportation Systems Magazine* (2024).

- [39] Eleanor Watson, Thiago Viana, and Shujun Zhang. 2023. Augmented Behavioral Annotation Tools, with Application to Multimodal Datasets and Models: A Systematic Review. *AI* 4, 1 (2023), 128–171.
- [40] Kai Yang, Xiaolin Tang, Jun Li, Hong Wang, Guichuan Zhong, Jiaxin Chen, and Dongpu Cao. 2023. Uncertainties in onboard algorithms for autonomous vehicles: Challenges, mitigation, and perspectives. *IEEE Transactions on Intelligent Transportation Systems* (2023).
- [41] Ziyuan Zhong, Zhisheng Hu, Shengjian Guo, Xinyang Zhang, Zhenyu Zhong, and Baishakhi Ray. 2022. Detecting multi-sensor fusion errors in advanced driver-assistance systems. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*. 493–505.