

## Autonomous Driving and Robotics

03/31/2026

The field of embodied artificial intelligence, as manifested in autonomous driving (AD) and service robotics, represents a transition from software-only AI systems to agents that interact with the physical world in real time. These domains share foundational reliance on perception, planning, and control pipelines, yet diverge in operational constraints, safety criticality, and deployment environments. This article analyzes the leading paradigm—end-to-end learning combined with modular safety architectures—by dissecting its core assumptions, implementation strategies, methods of field testing, approaches to safety assurance, and the reality of current deployments. The central argument is that while remarkable progress has been made in constrained operational design domains, the industry faces a fundamental tension between the statistical nature of learned competencies and the deterministic safety requirements of physical systems, a tension that current validation and safety architectures only partially resolve.

### Fundamental Assumptions

The development of autonomous driving and service robotics rests on several interlocking assumptions that shape both research directions and commercial strategies.

*The Sufficiency of Data-Centric Learning for Perception.* A primary assumption is that the perception problem—the ability to sense and interpret the environment—can be solved through supervised learning on massive, diverse datasets. Formally, let  $\mathcal{X}$  be the space of sensor inputs (camera images, lidar point clouds, radar returns) and  $\mathcal{Y}$  be the space of ground-truth world states (object positions, velocities, semantic labels). A perception model  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  is trained to minimize empirical risk

$$\frac{1}{N} \sum_{i=1}^N \mathcal{L}(f_\theta(x_i), y_i)$$

The assumption is that with sufficiently large  $N$  and model capacity,  $f_\theta$  will generalize to all scenarios encountered during deployment. This overlooks the long-tail problem: the distribution of real-world driving or service scenarios is heavy-tailed, with rare but safety-critical events (e.g., a tire bouncing across a highway, a child obscured by an occluding object) that may never appear in the training set. The assumption conflates *average-case* performance with *worst-case* safety.

*The Closed-Loop Learnability of Planning.* In many modern systems, planning and control are also learned, either through imitation learning (behavioral cloning from human demonstrations) or reinforcement learning in simulation. The underlying assumption is that these learned policies can capture the complexity of human-like driving or manipulation while generalizing beyond the training distribution. For imitation learning, the policy  $\pi$  is optimized to mimic an expert trajectory:

$$\pi^* = \arg \min_{\pi} \mathbb{E}_{s \sim d_{\text{expert}}} [\mathcal{L}(\pi(s), a_{\text{expert}})]$$

where  $s$  is the state and  $a$  the action. The critical flaw is the covariate shift problem: during deployment, the agent's own actions lead to states that were rarely or never seen in the expert demonstrations, causing compounding errors. This assumption underestimates the need for robust closed-loop controllers that can recover from novel situations.

*The Tractability of Operational Design Domain (ODD).* For both autonomous vehicles and service robots, the concept of an *operational design domain* (ODD) is central. The ODD specifies the conditions under which the system is designed to function (e.g., geofenced area, weather, lighting, road type). The assumption is that an ODD can be precisely defined, monitored, and enforced. In practice, ODD boundaries are often fuzzy and dynamic. For a delivery robot, “sidewalk” is a semantic category that varies across municipalities and changes with construction or temporary obstacles. The reliance on ODD as a safety boundary presumes a level of environmental predictability that does not exist in unstructured human spaces.

*The Homomorphism Between Simulation and Reality.* A foundational enabler is the belief that simulation can serve as a faithful proxy for the real world for the purposes of training and validation. This is expressed as an assumption of *sim-to-real transfer*: a policy trained in a simulator with dynamics  $T_{\text{sim}}$  and sensor models  $S_{\text{sim}}$  will perform effectively in the real world with dynamics  $T_{\text{real}}$  and  $S_{\text{real}}$ . The gap between  $T_{\text{sim}}$  and  $T_{\text{real}}$ —often termed the “reality gap”—is managed through domain randomization and system identification, but these techniques cannot eliminate the possibility of emergent behaviors in the real world that were not represented in the simulation distribution.

### Implementation Approaches

The implementation of autonomous driving and service robotics systems typically follows a hybrid architecture that combines learned components with traditional engineering.

*Perception Stack.* The perception system transforms raw sensor data into a structured representation of the environment. In leading implementations, this is accomplished using deep neural networks, often employing a multi-head architecture. For example, a single network might output object detection bounding boxes, lane markings, free-space segmentation, and depth estimates. These outputs are fused across modalities (cameras, lidar, radar) using sensor fusion algorithms, often based on *Kalman* filters or learned attention mechanisms. A critical concept is *semantic segmentation*, where each pixel (or point) is classified into categories such as “vehicle,” “pedestrian,” “curb.” The reliance on neural networks introduces epistemic uncertainty: the model outputs a point estimate without a principled measure of confidence, leading to challenges in downstream safety monitoring.

*Planning and Control.* The planning stack is typically hierarchical. A *route planner* computes a high-level path from origin to destination. A *behavioral planner* decides on tactical maneuvers (e.g., lane change, yield, stop) often using rule-based state machines or learned policies. Finally, a *local motion planner* generates a collision-free trajectory that satisfies kinematic constraints. In many service robots, this is implemented using model predictive control (MPC), which solves an optimization problem at each time step:

$$\min_{u_0, \dots, u_{N-1}} \sum_{k=0}^{N-1} \ell(x_k, u_k) + \ell_f(x_N) \quad \text{s.t.} \quad x_{k+1} = f(x_k, u_k), \quad x_k \in \mathcal{X}_{\text{free}}, \quad u_k \in \mathcal{U}$$

where  $f$  is the robot dynamics,  $\ell$  is a stage cost, and  $\mathcal{X}_{\text{free}}$  is the collision-free state space. The integration of learned components (e.g., a neural network predicting future trajectories of other agents) into the MPC formulation is an active area, but the coupling creates a challenge: the planner’s guarantees are only as reliable as the learned predictions.

*End-to-End Learning.* A more radical implementation approach seeks to replace the modular pipeline with a single neural network that maps raw sensors directly to control commands (e.g., steering,

acceleration). This end-to-end paradigm reduces engineering complexity but exacerbates the interpretability and safety verification problems. The network learns a function  $\pi_{e2e} : \mathcal{X} \rightarrow \mathcal{A}$ , where  $\mathcal{A}$  is the action space. The absence of intermediate representations makes it difficult to isolate the cause of failures or to provide formal safety assurances.

*Simulation Infrastructure.* Both AD and robotics companies invest heavily in high-fidelity simulators that model physics, sensors, and scenario diversity. These simulators are used for training (via reinforcement learning or imitation learning), for automated testing (running millions of simulated miles), and for generating synthetic training data. The implementation assumes that scenario generation can capture a great amount of variabilities of the real world, typically through combinatorial parameter sweeps or learned adversarial scenario generation.

## Field Testing

Field testing encompasses both closed-course evaluations and public-road or public-space deployments. The methodologies used reveal significant limitations in current validation paradigms.

*Miles-Driven Metric.* The industry has historically used “miles driven without disengagement” as a primary metric of progress. A disengagement occurs when the safety driver or the system’s autonomous supervisor takes control due to a failure. This metric is deeply flawed. It conflates ODD restrictiveness with capability: a system tested only on simple highways in fair weather will naturally achieve higher miles per disengagement than one tested in dense urban environments. Moreover, the metric does not account for the *severity* of potential failures; a system may drive millions of miles without a disengagement while still harboring a rare but catastrophic failure mode. Formally, the metric provides an estimate of the system’s *reliability* under the test distribution but offers no guarantee on the tail risk  $\Pr(\text{failure} \mid \text{scenario})$  for scenarios not adequately sampled.

*Scenario-Based Testing.* Recognizing the limitations of random miles, regulators and developers have adopted scenario-based testing. This involves identifying a set of critical scenarios (e.g., cut-ins, pedestrians running into the street, occluded intersections) and testing system performance on those scenarios, often in simulation or closed-course environments. The challenge lies in the completeness of the scenario catalog. For a system with  $n$  interacting agents, the state space grows combinatorially; exhaustively covering all critical interactions is computationally intractable. The assumption that a finite set of scenarios can represent all safety-relevant events is a necessary but unproven simplification.

*On-Road Pilots.* Commercial deployments often begin with supervised pilots: vehicles with safety operators, operating within geofenced areas. These pilots serve as both data-collection mechanisms and public demonstrations. However, the operational constraints (e.g., limited geographic area, time-of-day restrictions, favorable weather) mean that the system is never tested under the full range of conditions it would encounter in a true deployment. The transition from a pilot to a driver-less service (where no safety operator is present) requires a leap in confidence that current testing methodologies struggle to justify.

*Service Robotics Field Testing.* For service robots (e.g., delivery robots, warehouse autonomous mobile robots), field testing often occurs in controlled environments like warehouses or designated university campuses. These environments are simpler than public roads but still present challenges: unpredictable human behavior, narrow corridors, and dynamic obstacles. The test methodology typically relies on teleoperation fallback: when the robot encounters a situation it cannot handle, a

remote human operator takes over. This creates a validation paradox: the system is never forced to resolve the hardest edge cases autonomously, so its true capability remains untested.

## Safety Regulations

Safety regulations refers to the methods used to assure that the system will not cause harm. In autonomous driving and robotics, this is the most contested area, where statistical learning meets life-critical applications.

*Functional Safety (ISO 26262) and SOTIF (ISO 21448).* The automotive industry has adopted functional safety standards (ISO 26262) that address systematic and random hardware failures. For autonomous driving, a supplementary standard, ISO 21448 (Safety of the Intended Functionality, SOTIF), addresses the absence of unreasonable risk due to *functional insufficiencies*—*i.e.*, the AI system’s limitations. SOTIF introduces concepts such as *known unsafe scenarios* (which must be mitigated) and *unknown unsafe scenarios* (which must be discovered through validation). The standard’s applicability to learned systems is limited because it relies on hazard analysis and scenario identification that presume a degree of system transparency that neural networks do not provide.

*The Responsibility-Sensitive Safety (RSS) Model.* A prominent approach to safety resolution in autonomous driving is RSS, a formal model that defines “safe” behavior using a set of rules about following distance, right-of-way, and cautious driving. RSS provides a mathematical framework to prove that an autonomous vehicle will not cause an accident as long as it adheres to the rules and other actors behave reasonably. However, RSS does not address perception errors; it assumes that the vehicle has perfect knowledge of the positions and velocities of other actors. In practice, perception uncertainty renders the formal guarantees moot. Moreover, RSS’s rules are conservative and can lead to over-cautious behavior that makes the vehicle impractical in dense urban traffic.

*Runtime Monitoring and Fallback.* A widely adopted safety architecture uses a *monitor* that runs in parallel with the primary AI system, checking for conditions that violate safety constraints. If a violation is detected, a *fallback* controller (often simpler and formally verified) brings the system to a minimal risk condition (*e.g.*, stopping safely). Formally, the monitor implements a predicate  $\phi(s)$  that is true when the current state  $s$  is deemed safe. The safety guarantee is that the primary system is only allowed to operate while  $\phi(s)$  holds; otherwise, control is transferred. The challenge is designing  $\phi$  to be both sufficiently tight to prevent accidents and adequately loose to allow useful operation. For perception-based systems,  $\phi$  must reason about uncertain world states, making the monitor itself a complex system prone to false positives (nuisance disengagements) and false negatives (undetected unsafe states).

*Simulation-Based Safety Validation.* Given the impossibility of exhaustive real-world testing, companies attempt to validate safety through billions of simulated miles. The statistical argument is that if a system performs without failure over a sufficiently large number of simulated scenarios sampled from the expected deployment distribution, then the probability of a failure in deployment is acceptably low. This relies on the assumption that the simulation distribution is identical to the real-world distribution—an assumption that cannot be verified. Furthermore, the independence of scenarios is violated because real-world driving involves sequences of correlated events; a rare event in simulation may be a precondition for a subsequent failure.

## Deployments

*Autonomous Driving.* Publicly accessible autonomous driving services are largely limited to geofenced areas with significant operational constraints. *Waymo's* robotaxi service operates in parts of Phoenix, San Francisco, and a few other cities, with vehicles capable of driver-less operation but only under favorable weather, mapped roads, and relatively low-speed environments. Cruise similarly operated in San Francisco before a high-profile incident led to a suspension, illustrating the fragility of the safety case when confronted with unexpected scenarios (e.g., a vehicle becoming immobilized in an intersection, blocking emergency responders). These deployments are characterized by extensive remote assistance infrastructure: when the vehicle encounters a situation outside its ODD, a remote operator provides guidance. The actual level of autonomy is therefore conditional on the availability of human oversight.

*Robotaxis Versus Consumer Vehicles.* A critical distinction is between purpose-built robotaxis (with redundant sensors, compute, and safety systems) and consumer vehicles with “autopilot” or “full self-driving” capabilities. Consumer systems operate under the assumption of a human driver who remains responsible, which categorically changes the safety expectation. The deployment of these systems has led to numerous accidents, often due to the human’s failure to maintain supervision—a failure that the system’s design implicitly encourages by over-promising capability.

*Service Robotics.* In the service robotics sector, deployments are more varied and generally less safety-critical. Warehouse autonomous mobile robots (AMRs) from companies like Amazon Robotics operate in controlled environments with strict separation from humans, using centralized orchestration to prevent collisions. These deployments are successful because the ODD is tightly constrained and the environment can be instrumented with external safety systems (e.g., emergency stop buttons, physical barriers). In contrast, sidewalk delivery robots (e.g., *Starship*, *Kiwibot*) have been deployed in dozens of cities, operating at low speeds and with remote monitoring. Their safety record is generally good, but they face operational challenges: getting stuck on uneven terrain, being vandalized, or creating nuisances for pedestrians. These deployments have not achieved the scale or autonomy promised, often relying on a human-to-robot ratio far from the envisioned “fleet operates itself.”

## Concluding Remarks

The leading AI paradigms in autonomous driving and service robotics have enabled impressive demonstrations in constrained environments but have not yet delivered the general, robust autonomy that their foundational assumptions implied. The reliance on data-driven learning for perception and planning, while effective in the average case, struggles to provide the deterministic safety guarantees required for unconstrained deployment. Implementation approaches have gravitated toward hybrid architectures that use learning for perception and traditional control for safety-critical decisions, but the integration remains imperfect. Field testing methodologies, dominated by miles-driven metrics and simulation, cannot statistically validate the absence of rare but catastrophic failures. Safety regulation frameworks like SOTIF and RSS provide valuable structure but do not resolve the fundamental tension between statistical models and formal safety requirements. Actual deployments reveal a consistent pattern: successful systems operate within highly constrained operational design domains, often with human oversight or remote assistance. The gap between these deployments and the vision of universally capable autonomous agents remains wide, suggesting that future progress will depend less on scaling data and compute than on developing new architectures that integrate formal methods, uncertainty quantification, and robust safety envelopes into the core of learned systems.