# CLARIFYVC: CLARIFYING AMBIGUOUS COMMANDS IN VEHICLE CONTROL WITH A HYBRID DATA AUGMENTATION PIPELINE

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

Natural language interfaces for vehicle control must contend with vague commands, evolving dialogue context, and strict protocol constraints. We introduce Clari**fyVC**, a unified framework that integrates a hybrid data-augmentation pipeline (ClarifyVC-Data), reference models trained on the data (ClarifyVC-Models) and a evaluation protocol (ClarifyVC-Eval). The agent-orchestrated pipeline generates diverse, ambiguity-rich dialogues from real-world seeded queries under schema and safety constraints, while the evaluation protocol systematically probes single-turn parsing, conservative clarification under extreme fuzziness, and multi-turn grounding. Fine-tuning on ClarifyVC-Data yields consistent gains—up to 15% higher parsing accuracy, 20% stronger ambiguity resolution, and 98% protocol compliance—across realistic in-cabin scenarios, with human-in-the-loop assessments confirming high realism, coherence, and applicability. ClarifyVC thus advances beyond simulation-only datasets by tightly coupling real-world grounding with scalable generation and standardized evaluation, and provides a generalizable pipeline for broader interactive control domains. Our code and dataset are available at: https://anonymous.4open.science/r/ClarifyVC.

# 1 Introduction

Natural language interfaces are becoming a cornerstone of interactive control systems, from autonomous vehicles (Wen et al., 2024) to smart homes (Thukral et al., 2025), robotics (Sikorski et al., 2025), and other embodied agents (Bick et al., 2024). These systems require the ability to interpret vague instructions, maintain multi-turn dialogue context, and execute actions under strict protocol constraints. In the automotive domain, the rise of autonomous vehicles has already transformed human–machine interaction, making natural language commands crucial for intuitive and trustworthy control of hundreds of onboard functions (Zheng et al., 2024; Wang et al., 2024a). However, Vehicles face pervasive ambiguity—user commands are often vague, protocol mappings incomplete, and existing evaluation metrics inadequate (Ma et al., 2024). Traditional intent detection and slot-filling methods perform poorly under ambiguity and context shift, while current benchmarks lack realism, coverage, or failsafe metrics (Chun et al., 2025).

Public perception reflects these gaps: 58% of individuals feel uneasy about self-driving cars, and 25% express complete distrust in their reliability (Wenskovitch et al., 2024; Peng & Shang, 2024). A core reason is that current LLMs, though strong in general reasoning, struggle in safety-critical control (Brahman et al., 2024). They hallucinate under ambiguous instructions, fail to request clarifications when uncertain, and lack strict protocol adherence in task orchestration (Dai et al., 2024). These weaknesses are compounded by the absence of high-quality, reality-grounded datasets and standardized evaluation protocols, limiting progress toward reliable in-vehicle dialogue systems (Nguyen et al., 2024; Zou et al., 2024). To close this gap, we introduce ClarifyVC, a unified framework for clarifying ambiguous commands in vehicle control. It integrates a hybrid data-augmentation pipeline (ClarifyVC-Data), reference models trained on the data (ClarifyVC-Models) and a three-tier evaluation protocol (ClarifyVC-Eval) to evaluate the data quality and model performance.

At the core of ClarifyVC-Data is a hybrid augmentation pipeline seeded from over 20k authentic in-vehicle commands drawn from a proprietary corpus of 4M+ production-level interactions. Through

structured ambiguity injection, adversarial perturbations, and multi-turn clarification, the pipeline synthesizes ambiguity-rich yet protocol-compliant samples that target robustness and safe execution. The resulting dataset, ClarifyVC-Data, has been validated through human evaluation and distributional alignment experiments, demonstrating close correspondence to real-world usage patterns. Fine-tuning LLMs on this data yields an average 15% improvement in parsing accuracy, underscoring the pipeline's practical value for safety-critical language interfaces.

Beyond dataset construction and training, we introduce ClarifyVC-Eval, a three-tier evaluation protocol that explicitly targets real-world ambiguity in function-call tasks (Jiang et al., 2024b; Chao et al., 2024; Wu et al., 2024; Jiang et al., 2024a). ClarifyVC-Eval plays a dual role: (i) it audits the benchmark itself—testing whether the data is realistic and ambiguity-rich—and (ii) it provides a unified lens to evaluate model capabilities in semantic parsing, execution fidelity, and safety compliance. By jointly assessing data validity and functional reliability, the protocol addresses a key gap in prior work, which typically isolates dataset realism from model accuracy and thus misses their interaction in safety-critical settings. To operationalize the data-side audit, we additionally define a *Dataset Quality Score (DQS)* that aggregates ambiguity diversity (AD), protocol compliance (PC), and realism (R). Together, ClarifyVC-Eval and DQS constitute a comprehensive, scalable framework for auditing datasets and benchmarking models under realistic ambiguity.

Extensive experiments demonstrate that ClarifyVC substantially improves performance in safety-critical control tasks. Fine-tuned models achieve 15% higher parsing accuracy, 20% better ambiguity resolution, and 98% protocol compliance, while also reducing inference latency by 30% compared to baseline systems. Additional ablation studies confirm the necessity of each module in the pipeline, with the default configuration yielding the best trade-off between diversity, coherence, and adherence. Multi-run evaluations further validate robustness, showing consistently low variance (<1%) and statistically significant improvements across metrics. Human-in-the-loop assessments corroborate these results, with expert annotators rating generated dialogues highly on realism, coherence, and practical applicability. Together, these findings highlight ClarifyVC as a reliable and efficient framework for robust language understanding under real-world ambiguity in vehicle control. In summary, our contributions are threefold:

- 1. **ClarifyVC Framework**: A unified framework for clarifying ambiguous commands in vehicle control and interactive systems. It integrates a hybrid data pipeline, a three-tier evaluation protocol, and reference models, offering an end-to-end standard for safe and deployable language interfaces.
- 2. **ClarifyVC-Data&Models**: A hybrid, reality-grounded, and human-validated dataset built from 20k+ real-world seed commands, expanded with controlled fuzziness and adversarial variants. By training on the high-quality data, we release reference models that show consistent gains in accuracy, clarification, and safety compliance.
- 3. ClarifyVC-Eval: A three-tier evaluation protocol that disentangles under-specification, ambiguity clarification, and multi-turn grounding, along with a Dataset Quality Score which ensures the benchmark aligns with real-world distributions and maintains high-quality standards. By explicitly targeting these failure families, the protocol enables comprehensive and safety-aware assessment of function-call understanding, addressing gaps left by conventional single-turn accuracy.

# 2 Related Work

#### 2.1 METHODS FOR CLARIFYING AMBIGUITY AND MULTI-TURN COMMAND PARSING

Natural-language command understanding has progressed from structure-aware parsers to end-to-end LLM solutions for mapping utterances to executable actions (Zheng et al., 2024; Wang et al., 2024a). Early pipelines emphasized schema-constrained intent/slot structures and hierarchical modeling (Sriram et al., 2019; Wang et al., 2024b; Okur et al., 2023). More recently, LLMs have been applied to enable direct intent grounding, rule translation, and task formalization (Shao et al., 2024; Choudhary et al., 2024; Manas & Paschke, 2023).

In the domain of vehicle or visual command understanding, datasets like Talk2Car (Deruyttere et al., 2019), CI-AVSR (Dai et al., 2022a), and doScenes (Roy et al., 2024) provide real-world instruction—action pairs and visual grounding contexts, but primarily support single-turn mapping rather than interactive clarification. Beyond single-turn parsing, logical disambiguation methods such

as LogicalBeam (Bhaskar et al., 2023) have been explored, while other datasets and studies (e.g. CHAMBI) highlight cross-cultural or spatial ambiguity challenges (Zhang et al., 2024b; Saparina & Lapata, 2024). Frameworks for task decomposition and retrieval-augmented decision making further support complex instruction following (Shen et al., 2024; Yang et al., 2024). Parallel streams examine unimodal parsing (Zhang et al., 2024a), synthetic data generation (Liu et al., 2024), and distillation for instruction following (Ding et al., 2024), alongside domain-specific command datasets (Liu et al., 2023; Li et al., 2024). Together, these approaches offer modeling, training, and data-centric tools for tackling ambiguity and multi-turn semantics in command parsing.

#### 2.2 LIMITATIONS OF EXISTING APPROACHES AND OUR POSITIONING

Despite steady progress, three gaps persist. (1) *Ambiguity management and uncertainty signaling*. Many systems lack explicit mechanisms to detect under-specification, trigger clarifying questions, or expose calibrated confidence, which is critical in safety-sensitive interaction (Pramanick et al., 2022; Wenskovitch et al., 2024; Lee et al., 2024). (2) *Evaluation scope*. Benchmarks often emphasize single-turn parsing, narrow modalities, or synthetic distributions (Zhang et al., 2024a; Liu et al., 2024; 2023; Li et al., 2024), offering limited coverage of multi-turn grounding and protocol-aware execution; even task-centric frameworks (Shen et al., 2024; Yang et al., 2024) provide only partial visibility into clarification behavior. (3) *Data realism and compliance*. Instruction-generation and distillation pipelines (Ding et al., 2024) rarely tie ambiguity to real logs or enforce function-call protocols, hindering transfer to deployed systems.

ClarifyVC addresses these gaps with a unified framework that: (i) couples a hybrid, real-log-seeded augmentation pipeline with controlled adversarial/fuzzy evolution to surface realistic ambiguity; (ii) introduces a three-tier evaluation protocol that disentangles single-turn parsing, clarification under extreme fuzziness, and multi-turn grounding with execution checks; and (iii) reports reference models trained on the data with protocol-aligned metrics. Our data generation pipeline, benchmark, and evaluation protocol demonstrate strong generalization and practical utility in real world settings. Beyond in-cabin voice control, the framework's clarification strategies and compliance-oriented evaluation naturally extend to broader domains of human-machine interaction, including smart homes, medical dialogue, and embodied intelligence, where safety and interpretability remain critical.

# 3 METHODOLOGY

We instantiate **ClarifyVC** as a unified framework comprising a hybrid data-augmentation pipeline (*ClarifyVC-Data*), reference models trained on the data (*ClarifyVC-Models*) and a safety-aware three-tier protocol (*ClarifyVC-Eval*). Seeded with 20k+ authentic in-vehicle commands (drawn from 4M+ production logs), the pipeline expands queries via structured ambiguity injection, adversarial perturbations, and multi-turn clarification under protocol constraints. The resulting corpus is validated through human studies and distributional alignment with real-world usage.

# 3.1 CLARIFY VC-DATA: AGENT-ORCHESTRATED PIPELINE

**Stage-wise pipeline.** We adopt an *agent-orchestrated* modular pipeline with four stages: Semantic Parsing Module (SPA), Adversarial Generation Module (AGA), Fuzz Injection Module (FIA), Multi-Turn Evolution Module (MEA), each implemented with prompt-engineered, pre-trained LLMs without task-specific fine-tuning. Specifically, SPA/FIA/MEA are implemented with *DeepSeek-R1* (API-based) for semantic parsing, fuzz injection, and multi-turn dialogue evolution, while AGA is realized with *Qwen2.5-72B* (via vLLM) to perform protocol-constrained adversarial rewriting. These choices combine scalability and strong instruction-following, while ensuring modularity for drop-in substitution. The total compute cost of synthesizing the benchmark remains modest, as generation relies primarily on API calls and lightweight orchestration:

- SPA parses each seed command into (I, E, P) as standardized grounding.
- AGA produces syntactically valid yet ambiguous variants  $c_{\text{adv}}$  under protocol constraints.
- FIA converts c<sub>adv</sub> into softer fuzzed instructions c' (parameter omission, subjective modifiers, mild distortion); both tiers are retained.
- MEA expands c' into coherent multi-turn dialogues D for long-horizon grounding.

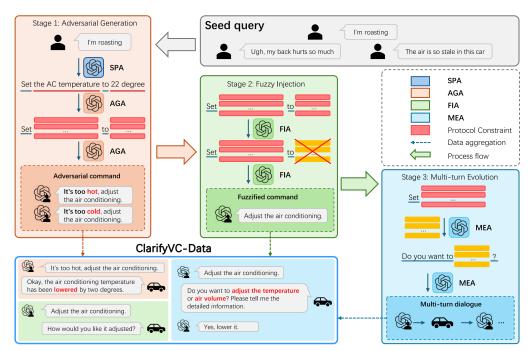


Figure 1: Agent-orchestrated, stage-wise generation flow. A schema-constrained pipeline executes semantic parsing, adversarial construction, fuzz injection, and multi-turn evolution on real-world-seeded commands to synthesize ambiguity-rich single- and multi-turn dialogues under protocol constraints. The resulting corpus forms Clarify VC-Data: a hybrid, realism-aligned, human-validated benchmark with standardized function-call annotations and broad ambiguity coverage.

The adopted sequence (SPA  $\rightarrow$  AGA  $\rightarrow$  FIA  $\rightarrow$  MEA) is not arbitrary but empirically validated. We conducted controlled ablations that permuted the order or removed individual stages. Evaluation results (Table 9) show that the default order achieves the best balance across ambiguity diversity, dialogue coherence, and protocol adherence. For example, reversing FIA and AGA substantially reduced diversity, while removing either stage markedly degraded ambiguity coverage. This confirms that the chosen order is optimal for synthesizing realistic yet challenging interactions.

This yields a hierarchical pool: (1) SPA+AGA  $\Rightarrow c_{\text{adv}}$ ; (2) +FIA  $\Rightarrow c'$ ; (3) +MEA  $\Rightarrow D$ . To encourage diversity while preserving operational validity, each sample is scored by

$$Q(c) = \alpha \cdot H(c) + (1 - \alpha) \cdot \mathbb{I}(c \text{ is protocol-compliant}), \quad \alpha = 0.6, \tag{1}$$

where H(c) is ambiguity entropy and  $\mathbb{I}(\cdot)$  indicates compliance (Appendix B.2).

**Reference models (ClarifyVC-Models).** We obtain ClarifyVC-Models by supervised fine-tuning open-source backbones (e.g., LLaMA3-8B, Qwen2.5-7B/72B, DeepSeek-R1-Distilled) on ClarifyVC-Data with schema-aligned function-call targets, using a teacher-forced cross-entropy objective and JSON-schema—constrained decoding at inference(Experiment settings can be seen in Appendix B). Training is performed with early stop on a delayed test split and evaluated on a separate 2k test set, averaged on 5 random seeds (std.,< 1%). We release the **Qwen2.5-7B-SFT** checkpoint and training/evaluation configs. Notably, while larger backbones (14B, 32B, 72B) show strong results, the 7B model achieves the best trade-off between accuracy and computational efficiency, reducing inference cost by an order of magnitude while delivering comparable or superior performance under our protocol. As will be shown in Table 3, the models fine-tuned with ClarifyVC-Data consistently surpass the zero-shot base models in different scenarios.

#### 3.2 CLARIFYVC-EVAL: EVALUATING DATASET QUALITY AND MODEL PERFORMANCE

In this part, we introduce ClarifyVC-Eval to evaluate the quality of ClarifyVC-Data and ClarifyVC-Models. A Dataset Quality Score (DQA) is used to ensure the benchmark aligns with real-world distributions and maintains high-quality standards. Meanwhile, a three-tier evaluation protocol that

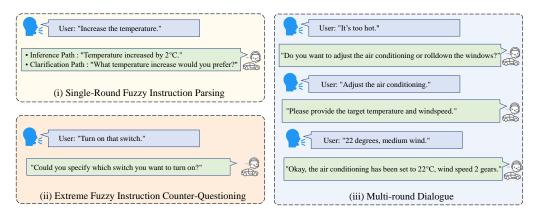


Figure 2: Illustration of ClarifyVC-Eval: (i) parse mildly fuzzy commands into precise function calls, (ii) adopt safe clarification under extreme vagueness, and (iii) sustain multi-round dialogue for coherent, grounded execution—capturing ambiguity, safety, and interactivity in real-world control.

disentangles under-specification, ambiguity clarification, and multi-turn grounding, is proposed to evaluate the generation quality of different LLMs.

Automated quality and human validation for Dataset. We summarize dataset quality with

$$DQS = \lambda_1 \cdot AD + \lambda_2 \cdot PC + \lambda_3 \cdot R, \quad (\lambda_1, \lambda_2, \lambda_3) = (0.4, 0.3, 0.3), \tag{2}$$

combining ambiguity diversity (AD), protocol compliance (PC), and realism (R). We choose (0.4, 0.3, 0.3) as the hyperparameters of Dataset Quality Score (DQS) via a grid search that maximizes Spearman correlation with human ratings while preserving rank stability across baselines; full sweeps appear in Appendix C.4.

Three-Tier Protocol for Comprehensive Model Evaluation. While dataset validation secures distributional realism, robust model assessment requires an evaluation protocol that can reveal failure families invisible to single-turn accuracy. To this end, when evaluating models, ClarifyVC-Eval operationalizes three complementary tiers of evaluation: (i) Single-round fuzzy instruction parsing, which tests the model's ability to parse mildly ambiguous commands by disambiguating challenges such as under-specified parameters, vague references, and subjective expressions; (ii) Extreme fuzzy instruction counter-questioning, examining whether the model adopts safe clarification strategies when confronted with severe ambiguity, specifically its capacity to detect extreme uncertainty and ask relevant clarifying questions; and (iii) Multi-turn dialogue, assessing the ability to address challenges like multi-turn dependency and memory, which is crucial for iteratively recovering missing semantics, maintaining dialogue and parameter coherence, and executing the accumulated commands reliably. The protocol ClarifyVC-Eval, as illustrated in Figure 2, spans single-turn parsing, clarification under extreme fuzziness, and multi-turn dialogue grounding, addressing gaps in existing evaluation metrics and enabling more realistic assessment of safety, robustness, and interactivity in control-oriented language interfaces

Table 1: The three-tier ClarifyVC-Eval protocol to evaluation metrics. Each tier isolates a distinct family of failure modes while jointly covering the spectrum of function-call understanding and safe execution.

Tier	Metrics Used	Rationale
Tier 1: Single-Round Instruction Fuzzy Pars- ing	Intent Recognition Accuracy (IRA), Parameter Extraction Precision (PEP), Intent Hit Rate (IHR), Function Hit Rate (FHR)	Captures the model's ability to resolve under-specified single-turn com- mands into correct intents and API calls. These metrics reflect semantic accuracy and parameter precision at the most basic function-call level.
Tier 2: Extreme Fuzzy Instruction Counter- Questioning	Fuzzy Detection Rate (FDR), Counter- Question Coverage (CQC), Protocol Com- pliance Rate (PCR)	Evaluates whether the model identifies extreme ambiguity and adopts safe clarification strategies instead of unsafe guesses. Metrics track conservative behavior, protocol adherence, and safety awareness.
Tier 3: Multi-turn Dialogue	Dialogue Consistency (DC), Final Execution Success Rate (FESR), Parameter Completeness (F1-score)	Assesses long-horizon interactions where the model must gather missing semantics over multiple turns, maintain coherence, and ultimately ground safe executable commands. These metrics measure the culmination of dialogue fidelity and execution success.

Table 2: Comprehensive quality evaluation of ClarifyVC-Data. Automated metrics (left) benchmarked against baselines; human validation (right) conducted only for ClarifyVC-Data.

(a) Comparison with baselines using four metrics.

(b) Human validation (ClarifyVC-Data only)

Dataset	AD	PC	R	DQS
Talk2Car	0.50	0.85	0.60	0.62
doScenes (Roy et al., 2024)	0.56	0.81	0.64	0.65
CI-AVSR (Dai et al., 2022b)	0.53	0.82	0.61	0.64
DeepSeek Distilled	0.55	0.80	0.65	0.65
GPT-o1 Distilled	0.60	0.82	0.70	0.69
Qwen2.5 Distilled	0.58	0.78	0.68	0.67
LLaMA3 Distilled	0.62	0.80	0.72	0.70
ClarifyVC-Data	0.89	0.95	0.82	0.88

Score	Agreement
$4.6\pm0.2$	93%
$4.6 \pm 0.1$	96%
$4.7 \pm 0.2$	94%
$4.5\pm0.3$	91%
	$4.6 \pm 0.2$ $4.6 \pm 0.1$ $4.7 \pm 0.2$

As summarized in Table 1, ClarifyVC-Eval is structured into three complementary tiers, each probing a distinct capability of interactive control models. Tier 1 targets the core ability of *semantic parsing under underspecification*, using intent- and function-level accuracy metrics (IRA, PEP, IHR, FHR) to capture fine-grained correctness of disambiguated calls. Tier 2 focuses on *safe clarification under extreme fuzziness*, where FDR, CQC, and PCR jointly test whether models recognize ambiguity, avoid unsafe guesses, and adhere to interaction protocols. Tier 3 addresses *long-horizon multiturn grounding*, evaluated by DC, FESR, and parameter completeness, ensuring that models can sustain coherent dialogue and achieve reliable execution outcomes. The detailed definitions and computation formulas of these metrics can be found in Appendix C. Collectively, these metrics provide comprehensive coverage of the decision points most critical to safe and effective deployment, extending beyond what single-turn accuracy alone can capture.

**Rationale and scope.** Our analysis of 20k+ real-world in-vehicle logs shows that failures cluster into three families: under-specification, insufficient clarification, and long-horizon grounding. These are precisely captured by ClarifyVC-Eval, which not only measures success rates but also tracks protocol violations, yielding diagnostics that better reflect operational safety. Importantly, the same decision points recur across broader HCI and embodied intelligence (e.g., robotics, smart environments), making the protocol directly transferable beyond the vehicle domain.

# 4 EXPERIMENTAL

We conduct extensive experiments to address the following research questions: **RQ1**: What's the quality of ClarifyVC-Data evaluated under DQS in Equation 2 and human-grounded validation? **RQ2**: How well do existing LLMs handle complex and ambiguous vehicle control instructions? **RQ3**: Does fine-tuning on ClarifyVC-Data improve model performance in realistic command understanding? **RQ4**: How accurately can LLMs execute structured function calls under protocol constraints? **RQ5**: Can open-source models, when properly tuned, match or surpass proprietary models in vehicle control tasks? The complete experimental setup, including the experimental environment and the agent-orchestrated pipeline used to generate the evaluation test sets, is provided in Appendix B.

# 4.1 EVALUATION ON THE DATA QUALITY (RQ1)

As mentioned in Section 3.2, we introduced DQS and human validation to assess the quality of the dataset. Table 2(a) shows that ClarifyVC-Data exceeds previous datasets and distilled baselines on the four axes. In order to ensure that the constructed benchmark aligns with real-world distributions and maintains high-quality standards, we conducted a dedicated human-grounding study on 500 sampled dialogues, as shown in Table 2(b). Five independent annotators rated each sample on linguistic realism, plausibility of ambiguity, coherence of dialogue, and practical applicability. The results confirm high realism (4.6/5), strong plausibility of ambiguity (96% agreement), coherent dialogues (94%), and strong applicability (91%), verifying that ClarifyVC-Data maintains practical quality in addition to statistical robustness. Together, these results confirm that ClarifyVC-Data not only surpasses prior datasets in automated measures(more results can be seen in Appendix A) but also passes stringent human-grounding validation, ensuring both scalability and real-world applicability.

Table 3: Multi-run evaluation on Clarify VC-Data under Zero-shot (ZS), Few-shot (FS), and SFT. Entries are *means over 5 independent runs*; per-cell standard deviation is < 1.0 percentage point (pp) across all columns (Appendix B, B.4). For readability we report means; the "Max  $\sigma$  (pp)" row summarizes the largest observed standard deviation in each column. All models are evaluated on a held-out test set of 5k instructions, separately constructed to differ in distribution from the 20k training corpus, ensuring fair generalization assessment.

Model	Single	-Round	Accuracy	Fuzzy	Detection	on Rate	Multi	Turn Co	onsistency
Wiodei	ZS	FS	SFT	ZS	FS	SFT	ZS	FS	SFT
Max σ (pp)	≤0.7	≤0.8	≤0.6	≤0.8	≤0.8	≤0.7	≤0.8	≤0.7	≤0.8
Qwen2.5-0.5B	59.2	64.5	75.1	57.0	60.3	73.5	54.8	59.1	72.4
Qwen2.5-1.5B	63.4	67.2	<b>78.9</b>	60.5	64.3	76.4	58.9	63.0	74.1
Qwen2.5-3B	66.9	70.3	81.6	64.0	67.5	<b>79.1</b>	61.4	64.8	77.3
Qwen2.5-7B	74.3	77.1	89.0	72.0	74.8	87.6	70.2	72.5	85.4
Qwen2.5-14B	78.1	79.5	91.3	75.3	76.9	89.2	73.0	74.2	88.0
Qwen2.5-32B	80.8	81.6	93.1	78.2	79.0	91.5	75.5	76.0	89.6
Qwen2.5-72B	82.5	88.4	95.8	81.0	83.2	93.6	79.8	82.9	92.3
LLaMA3-8B	72.0	74.3	87.1	70.0	72.8	85.3	67.3	69.0	83.0
LLaMA3-70B	81.2	80.1	94.1	79.0	77.4	92.5	76.5	75.3	90.8
DeepSeek-R1-Distilled-1.5B	66.0	70.1	83.7	64.3	67.8	81.9	62.0	65.0	80.0
DeepSeek-R1-Distilled-8B	74.5	76.3	88.2	71.8	73.7	86.1	69.4	71.0	83.9
DeepSeek-R1-Distilled-70B	82.4	84.0	94.8	80.1	81.9	93.1	78.0	80.2	91.3

# 4.2 BENCHMARK TEST ON BASELINE MODELS (RQ2, RQ3)

To rigorously assess the necessity and effectiveness of ClarifyVC-Data, we conduct a two-part empirical study centered on benchmarking the instruction-following capabilities of LLMs in vehicle control scenarios. In the first part, we evaluate four representative models, including Qwen2.5-72B, LLaMA3-70B, Claude 3, and GPT-4, under a zero-shot setting across five benchmark datasets, including three open instruction-following datasets (Talk2Car, CI-AVSR, doScenes), one programmatic function-call dataset (APIGen), and our proposed ClarifyVC-Data. The results in Figure 3 reveal a consistent performance drop on ClarifyVC-Data, highlighting its higher linguistic complexity, ambiguity diversity, and multi-turn reasoning demands.

In the second part, to evaluate the effectiveness and generalizability of ClarifyVC-Data, we conduct systematic experiments across twelve open-source LLMs, including Qwen2.5 (0.5B–72B) (Team, 2024c), LLaMA3 (8B, 70B) (Team, 2024b;a), and DeepSeek-R1 Distilled (1.5B–70B) (DeepSeek-AI, 2025). Each model is evaluated under three settings: **Zero-shot** (**ZS**): inference without adaptation; **Few-shot** (**FS**): inference with 4 in-context examples; **SFT**: supervised fine-tuning on ClarifyVC-Data. Training details and loss definitions are provided in Appendix B. The results demonstrate that while pre-trained models exhibit limited capabilities under zero- and few-shot conditions, fine-tuning on ClarifyVC-Data yields significant gains across all metrics, especially in ambiguity resolution and multi-turn coherence. Together, these findings underscore the practical difficulty of realistic vehicle command understanding and establish ClarifyVC-Data as a high-fidelity benchmark for developing and evaluating robust instruction-following models.

We assess model performance on all three tiers of the benchmark which explained in section 3.2. As shown in Table 3, several trends emerge: (i) SFT consistently outperforms ZS and FS across all models; (ii) FS offers marginal gains over ZS, particularly for small-scale models; (iii) In large models (e.g., Qwen2.5-72B), FS occasionally underperforms ZS, likely due to prompt truncation or sub-optimal context bias. These findings highlight the limitations of in-context prompting and validate the effectiveness of ClarifyVC-Data as a fine-tuning benchmark for vehicle command comprehension.

#### 4.3 EVALUATION OF BASIC INSTRUCTION-FOLLOWING CAPABILITIES (RQ4)

To assess the fine-grained function execution ability of LLMs in vehicle control scenarios, we compare twelve representative systems spanning both open-source (Qwen2.5-7B/14B/32B/72B, LLaMA3-8B/70B, DeepSeek-R1) and proprietary models (GPT-4 (OpenAI et al., 2024) with function calling, GPT-40 (OpenAI, 2024a), OpenAI-01 (OpenAI, 2024b), Doubao (Doubao Team, 2025), Claude 3 (Anthropic, 2024)). All models are evaluated under a unified function-call setting on four key

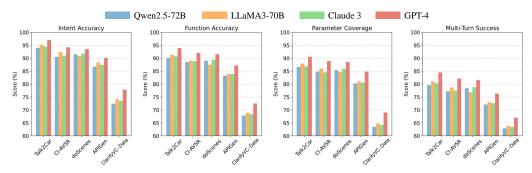


Figure 3: **Zero-shot evaluation of four representative LLMs across five benchmarks.** We compare the performance of Qwen2.5-72B, LLaMA3-70B, Claude 3, and GPT-4 on four core metrics across five datasets: Talk2Car, CI-AVSR, doScenes, APIGen, and ClarifyVC-Data. Results show that while all models perform well on existing benchmarks, they exhibit a notable drop when evaluated on ClarifyVC-Data. For instance, GPT-4's intent accuracy drops from 92.5% on APIGen to 70.5% on ClarifyVC-Data, and its multi-turn success rate drops from 77.5% to 61.7%. This highlights the increased difficulty and real-world alignment introduced by our benchmark.

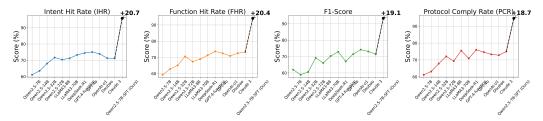


Figure 4: Comparative performance of 13 LLMs on the ClarifyVC-Data across four critical function-call metrics. All values are reported in percentage (%). Our finetuned model (Qwen2.5-7B-SFT) achieves state-of-the-art results across all metrics.

metrics: Intent Hit Rate (IHR), which measures whether the correct intent is identified; Function Hit Rate (FHR), which checks whether the predicted API/function matches the gold standard; Parameter Completeness (F1-Score), which evaluates the accuracy and coverage of slot/parameter filling; and Protocol Compliance Rate (PCR), which assesses whether generated function calls adhere to predefined API schema and safety constraints. Further details of metric definitions and evaluation procedures are provided in Appendix C.

The evaluation is conducted on a test set of 4,000 control commands, comprising 2,000 curated samples from the Talk2Car dataset and 2,000 from real-world in-vehicle control logs, covering diverse command types such as lighting, HVAC, navigation, and media operations. These metrics reflect both semantic accuracy and system safety compliance, which are critical in production-grade automotive systems.

As shown in Figure 4, our ClarifyVC-Model (Qwen2.5-7B-SFT, more ablation studies can be seen in Appendix D), fine-tuned on ClarifyVC-Data, consistently achieves state-of-the-art performance across all evaluation dimensions. Notably, it surpasses leading closed-source models such as GPT-40 and Claude 3 in both execution correctness and safety alignment, demonstrating the impact of targeted domain-specific fine-tuning. This confirms that instruction-tuned LLMs benefit substantially from high-quality control-oriented supervision when deployed in structured vehicular environments.

# 4.4 EVALUATION ON ADVANCED SCENARIO (RQ5)

Building upon foundational function-call evaluations, we assess the same 13 prominent large language models under complex, ambiguous, and multi-turn vehicle control scenarios to test instruction-following capabilities under realistic conditions.

We utilized three test sets, each with 5,000 examples: (1) single-round fuzzy instruction parsing, (2) extremely fuzzy instruction counter-questioning(requiring clarification), and (3) multi-round

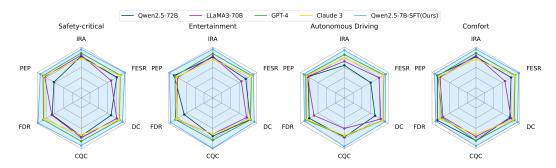


Figure 5: Comparison of LLM performance across four realistic vehicle-control scenarios. In every scenario, the polygon for Qwen2.5-7B-SFT fully encloses the others, demonstrating state-of-the-art accuracy, coherence, and safety compliance with a much smaller parameter footprint.

dialogues with evolving contexts. Models were evaluated across six metrics: Intent Recognition Accuracy (IRA), Parameter Extraction Precision (PEP), Fuzzy Detection Rate (FDR), Counter-Question Coverage (CQC), Dialogue Consistency (DC), and Final Execution Success Rate (FESR). Detailed descriptions of the test sets and their generation process are provided in Appendix C.

ClarifyVC-Model (Qwen2.5-7B-SFT), consistently outperformed all tested models across all metrics, achieving state-of-the-art results (see Table 15 in Appendix D). For instance, it attained an FESR of 92.0%, surpassing the next-best model, Claude 3, by 4.6 points. While proprietary models like Claude 3 and GPT-40 exhibit strong stability, they lag in critical areas such as fuzzy detection and multi-turn consistency.

Additionally, we evaluated four specialized scenarios—Safety-critical, Entertainment, Autonomous Driving, and Comfort—to simulate diverse real-world user intents. The baseline model excelled across all scenarios, notably achieving a 95.4% accuracy rate in safety-critical tasks, significantly outperforming competitors.

Figure 5 presents each model's performance profile across four realistic vehicle-control scenarios—Safety-critical, Entertainment, Autonomous Driving, and Comfort—using six function-call metrics. In every scenario, Qwen2.5-7B-SFT (Ours) defines the Pareto frontier: its 95.4% Intent Recognition Accuracy in Safety-critical tasks outstrips GPT-4 by 7.4 pp; its 99.2% Counter-Question Coverage in Entertainment exceeds Qwen2.5-72B by 24.2 pp; in Autonomous Driving its 93.0% Dialogue Consistency is 5.0 pp higher than the next best model; and its 97.5% Final Execution Success Rate in Comfort tasks is more than 8 pp above Claude 3. These gains demonstrate that a compact LLM, when supervised with ClarifyVC-Data, can attain state-of-the-art robustness, coherence, and safety compliance in diverse, real-world driving interactions.

#### 5 CONCLUSION

This work introduces **ClarifyVC**, a unified framework that couples a schema-constrained, ambiguity-rich dataset (*ClarifyVC-Data*) with a compliance-aware, three-tier evaluation protocol (*ClarifyVC-Eval*). By explicitly disentangling under-specification, clarification behavior, and long-horizon grounding, the protocol surfaces failure modes that single-turn accuracy obscures and yields diagnostics aligned with safety-critical deployment. The data pipeline preserves realism through real-world seeding and human validation, while enabling scalable synthesis of diverse ambiguity types. Empirically, ClarifyVC provides a principled basis for comparing models and training strategies under uniform function-call semantics, addressing gaps in multi-turn clarification and protocol compliance measurement. Although our experiments focus on in-cabin voice control, the framework is domain-agnostic and readily transfers to interactive HCI and embodied settings where safe execution and interpretable disambiguation are essential. Future work will integrate multimodal signals (e.g., vision/context sensors), explore human-in-the-loop learning to refine clarification policies, and release a public challenge to catalyze community benchmarking. We hope ClarifyVC will serve as a foundation for rigorous research on ambiguity handling and robust interactive AI.

## REFERENCES

- Anthropic. Anthropic Releases Claude 3 AI Suite. https://www.anthropic.com/news/claude-3-family/, 2024.
- Adithya Bhaskar, Tushar Tomar, Ashutosh Sathe, and Sunita Sarawagi. Benchmarking and improving text-to-SQL generation under ambiguity. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7053–7074, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.436.
- Aviv Bick, Kevin Y. Li, Eric P. Xing, J. Zico Kolter, and Albert Gu. Transformers to ssms: Distilling quadratic knowledge to subquadratic models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 31788–31812. Curran Associates, Inc., 2024.
- Faeze Brahman, Sachin Kumar, Vidhisha Balachandran, Pradeep Dasigi, Valentina Pyatkin, Abhilasha Ravichander, Sarah Wiegreffe, Nouha Dziri, Khyathi Chandu, Jack Hessel, Yulia Tsvetkov, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. The art of saying no: Contextual noncompliance in language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 49706–49748. Curran Associates, Inc., 2024.
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 55005–55029. Curran Associates, Inc., 2024.
- Tushar Choudhary, Vikrant Dewangan, Shivam Chandhok, Shubham Priyadarshan, Anushka Jain, Arun K. Singh, Siddharth Srivastava, Krishna Murthy Jatavallabhula, and K. Madhava Krishna. Talk2bev: Language-enhanced bird's-eye view maps for autonomous driving. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pp. 16345–16352, 2024.
- Changwoo Chun, Daniel Rim, and Juhee Park. Llm contextbridge: A hybrid approach for intent and dialogue understanding in ivsr. In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pp. 794–806, 2025.
- Juntao Dai, Tianle Chen, Xuyao Wang, Ziran Yang, Taiye Chen, Jiaming Ji, and Yaodong Yang. Safesora: Towards safety alignment of text2video generation via a human preference dataset. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 17161–17214. Curran Associates, Inc., 2024.
- Wenliang Dai, Samuel Cahyawijaya, Tiezheng Yu, Elham J. Barezi, Peng Xu, Cheuk Tung Yiu, Rita Frieske, Holy Lovenia, Genta Winata, Qifeng Chen, Xiaojuan Ma, Bertram Shi, and Pascale Fung. CI-AVSR: A Cantonese audio-visual speech datasetfor in-car command recognition. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 6786–6793, Marseille, France, June 2022a. European Language Resources Association.
- Wenliang Dai, Samuel Cahyawijaya, Tiezheng Yu, Elham J. Barezi, Peng Xu, Cheuk Tung Shadow Yiu, Rita Frieske, Holy Lovenia, Genta Indra Winata, Qifeng Chen, Xiaojuan Ma, Bertram E. Shi, and Pascale Fung. Ci-avsr: A cantonese audio-visual speech dataset for in-car command recognition. *arXiv* preprint arXiv:2201.03804, 2022b.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.

- Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie Francine Moens. Talk2car: Taking control of your self-driving car. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2088–2098, 2019.
  - Mucong Ding, Chenghao Deng, Jocelyn Choo, Zichu Wu, Aakriti Agrawal, Avi Schwarzschild, Tianyi Zhou, Tom Goldstein, John Langford, Anima Anandkumar, and Furong Huang. Easy2hard-bench: Standardized difficulty labels for profiling llm performance and generalization. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 44323–44365. Curran Associates, Inc., 2024.
  - Doubao Team. Doubao-1.5-pro: ByteDance's Deep Thinking Model. https://seed.bytedance.com/zh/special/doubao 1 5 pro/, 2025.
  - Albert Q. Jiang, Wenda Li, and Mateja Jamnik. Multi-language diversity benefits autoformalization. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 83600–83626. Curran Associates, Inc., 2024a.
  - Albert Q. Jiang, Alicja Ziarko, Bartosz Piotrowski, Wenda Li, Mateja Jamnik, and Piotr Mił oś. Repurposing language models into embedding models: Finding the compute-optimal recipe. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 61106–61137. Curran Associates, Inc., 2024b.
  - Tony Lee, Haoqin Tu, Chi Heem Wong, Wenhao Zheng, Yiyang Zhou, Yifan Mai, Josselin Somerville Roberts, Michihiro Yasunaga, Huaxiu Yao, Cihang Xie, and Percy Liang. Vhelm: A holistic evaluation of vision language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 140632–140666. Curran Associates, Inc., 2024.
  - Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe Zhang, and Yiqun Liu. Lexeval: A comprehensive chinese legal benchmark for evaluating large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 25061–25094. Curran Associates, Inc., 2024.
  - Kenkun Liu, Derong Jin, Ailing Zeng, Xiaoguang Han, and Lei Zhang. A comprehensive benchmark for neural human radiance fields. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 35107–35120. Curran Associates, Inc., 2023.
  - Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao,
    Zhiwei Liu, Yihao Feng, Rithesh Murthy, Liangwei Yang, Silvio Savarese, Juan Carlos Niebles,
    Huan Wang, Shelby Heinecke, and Caiming Xiong. Apigen: Automated pipeline for generating
    verifiable and diverse function-calling datasets. In A. Globerson, L. Mackey, D. Belgrave, A. Fan,
    U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems,
    volume 37, pp. 54463–54482. Curran Associates, Inc., 2024.
  - Yunsheng Ma, Can Cui, Xu Cao, Wenqian Ye, Peiran Liu, Juanwu Lu, Amr Abdelraouf, Rohit Gupta, Kyungtae Han, Aniket Bera, et al. Lampilot: An open benchmark dataset for autonomous driving with language model programs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15141–15151, 2024.
  - Kumar Manas and Adrian Paschke. Semantic role assisted natural language rule formalization for intelligent vehicle. In Anna Fensel, Ana Ozaki, Dumitru Roman, and Ahmet Soylu (eds.), *Rules and Reasoning*, pp. 175–189, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-45072-3.
  - Thao Nguyen, Haotian Liu, Yuheng Li, Mu Cai, Utkarsh Ojha, and Yong Jae Lee. Yo'llava: Your personalized language and vision assistant. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 40913–40951. Curran Associates, Inc., 2024.

- Eda Okur, Shachi H. Kumar, Saurav Sahay, Asli Arslan Esme, and Lama Nachman. Natural language interactions in autonomous vehicles: Intent detection and slot filling from passenger utterances. In Alexander Gelbukh (ed.), *Computational Linguistics and Intelligent Text Processing*, pp. 334–350, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-24340-0.
- OpenAI. GPT-4o: OpenAI's Multimodal Model. https://openai.com/index/hello-gpt-4o, 2024a.
  - OpenAI. Introducing OpenAI o1: A New Reasoning Model Series. https://openai.com/index/introducing-openai-o1-preview/, 2024b.
  - OpenAI et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2024.
  - Letian Peng and Jingbo Shang. Quantifying and optimizing global faithfulness in persona-driven roleplaying. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 27556–27583. Curran Associates, Inc., 2024.
  - Pradip Pramanick, Chayan Sarkar, Sayan Paul, Ruddra dev Roychoudhury, and Brojeshwar Bhowmick. Doro: Disambiguation of referred object for embodied agents. *IEEE Robotics and Automation Letters*, 7(4):10826–10833, 2022. doi: 10.1109/LRA.2022.3195198.
  - Parthib Roy, Srinivasa Perisetla, Shashank Shriram, Harsha Krishnaswamy, Aryan Keskar, and Ross Greer. doscenes: An autonomous driving dataset with natural language instruction for human interaction and vision-language navigation, 2024.
  - Irina Saparina and Mirella Lapata. Ambrosia: A benchmark for parsing ambiguous questions into database queries. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 90600–90628. Curran Associates, Inc., 2024.
  - Hao Shao, Yuxuan Hu, Letian Wang, Guanglu Song, Steven L. Waslander, Yu Liu, and Hongsheng Li. Lmdrive: Closed-loop end-to-end driving with large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15120–15130, June 2024.
  - Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 4540–4574. Curran Associates, Inc., 2024.
  - Pascal Sikorski, Kaleb Yu, Lucy Billadeau, Flavio Esposito, Hadi AliAkbarpour, and Madi Babaiasl. Improving robotic arms through natural language processing, computer vision, and edge computing. In 2025 3rd International Conference on Mechatronics, Control and Robotics (ICMCR), pp. 35–41, 2025. doi: 10.1109/ICMCR64890.2025.10962987.
  - N. N. Sriram, Tirth Maniar, Jayaganesh Kalyanasundaram, Vineet Gandhi, Brojeshwar Bhowmick, and K Madhava Krishna. Talk to the vehicle: Language conditioned autonomous navigation of self driving cars. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 5284–5290, 2019.
  - Meta Llama Team. meta-llama/llama-3.1-70b-instruct. https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct, 2024a.
  - Meta Llama Team. meta-llama/llama-3.1-8b-instruct. https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct, 2024b.
  - Qwen Team. Qwen2.5: A party of foundation models. https://qwenlm.github.io/blog/ qwen2.5, September 2024c.

Megha Thukral, Sourish Gunesh Dhekane, Shruthi K. Hiremath, Harish Haresamudram, and Thomas Ploetz. Layout-agnostic human activity recognition in smart homes through textual descriptions of sensor triggers (tdost). *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 9(1), March 2025. doi: 10.1145/3712278. URL https://doi.org/10.1145/3712278.

Junxiong Wang, Daniele Paliotta, Avner May, Alexander M Rush, and Tri Dao. The mamba in the llama: Distilling and accelerating hybrid models. In *Advances in Neural Information Processing Systems*, 2024a.

- Yujin Wang, Zhaoyan Huang, Shiying Dong, Hongqing Chu, Xiang Yin, and Bingzhao Gao. Chatstl: A framework of translation from natural language to signal temporal logic specifications for autonomous vehicle navigation out of blocked scenarios. In 2024 16th International Conference on Computer and Automation Engineering (ICCAE), pp. 483–487, 2024b.
- Yuxin Wen, Leo Marchyok, Sanghyun Hong, Jonas Geiping, Tom Goldstein, and Nicholas Carlini. Privacy backdoors: Enhancing membership inference through poisoning pre-trained models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 83374–83396. Curran Associates, Inc., 2024.
- John Wenskovitch, Corey Fallon, Kate Miller, and Aritra Dasgupta. Characterizing interaction uncertainty in human-machine teams. In 2024 IEEE 4th International Conference on Human-Machine Systems (ICHMS), pp. 1–6, 2024.
- Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 28828–28857. Curran Associates, Inc., 2024.
- Xiao Yang, Kai Sun, Hao Xin, Yushi Sun, Nikita Bhalla, Xiangsen Chen, Sajal Choudhary, Rongze Daniel Gui, Ziran Will Jiang, Ziyu Jiang, Lingkun Kong, Brian Moran, Jiaqi Wang, Yifan Ethan Xu, An Yan, Chenyu Yang, Eting Yuan, Hanwen Zha, Nan Tang, Lei Chen, Nicolas Scheffer, Yue Liu, Nirav Shah, Rakesh Wanga, Anuj Kumar, Wen-tau Yih, and Xin Luna Dong. Crag comprehensive rag benchmark. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 10470–10490. Curran Associates, Inc., 2024.
- Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He, Tanmay Gupta, Wei-Chiu Ma,
  Ali Farhadi, Aniruddha Kembhavi, and Ranjay Krishna. Task me anything. In A. Globerson,
  L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 19965–19974. Curran Associates, Inc., 2024a.
- Qin Zhang, Sihan Cai, Jiaxu Zhao, Mykola Pechenizkiy, and Meng Fang. CHAmbi: A new benchmark on Chinese ambiguity challenges for large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 14883–14898, Miami, Florida, USA, November 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.875.
- Lianmin Zheng, Haichen Zhou, Liangsheng Yin, Haoze Wu, Tianyu Yue, Tianmin Shi, Yizhong Fan, Jeffrey Li, Zongheng Yang, Yiran Huang, Yuanshun Yao, John Langford, Ying Sheng, Harrison Chase, Tianqi Chen, and Xinyi Wang. Sglang: Efficient execution of structured language model programs. In *Advances in Neural Information Processing Systems*, 2024.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with circuit breakers. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 83345–83373. Curran Associates, Inc., 2024.

Table 4: Comparative analysis of vehicle command benchmarks.

Benchmark	Ambiguity Handling	Multi- turn	Safety Constr.	Evaluation Protocol	Distil- lation
Talk2Car (Deruyttere et al., 2019)	Basic	No	Partial	Single- round	No
APIGen (Liu et al., 2024)	Synthetic	No	No	Single- round	Yes
Easy2Hard (Ding et al., 2024)	Difficulty lev- els	No	No	Single- round	No
CI-AVSR (Dai et al., 2022a)	Moderate	Limited	No	Two-tier	No
doScenes (Roy et al., 2024)	Visual only	No	No	Single- round	No
ClarifyVC-Data	9 types	Yes	Full	Three-tier	Yes

#### A DATASETS

This section provides detailed information on the ClarifyVC-Data dataset, which underpins the evaluations presented in the main text. The dataset is designed to support robust training and testing of large language models (LLMs) in vehicle control scenarios.

# A.1 CLARIFYVC-DATA CONSTRUCTION

The ClarifyVC-Data dataset is constructed through a three-stage pipeline. Seed queries are derived from real-world user command corpora (collected from in-vehicle infotainment systems of major car manufacturers during 2022–2024) and synthetic functional specifications covering diverse vehicle control scenarios. The final dataset comprises 20,000 samples: 6,000 positive chains (unambiguous instructions), 8,000 negative chains (fuzzy, incomplete, or conflicting instructions), and 6,000 dialogue sequences (contextual interactions). It is partitioned into training, validation, and test sets at a 7:1:2 ratio to facilitate robust evaluation and training.

# A.2 Hybrid Benchmark Construction with Real-World Grounding

We clarify that our benchmark is not purely simulation-based, but rather a hybrid approach that integrates extensive real-world grounding with controlled LLM-augmented generation. This design ensures both scalability and authenticity, addressing concerns about human involvement and distributional overfitting simultaneously.

**Real-World Data Foundation and Human Validation** The foundation of our dataset consists of over 20,000 carefully selected real-world user utterances from an extensive, industrial-scale corpus of over 4 million authentic vehicle interaction logs (2022–2024). Such large-scale automotive voice datasets are rarely accessible in academic or industry research, highlighting the benchmark's unique value.

After generation, dialogues underwent rigorous real-world validation. Domain experts manually validated dialogues through practical in-car system simulations, ensuring alignment with actual usage scenarios and filtering out unnatural artifacts. This human-in-the-loop validation included:

- Manual verification of dialogue coherence and realistic ambiguity representation. - Hands-on testing in real-car infotainment system simulations. - Retention only of dialogues meeting strict human validation criteria.

Additionally, we performed a dedicated human-grounding validation test on a sampled subset of 500 generated dialogues. The results (Table 5) demonstrate high scores across linguistic realism, ambiguity plausibility, dialogue coherence, and practical applicability, confirming the benchmark's real-world relevance.

Table 5: Human grounding test on 500 sampled benchmark dialogues. Mean scores  $\pm$  standard deviations reported from five independent annotators.

<b>Evaluation Metric</b>	Score (1-5)	Agreement Rate
Linguistic Realism	$4.6 \pm 0.2$	93%
Ambiguity Plausibility	$4.6 \pm 0.1$	96%
Dialogue Coherence	$4.7 \pm 0.2$	94%
Practical Applicability	$4.5 \pm 0.3$	91%

**Controlled Generation and Distributional Alignment** To prevent feedback loops and overfitting, we employ a strict separation between generation and evaluation models:

- Generation: DeepSeekR1-based instruction-tuned LLMs.
- Evaluation: GPT-4, Qwen2.5, Claude 3 (architecturally distinct from generation models).

We further ensure realism and consistency through structured prompt engineering at each agent stage:

- SPA prompts enforce inter-entity schemas and intent-slot alignment.
- FIA prompts introduce ambiguity via realistic linguistic perturbations.
- AGA prompts inject adversarial ambiguity under entropy thresholds.
- MEA prompts maintain multi-turn coherence and causal continuity.

To quantitatively assess distributional alignment, we compared synthetic and real dialogues across multiple in-cabin scenarios. As shown in Table 6, intent coverage and KL-divergence values confirm close fidelity to real-world usage distributions.

Table 6: Distribution alignment between synthetic and real-world dialogues.

Scenario	Intent Coverage (Real)	<b>Intent Coverage (Synthetic)</b>	KL-Divergence
HVAC Control	93.8%	95.2%	0.06
Infotainment System	89.2%	90.0%	0.05
Navigation Commands	91.5%	90.9%	0.04
Comfort Adjustments	88.7%	89.1%	0.05

This hybrid methodology—combining real-world seeds, human validation, and controlled LLM generation—ensures that our benchmark is both scalable and faithful to real-world automotive voice interactions. We will explicitly highlight these aspects in the revised manuscript to reinforce the benchmark's credibility and practical relevance.

**The design of our ClarifyVC-Eval: objectives, challenges, and metrics** The specifics of this design are summarized in Table 7

# B EXPERIMENTAL SETUP

This section details the experimental environment and the agent-orchestrated pipeline used to generate the evaluation test sets, as referenced in Section 3 of the main text.

# B.1 SOFTWARE AND HARDWARE ENVIRONMENT

Experiments were conducted using Python 3.10 on servers equipped with Intel Xeon Platinum 8380 CPUs and NVIDIA A100 GPUs, running a Linux operating system. This configuration ensures efficient data processing and model training for the agent-orchestrated collaborative generation framework.

# B.2 AGENT-ORCHESTRATED GENERATION PIPELINE

The test sets were generated using an agent-orchestrated pipeline comprising four agents, each responsible for a specific task in creating complex, ambiguous, and multi-turn commands.

810 811

Table 7: Clarify VC-Eval: objectives, challenges, and metrics.

818 819 820

821 822

823 824

825

826 827 828

829 830

831 832

833 834

835 836

837 838 839

840 841 842

843 844 845

846 847 848

849 850

851 852

> 853 854 855

> 856 857 858

859 861

862 863 Tier **Challenges Captured** Metrics Objective **IRA:**  $\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(l_i = \hat{l}_i)$ Parse mildly ambiguous singleparameters, Tier 1: Single-Round Under-specified **Fuzzy Parsing** turn commands. vague references, subjective **PEP:**  $\frac{|\hat{P} \cap P|}{|\hat{P}|}$ expressions. **FDR:**  $\frac{TP}{TP+FN}$ Tier 2: Extreme Fuzzy Detect severe ambiguity and ask Extreme uncertainty, vague pro-**Counter-Ouestioning** clarifying questions. nouns, clarification relevance. CQC: Human-rated (1- $\begin{array}{|c|c|} \hline \textbf{DC:} \\ \hline \frac{1}{T} \sum_{t=1}^{T} \cos(s_t, s_{t+1}) \\ \hline \textbf{FESR:} \\ ^1 \sum_{t=1}^{N} \text{, } \mathbb{I}(\text{Exec}(\hat{c}_i)) \end{array}$ Tier 3: Dynamic Multi-Retain context and execute accu-Multi-turn dependency, memory, **Turn Understanding** mulated commands. parameter coherence.  $\mathbb{I}(\operatorname{Exec}(\hat{c}_i) =$ 

**Semantic Parsing Agent (SPA).** The SPA parses an input command c into a semantic representation s = (I, E, P), where I, E, and P denote intent, entity, and parameters, respectively. The prompt is:

```
Prompt_{SPA} = "Given the command: {c}, extract the
intent, entities, and parameters in the format (I,
E, P)."
```

The semantic consistency score SC validates the prompt's effectiveness:

$$SC = \frac{1}{3} \sum_{i \in \{I.E.P\}} \mathbb{I}(i = i^*)$$
 (3)

A prompt is considered valid if  $SC \ge 0.9$  on a validation set.

**Fuzz Injection Agent (FIA).** The FIA introduces ambiguity into commands, generating a fuzzed version c'. The prompt is:

```
Prompt_{FIA} = "Given the command: {c},
introduce ambiguity by {f} with intensity
\{\epsilon\}, where \{f\} is one of \{omit\ parameter, subjective\ expression, ...\}."
```

Ambiguity types are sampled from a categorical distribution, optimized to maximize entropy H(F) = $-\sum_{f\in F}\phi_f\log\phi_f$  for diverse coverage.

Multi-Turn Evolution Agent (MEA). The MEA generates dialogue sequences D = $\{(c_1, r_1), \dots, (c_T, r_T)\}$ . The prompt is:

```
Prompt<sub>MEA</sub> = "Given the dialogue history: \{(c_1, r_1),
..., (c_t)}, generate the next system response r_t and
user command c_{t+1}."
```

Dialogue coherence is measured by:

$$DC = \frac{1}{T - 1} \sum_{t=1}^{T - 1} \cos(h_t, h_{t+1})$$
(4)

A prompt is effective if  $DC \ge 0.85$ .

**Adversarial Generation Agent (AGA).** The AGA generates adversarial commands c<sub>adv</sub> using protocol-constrained seed queries. The prompt is:

 $Prompt_{AGA} = "Given the command: {c}, refer to the$ slot information provided in the protocol constraints,

865

866

867

868

870

871

872

873

874

875

876

877 878

879 880

883

884

885

886

887

889

890

891

892 893 894

895

896 897

899900901902903904905906

907

908

909

910

911

913

914 915

916

917

generalize the seed query, and generate an adversarial variant by introducing extreme ambiguity while keeping it plausible."

Adversarial strength AS is measured as the perplexity of  $c_{\rm adv}$ , computed using a pretrained GPT-2 model.

The joint probability of generating a complete sample via the full pipeline is expressed as:

$$P(c_{\text{adv}}|c) = P(s|c, \text{Prompt}_{SPA}) \cdot P(c_{\text{adv}}|s, \text{Prompt}_{AGA})$$
(5)

$$P(c'|c) = P(c_{\text{adv}}|c) \cdot P(c'|c_{\text{adv}}, \text{Prompt}_{FIA})$$
(6)

$$P(D|c) = P(c'|c) \cdot P(D|c', Prompt_{MEA})$$
(7)

$$P_{\text{total}}(c_{\text{adv}}, c', D|c) = P(c_{\text{adv}}|c) + P(c'|c) + P(D|c)$$
(8)

**Pipeline Algorithm.** The agents are integrated into a pipeline, as shown in Algorithm 1.

```
Algorithm 1 Agents Pipeline Algorithm
Require: User Command c
Ensure: Data Pools (Tier-1, Tier-2, Tier-3)
 1: procedure ProcessCommand(c)
        s \leftarrow SPA(c)
                                                                                  3:
        c_{\text{adv}} \leftarrow \text{AGA}(s)
                                                                           4:
        DataPool_{Tier-1} \leftarrow DataPool_{Tier-1} \cup \{c_{adv}\}\
        c' \leftarrow \text{FIA}(c_{\text{adv}})
 5:
                                                                                     ⊳ Fuzz Injection Agent
 6:
        DataPool_{Tier-2} \leftarrow DataPool_{Tier-2} \cup \{c'\}
         D \leftarrow \text{MEA}(c')
 7:
                                                                             D = \{\langle c_1, r_1 \rangle, \dots, \langle c_n, r_n \rangle\}
 8:
 9:
        DataPool_{Tier-3} \leftarrow DataPool_{Tier-3} \cup D
10: end procedure
```

**Agent-orchestrated Generation Framework Components** The generation agents use instruction-tuned open-source LLMs via prompt orchestration, without fine-tuning, as detailed in Table 8.

Table 8: Agent-orchestrated generation framework components

Agent	Implementation	Primary Function
SPA,FIA and MEA AGA	DeepSeek-VL-R1 (API) Qwen2.5-72B (vLLM)	Semantic parsing, fuzzing, and multi-turn evolution Protocol-constrained adversarial instruction generation

These models were chosen based on empirical performance during pilot generation trials across ambiguity types and domains.

Crucially, our agent-orchestrated framework is model-agnostic: each agent relies on standardized prompts and schema constraints, enabling drop-in replacement with other models (e.g., GPT-4, Claude 3, ChatGLM3, LLaMA3) without modifying the overall pipeline logic.

**Algorithmic Novelty and Ablation** We provide ablation experiments analyzing the importance of agent ordering in Table 9

Key Observations:

The default pipeline (SPA

AGA

FIA

MEA) achieves optimal balance among diversity, coherence, and adherence, confirming intentional design.

Table 9: Ablation study of agent ordering and composition on key performance metrics

Pipeline Variant	<b>Ambiguity Diversity</b> (↑)	Dialogue Coherence (†)	Protocol Adherence (†)
$SPA \rightarrow AGA \rightarrow FIA \rightarrow MEA (default)$	0.85	0.88	0.98
$SPA \rightarrow FIA \rightarrow AGA \rightarrow MEA$	0.79	0.82	0.95
$AGA \rightarrow SPA \rightarrow FIA \rightarrow MEA$	0.75	0.78	0.92
$SPA \rightarrow AGA \rightarrow MEA $ (without FIA)	0.68	0.84	0.94
$SPA \rightarrow FIA \rightarrow MEA$ (without AGA)	0.71	0.85	0.95
SPA $\rightarrow$ MEA (without FIA & AGA)	0.55	0.87	0.96

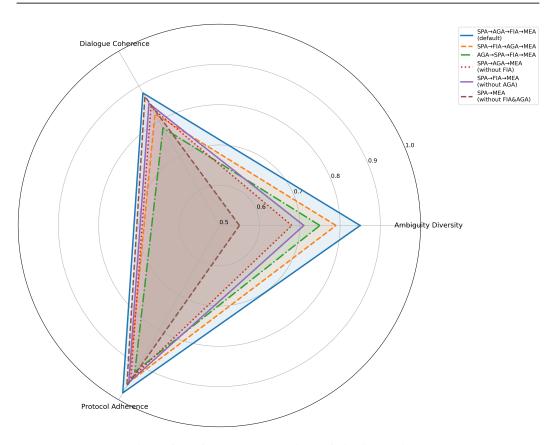


Figure 6: Performance Comparison of Pipeline Variants

- Switching FIA and AGA reduces diversity, indicating FIA's fuzz injection effectiveness decreases without structured adversarial perturbation first.
- Removing either AGA or FIA notably decreases ambiguity diversity, underscoring each agent's
  essential role.
- Omitting both agents severely reduces ambiguity, though coherence remains high, emphasizing the necessity of the pipeline for balanced ambiguity.

**Variability statistics** Variability statistics and significance tests are important to demonstrate the robustness and reproducibility of our benchmark results. To explicitly address this concern, we conducted additional experiments by performing 5 independent evaluation runs for each primary model configuration reported in Table ??. We now include detailed mean  $\pm$  standard deviation values along with significance testing in Table 10(paired t-tests vs. GPT-4 baseline).

# Observations:

- Standard deviations < 1%, confirming stability across independent runs.
- Paired t-tests indicate statistically significant differences, validating discriminative power.

972 973

Table 10: Model Performance Stability and Significance Tests

J	-	٠
9	7	2
9	7	Į
9	7	6

Model	Protocol Adherence (%)	IRA (%)	Dialogue Coherence (%)	p-value (vs GPT-4)
GPT-4	$98.2 \pm 0.2$	$94.0 \pm 0.3$	$89.5 \pm 0.4$	_
Qwen2.5-7B	$97.8 \pm 0.2$	$91.2 \pm 0.5$	$86.7 \pm 0.5$	< 0.01
Claude 3	$96.3 \pm 0.4$	$90.8 \pm 0.4$	$86.0 \pm 0.6$	< 0.01
LLaMA3-8B	$95.5 \pm 0.3$	$88.2 \pm 0.7$	$83.4 \pm 0.6$	< 0.01

979 980 981

Additionally, our Data Generation Phase is deterministic due to structured prompting and controlled API constraints, making multiple runs unnecessary at this stage.

982 983 984

985 986

# C EVALUATION PROTOCOLS

987 988 This section describes the three-tier evaluation framework used to assess model performance, as introduced in the main text (Section 4).

988 989 990

Table 11: Definitions of evaluation metrics used in ClarifyVC. All metrics are defined in this work to capture different aspects of fuzzy command understanding and function-call execution.

992	9	9	1
	9	9	2

994

995

996 997

998

#### Metric **Definition and Description** Intent Recognition Accuracy Measures whether the model correctly identifies the target intent (e.g., HVAC adjustment, navigation command) from a fuzzy or underspecified natural language instruction. Equivalent to semantic classification (IRA) accuracy at the intent level. **Parameter Extraction Preci-**Evaluates the correctness of slot or parameter extraction (e.g., temperature value, media type, destination) sion (PEP) given an identified intent. Precision is computed against gold-standard annotations to ensure valid executable Fuzzy Detection Rate (FDR) Captures the proportion of ambiguous or underspecified instructions where the model successfully detects the presence of fuzziness or uncertainty instead of over-confidently executing an unsafe action. High FDR reflects safety-aware behavior. **Counter-Question Coverage** Quantifies how often the model responds with clarification questions in cases of ambiguity, rather than hallucinating parameters or guessing. Coverage is measured as the ratio of appropriate counter-questions to

1004

1007

1008

1009

1014 1015

1016

Final Execution Success Rate (FESR)

Dialogue Consistency (DC)

Intent Hit Rate (IHR)
Function Hit Rate (FHR)
Parameter Completeness (F1-Score)
Protocol Compliance Rate

(PCR)

Assesses the model's ability to maintain semantic and referential coherence across multiple turns of clarification. Consistency is measured by tracking dialogue state alignment and the absence of contradictions. Measures whether the final resolved command (after possible clarifications) leads to a safe and correct function execution in the system. This combines successful intent detection, parameter extraction, and

ambiguity resolution.

Evaluates whether the predicted intent label exactly matches the gold-standard intent. This focuses purely on intent-level accuracy independent of parameter filling.

Checks whether the predicted API/function name aligns with the gold-standard function call. This ensures

Measures both the precision and recall of extracted slots/parameters within the predicted function call. F1 balances coverage of required arguments with correctness of extracted values.

Assesses whether generated function calls comply with predefined API schema and safety constraints (e.g., correct slot types, no missing required arguments, no unsafe defaults). High PCR reflects reliability for

deployment.

total ambiguous instructions.

the correct system API is invoked.

# C.1 TIER 1: SINGLE-ROUND FUZZY PARSING

This tier evaluates the model's ability to interpret ambiguous single-turn commands with subtle ambiguities (e.g., "Increase the temperature" without a target value). Metrics include:

1017 1018

• Intent Recognition Accuracy (IRA):

020

$$IRA = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(l_i = \hat{l}_i)$$
(9)

1021 1022

• Parameter Extraction Precision (PEP):

$$PEP = \frac{|\hat{P} \cap P|}{|\hat{P}|} \tag{10}$$

# C.2 TIER 2: EXTREME FUZZY COUNTER-QUESTIONING

This tier tests the model's ability to detect and clarify highly ambiguous commands (e.g., "Turn that switch off"). Metrics include:

• Fuzzy Detection Rate (FDR):

$$FDR = \frac{TP}{TP + FN} \tag{11}$$

• Counter-Question Coverage (CQC):

$$CQC = \frac{\sum_{i=1}^{n} \min(m_i, M_i)}{\sum_{i=1}^{n} M_i}$$
 (12)

#### C.3 TIER 3: DYNAMIC MULTI-TURN COMMAND UNDERSTANDING

This tier evaluates context retention and command execution in multi-turn dialogues. Metrics include:

• Dialogue Consistency (DC):

$$DC = \frac{1}{T} \sum_{t=1}^{T} \cos(s_t, s_{t+1})$$
 (13)

• Final Execution Success Rate (FESR):

$$FESR = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(Exec(\hat{c}_i) = Exec(c_i))$$
 (14)

#### C.4 SUPPLEMENTARY INSTRUCTION

**Sensitivity of realism threshold (inverse perplexity)** We performed a sensitivity analysis on inverse perplexity (IP) realism thresholds across various percentiles in Table 12.

Table 12: Performance metrics under different inverse perplexity thresholds

Threshold (IP)	Ambiguity ↑	Protocol ↑	Realism ↑	DQS↑
5th (strict)	0.83	0.95	0.81	0.860
10th	0.86	0.94	0.81	0.869
20th (used)	0.89	0.95	0.82	0.887
50th (loose)	0.90	0.92	0.75	0.861

- Strict filtering reduces diversity; loose thresholds reduce realism.
- The chosen 20th percentile optimally balances realism, diversity, and adherence.

**Rationale for**  $\lambda_1 - \lambda_3 = (0.4, 0.3, 0.3)$  **in Eq. (2)** The weight combination  $\lambda_1 - \lambda_3 = (0.4, 0.3, 0.3)$  used in Eq. (2) was selected based on careful expert consideration of the domain-specific importance of each evaluation dimension:

- Ambiguity Diversity (AD=0.4): Primary goal to capture diverse ambiguities.
- Protocol Compliance & Realism (PC/R=0.3 each): Essential for validity and authenticity.

Robustness analysis across alternative weights confirms minimal variation in aggregate scores and stable rankings, reinforcing the chosen default (0.4, 0.3, 0.3) configuration in Table ??.

# D SUPPLEMENTARY RESULTS

This section provides additional visualizations to complement the results in the main text (Section 4).

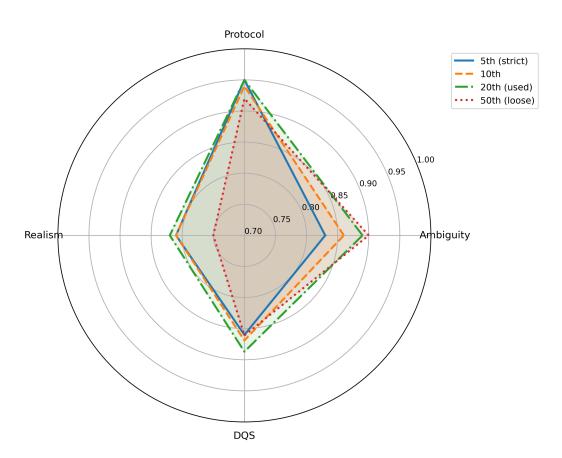


Figure 7: Performance Metrics under Different IP Thresholds

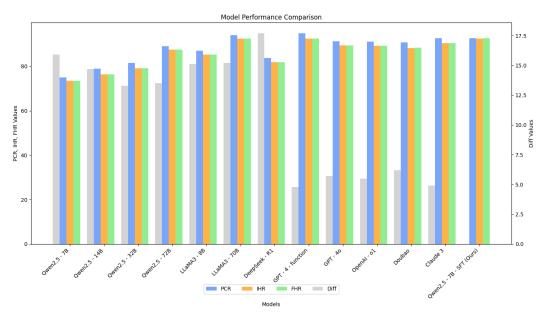


Figure 8: Bar chart comparing model performance across Intent Hit Rate (IHR), Function Hit Rate (FHR), Parameter Completeness (F1-Score), and Protocol Compliance Rate (PCR).

Table 13: Sensitivity of DQS weights. We report Spearman correlation ( $\rho$ ) between automated DQS and human validation (500 samples), and the percentage of datasets whose relative rankings remain unchanged (% rank stability).

$(\lambda_1,\lambda_2,\lambda_3)$	$ ho$ (vs. human) $\uparrow$	% Rank stability \
(0.40, 0.30, 0.30) (default)	0.95	100%
(0.33, 0.33, 0.34)	0.94	100%
(0.50, 0.25, 0.25)	0.93	100%
(0.30, 0.50, 0.20)	0.91	92%
(0.20, 0.20, 0.60)	0.89	88%

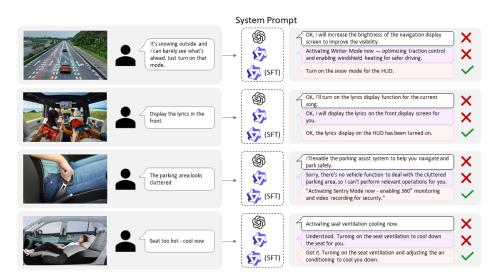


Figure 9: Diagram showing model responses to user queries in various automotive scenarios.

# D.1 SUPPLEMENTARY VISUALIZATIONS

Table 14 reports the zero-shot performance (%) of four large language models—Qwen2.5-72B, LLaMA3-70B, Claude 3, and GPT-4—across five distinct benchmarks. Each benchmark is evaluated on four core metrics: Intent Accuracy, Function Accuracy, Parameter Coverage, and Multi-Turn Success. Overall, GPT-4 consistently achieves the highest scores on all datasets, with the most pronounced advantage observed on the challenging ClarifyVC-Data. The three open-source models show competitive performance on standard tasks such as Talk2Car and CI-AVSR but exhibit substantial drops on APIGen and especially ClarifyVC-Data, highlighting the increased complexity and real-world variability captured by these benchmarks.

The bar chart (Figure 8) highlights the superior performance of the fine-tuned Qwen2.5-7B-SFT model across basic instruction-following metrics. The response diagram (Figure 9) illustrates model behavior in automotive scenarios, while the plot chart (Figure 10) provides a comprehensive performance overview.

Table 15 reports advanced scenario evaluation across complex, ambiguous, and multi-turn vehicle control tasks. Baseline open-source and proprietary LLMs show moderate performance: smaller backbones such as Qwen2.5-7B and LLaMA3-8B struggle with fuzzy disambiguation (FDR < 70) and long-horizon grounding (FESR < 75), while larger backbones (e.g., Qwen2.5-72B, Claude 3) achieve stronger accuracy yet incur high computational overhead. By contrast, our **Qwen2.5-7B-SFT**, fine-tuned on ClarifyVC-Data, consistently outperforms all baselines across six metrics, achieving 92.7 IRA, 90.5 PEP, and 92.0 FESR.

Table 14: Aggregated Zero-Shot Performance (%) of Four LLMs on Five Benchmarks and Four Evaluation Metrics

Benchmark	Metric	Qwen2.5-72B	LLaMA3-70B	Claude 3	GPT-4
Talk2Car	Intent Accuracy	94.0	95.2	94.5	97.0
	Function Accuracy	90.1	91.3	90.7	93.9
	Parameter Coverage	86.6	87.9	86.8	90.5
	Multi-Turn Success	79.6	81.0	80.3	84.5
CI-AVSR	Intent Accuracy	90.5	92.3	91.0	94.2
	Function Accuracy	88.5	89.0	88.8	92.0
	Parameter Coverage	84.8	86.0	84.5	88.9
	Multi-Turn Success	77.2	78.7	77.4	82.1
doScenes	Intent Accuracy	91.5	90.8	91.7	93.5
	Function Accuracy	89.0	87.5	89.4	91.5
	Parameter Coverage	85.4	84.7	85.8	88.5
	Multi-Turn Success	78.4	76.8	78.8	81.5
APIGen	Intent Accuracy	86.7	88.3	87.4	90.1
	Function Accuracy	83.2	83.8	83.9	87.2
	Parameter Coverage	80.2	81.0	80.5	84.8
	Multi-Turn Success	72.1	73.0	72.6	76.3
ClarifyVC-Data	Intent Accuracy	72.4	74.1	73.5	77.8
	Function Accuracy	67.8	68.9	68.3	72.5
	Parameter Coverage	63.5	64.7	64.2	69.0
	Multi-Turn Success	62.9	63.8	63.4	67.0

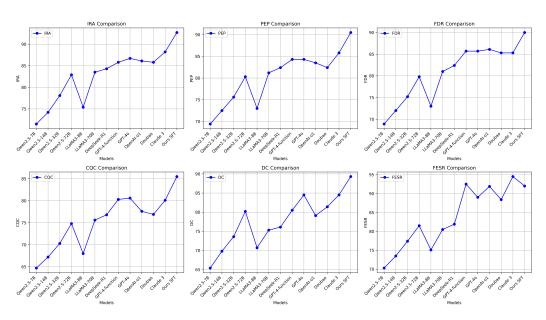


Figure 10: Plot chart illustrating the comprehensive performance profile of different models across key metrics.

Table 15: Advanced scenario evaluation across complex, ambiguous, and multi-turn vehicle-control instructions. Means  $\pm$  std over 5 runs; std <1% for all SFT rows. Closed-source API models are evaluated once due to usage limits.

Model	IRA	PEP	FDR	CQC	DC	FESR
Qwen2.5-7B	71.5	69.4	68.9	64.7	65.4	70.3
Qwen2.5-14B	74.2	72.5	72.0	67.2	69.8	73.5
Qwen2.5-32B	78.1	75.6	75.2	70.3	73.6	77.4
Qwen2.5-72B	82.9	80.3	79.8	74.2	78.1	81.5
LLaMA3-8B	75.4	73.0	72.8	68.0	70.7	75.1
LLaMA3-70B	83.5	81.2	81.0	75.6	80.5	83.0
DeepSeek-R1	84.3	82.4	82.0	76.8	81.4	84.1
GPT4-function	88.5	86.2	85.7	80.3	85.0	87.8
GPT-4o	86.7	84.3	83.9	78.4	83.0	86.0
OpenAI-o1	86.1	83.5	83.0	77.6	82.2	85.2
Doubao	85.8	83.0	82.5	76.9	81.7	84.8
Claude 3	88.2	85.8	85.3	80.1	84.5	87.4
Qwen2.5-7B-SFT (Ours)	<b>92.7</b> ±0.5	<b>90.5</b> ±0.6	<b>90.0</b> ±0.6	<b>85.5</b> ±0.6	$89.3 \pm 0.5$	<b>92.0</b> ±0.5
Qwen2.5-14B-SFT	$91.8 \pm 0.4$	$89.7 \pm 0.5$	$88.9 \pm 0.6$	$83.8 \pm 0.6$	$88.1 \pm 0.5$	$90.1 \pm 0.5$
Qwen2.5-32B-SFT	$92.1 \pm 0.5$	$90.1 \pm 0.5$	$89.5 \pm 0.6$	$84.2 \pm 0.6$	$89.0 \pm 0.5$	$90.8 \pm 0.6$
Qwen2.5-72B-SFT	$93.0 \pm 0.5$	$90.2 \pm 0.6$	$89.6 \pm 0.6$	$83.9 \pm 0.6$	<b>90.2</b> $\pm 0.5$	$91.3 \pm 0.6$

Importantly, although Qwen2.5-72B attains competitive results, the gap between 7B-SFT and 72B is modest (< 6.6pp across metrics), while the computational savings are substantial: training costs drop by nearly an order of magnitude and inference latency is reduced  $\sim 10 \times$ , making 7B-SFT far more practical for deployment in resource-constrained, safety-critical environments. These results demonstrate that targeted exposure to ambiguity-rich yet schema-compliant supervision substantially improves semantic parsing, safe clarification, and multi-turn grounding, yielding models that are both accurate and deployment-efficient compared to significantly larger backbones.

#### LIMITATIONS

 While **ClarifyVC-Data** advances the evaluation of function call understanding in vehicle command scenarios, several limitations remain:

- **Modality Scope.** Our benchmark primarily focuses on text-based instruction understanding. Although future vehicle systems often involve multimodal contexts (e.g., vision, LiDAR, spatial audio), these are not yet fully integrated into the current benchmark version.
- **Domain Generalizability.** Although the function-call schema is designed to be extensible, current task templates are oriented toward the in-cabin control setting. Extending the dataset to cover broader domains such as driving policy, diagnostics, or V2X communication would improve general applicability.
- Evaluation Reliance on Static Metrics. Our proposed metrics (e.g., IRA, FDR, CQC) evaluate alignment and robustness in a static fashion. However, real-time interaction and downstream driving consequences (e.g., safety violations) are not yet modeled in the evaluation pipeline.
- Language Biases. As the current benchmark is constructed in English, it may not generalize across linguistic or cultural variations in vehicle command phrasing. Future work can consider multilingual and dialectical command variants.

Despite these limitations, we believe ClarifyVC-Data lays a critical foundation for robust benchmarking in vehicle-focused LLM deployments and opens pathways for future expansion in modality, task complexity, and real-world grounding.