

# VLM-Driven Proactive Escalation for Tele-operations in Autonomous Vehicles

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Recent advances in autonomous driving enable vehicles*  
002 *to operate in increasingly complex environments, yet safe*  
003 *deployment still requires timely human intervention when*  
004 *the autonomous system encounters ambiguous or safety-*  
005 *critical situations. Existing remote driving systems rely*  
006 *on heuristic or confidence-based triggers to request assis-*  
007 *tance, which lack semantic understanding of driving con-*  
008 *text and often cause over- or under-escalation. We propose*  
009 *a vision–language reasoning-driven escalation framework*  
010 *that enables autonomous vehicles to request remote driv-*  
011 *ing assistance based on context-aware semantic reasoning*  
012 *rather than numerical uncertainty. Our approach uses a*  
013 *vision–language model (VLM) aligned using Direct Prefer-*  
014 *ence Optimization (DPO) on pairwise operator preferences*  
015 *to calibrate escalation decisions and reduce unnecessary*  
016 *or missed escalations. The model interprets visual obser-*  
017 *vations and driving context to generate structured explana-*  
018 *tions that help operators quickly understand and respond*  
019 *to situations. We evaluate the approach using NVIDIA*  
020 *Cosmos-2 across diverse uncertainty scenarios, including*  
021 *intersection negotiation and dynamic obstacle ambiguity.*  
022 *Results show improved intervention efficiency, reduced un-*  
023 *necessary operator engagement, and enhanced safety com-*  
024 *pared to heuristic triggering, highlighting a scalable path-*  
025 *way toward reliable human-in-the-loop autonomous driv-*  
026 *ing.*

## 027 1. Introduction

028 Level 4 automated vehicles (AVs) have been progressively  
029 deployed in both commercial and research settings in recent  
030 years. However, despite significant advances in perception,  
031 planning, and control, AVs still encounter situations that  
032 cannot be reliably resolved by automation alone. These sit-  
033 uations often lead to disengagements, where the automated  
034 driving system (ADS) must be deactivated, preventing mis-  
035 sion completion or safe continuation of the driving task.  
036 Such failures typically arise in long-tail scenarios, com-  
037 plex social interactions with other road users, or context-

dependent decision-making that exceed the capabilities of 038  
rule-based or confidence-driven autonomy [24, 34]. As a 039  
result, remote driving and teleoperation have emerged as 040  
critical human-in-the-loop mechanisms for ensuring safety, 041  
reliability, and operational continuity in real-world AV de- 042  
ployments [1, 10, 11, 21]. Beyond serving as a fallback 043  
control interface, remote operators can interpret semanti- 044  
cally ambiguous traffic situations and provide guidance or 045  
direct control when autonomy becomes uncertain. In prac- 046  
tical deployments such as fleet-based autonomous mobility 047  
services, logistics vehicles, and supervised AV operations, 048  
timely human intervention becomes particularly important 049  
when a limited number of operators must oversee multiple 050  
vehicles under constraints of bandwidth, latency, and atten- 051  
tion [5, 7, 12]. 052

In large-scale deployments of autonomous vehicles, par- 053  
ticularly in robotaxi and autonomous mobility services, a 054  
fundamental challenge lies in determining when an auto- 055  
mated vehicle should request remote operator assistance 056  
and ensuring that such requests are accompanied by clear 057  
and interpretable contextual information. Autonomous 058  
driving remains inherently difficult due to dynamic inter- 059  
actions with other road users, evolving traffic patterns, en- 060  
vironmental uncertainty, and rare long-tail events that can- 061  
not be reliably handled by pre-programmed behaviors or 062  
statistical confidence estimates. Conventional autonomous 063  
driving stacks, typically composed of perception, predic- 064  
tion, and planning modules, are prone to cascading errors 065  
and limited robustness when confronted with unfamiliar or 066  
socially complex traffic situations [15, 32]. As a result, au- 067  
tomated driving performance may degrade for multiple rea- 068  
sons, broadly including perception ambiguity, planning un- 069  
certainty, and control-level difficulties, each requiring dif- 070  
ferent levels of urgency and types of remote assistance. 071  
In fleet-scale robotaxi operations, requesting intervention 072  
too frequently can overwhelm remote operators and reduce 073  
system efficiency, whereas delayed escalation may lead to 074  
safety risks or service disruptions. Existing remote driving 075  
and tele-assistance frameworks often rely on heuristic rules, 076  
confidence thresholds, or handcrafted triggers to determine 077  
when a vehicle should request human intervention [11, 28]. 078

079 However, these mechanisms lack semantic understanding of  
080 the surrounding driving context and are therefore limited in  
081 their ability to prioritize assistance requests or communicate  
082 the underlying reason for escalation [25].

083 Motivated by these challenges and emerging needs, to  
084 the best of our knowledge, we are the first to introduce a  
085 VLM-based reasoning framework for proactive escalation  
086 of remote driving requests. Our key contributions are sum-  
087 marized as follows:

088 **Reasoning-Driven Escalation Framework:** We introduce  
089 a vision–language reasoning-based escalation framework  
090 that enables autonomous vehicles to proactively request re-  
091 mote driving assistance using context-aware semantic un-  
092 derstanding, moving beyond heuristic and confidence-based  
093 triggers.

094 **Preference-Aligned Decision Calibration via DPO:** We  
095 develop a task-conditioned post-training strategy using Di-  
096 rect Preference Optimization (DPO) to align the VLM with  
097 operator preferences, improving escalation accuracy and re-  
098 ducing both over- and under-intervention.

099 **Diverse Scenario Generation and Evaluation:** We ad-  
100 dress dataset limitations by generating and evaluating di-  
101 verse driving scenarios using NVIDIA Cosmos-2, enabling  
102 comprehensive assessment under varied uncertainty condi-  
103 tions.

## 104 2. Related Work

105 **Vision-Language Models in Autonomous Driving.** The  
106 integration of vision-language models (VLMs) into au-  
107 tonomous driving has recently gained significant traction  
108 [15]. Recent works such as DriveVLM [33] and Drive-  
109 eLM [30] use large vision-language models to interpret  
110 traffic scenes, perform structured reasoning, and support  
111 driving-related decision making or explanation generation  
112 [30, 33, 35]. Related multimodal driving frameworks such  
113 as DriveLLaVA and Dolphins further show that language-  
114 guided models can support behavior-level driving deci-  
115 sions and scene understanding in complex road environ-  
116 ments [22, 35]. More broadly, embodied multimodal sys-  
117 tems such as PaLM-E and RT-2 demonstrate that large  
118 models can connect visual perception, language reasoning,  
119 and action-oriented outputs in real-world decision-making  
120 settings [3, 8]. Collectively, these studies suggest that  
121 VLMs can contribute not only to perception enhancement,  
122 but also to higher-level semantic interpretation and human-  
123 interpretable reasoning in autonomous systems [3, 8, 15].

124 This capability is particularly relevant in driving be-  
125 cause some challenging situations arise not only from sen-  
126 sor noise or low confidence, but from ambiguity in scene  
127 semantics, intent, and interaction context [22, 30, 33]. Ex-  
128 amples include unusual construction layouts, temporary  
129 traffic control, and other long-tail scenarios where purely  
130 rule-based or perception-only reasoning may be insufficient

[15, 30, 33]. However, most prior VLM-based driving sys- 131  
tems are designed to improve the autonomy stack itself, 132  
such as through scene understanding, driving assistance, 133  
decision making, or planning [18, 22, 26, 30, 33, 35]. In con- 134  
trast, our work focuses specifically on the human-robot in- 135  
teraction layer, using VLM-based reasoning not to replace 136  
vehicle control, but to determine when autonomy should re- 137  
quest remote human assistance. This framing shifts the role 138  
of the VLM from direct driving support to context-aware, 139  
explainable escalation for teleoperation. 140

**Teleoperation and Remote Driving.** Teleoperation re- 141  
mains an important fallback mechanism for autonomous ve- 142  
hicle deployment, particularly when the onboard system en- 143  
counters situations it cannot safely or reliably resolve on 144  
its own [6, 20]. Prior work on remote driving and remote 145  
assistance has largely focused on communication latency, 146  
predictive displays, interface design, and operator perfor- 147  
mance, reflecting the practical difficulty of supervising or 148  
controlling a vehicle from a distance [5, 23, 32]. More 149  
recent studies also emphasize human-centered workstation 150  
design and operator situation awareness as critical require- 151  
ments for safe and effective remote assistance [29]. At the 152  
system level, latency mitigation and communication robust- 153  
ness continue to be treated as major technical bottlenecks 154  
for scalable teleoperation [11, 16, 20]. 155

However, the question of *when* a vehicle should re- 156  
quest teleoperation is often treated more simply than the 157  
question of *how* teleoperation should be delivered. In 158  
many existing systems, escalation is tied to heuristic trig- 159  
gers such as low confidence, low speed, operational design 160  
domain boundaries, or predefined remote-assistance work- 161  
flows [11, 28, 29]. While such rules are practical for deploy- 162  
ment, they do not fully address semantically ambiguous sit- 163  
uations in which the challenge lies not only in sensing fail- 164  
ure, but in interpreting intent, unusual roadway structure, 165  
or nonstandard human signaling [2, 31]. Examples include 166  
temporary construction control, uncommon edge-case traf- 167  
fic behavior, and traffic-directing gestures from pedestrians 168  
or police officers, all of which can challenge automated sys- 169  
tems even when no clear numerical failure signal is present 170  
[2, 31]. 171

**Reasoning Distillation.** Aligning smaller, deployable 172  
models via teacher supervision has evolved beyond simple 173  
label imitation [14]. Recent works demonstrate the efficacy 174  
of Chain-of-Thought (CoT) distillation and reasoning su- 175  
pervision, where student models are trained to output the 176  
intermediate logic traces generated by large teacher mod- 177  
els [13, 17]. Our framework relies heavily on this reasoning 178  
distillation paradigm to transfer complex driving context in- 179  
terpretation into a lightweight model. 180

**Preference Alignment and Safety.** Aligning language 181  
models with human preferences through Supervised Fine- 182  
Tuning (SFT) and methods such as Direct Preference Op- 183

184 timization (DPO) has become a common paradigm for  
185 improving model behavior in instruction-following sys-  
186 tems [27]. DPO is particularly attractive because it opti-  
187 mizes preferred outputs over rejected ones without requir-  
188 ing an explicit reward model [27]. However, applying DPO  
189 to safety-critical autonomous driving introduces additional  
190 challenges, especially under severe class imbalance [9]. In  
191 remote driving, truly safety-critical handoff events are rare  
192 compared with routine situations, which can bias standard  
193 preference optimization toward under-escalation in high-  
194 risk cases [9]. To address this, we build on standard DPO  
195 with an asymmetric, safety-aware loss that places greater  
196 penalty on misclassified edge cases. This design is moti-  
197 vated by the same principle as focal loss, which emphasizes  
198 hard and underrepresented examples during training [19],  
199 but we adapt it to pairwise preference learning for escala-  
200 tion alignment.

### 201 3. Problem Statement

202 We formulate remote driving request escalation as a dis-  
203 crete three-class triage problem, moving beyond heuristic  
204 or confidence-based approaches that lack contextual  
205 understanding and typically trigger intervention only after  
206 failures occur. Given a camera observation  $\mathcal{X}$ , a vision-  
207 language model (VLM) predicts an escalation level  
208  $y \in \text{AUTONOMY, SUPERVISION, TELEOPERATION}$ ,  
209 along with a natural-language reasoning trace  $\mathcal{R}$  that ex-  
210 plains the decision. Unlike reactive methods, our reasoning-  
211 driven approach enables proactive escalation before poten-  
212 tial incidents arise and provides context-aware guidance, al-  
213 lowing operators to select appropriate actions based on the  
214 evolving driving scenario.

215 We define these triage levels according to the degree of  
216 human involvement required for safe operation:

- 217 • **AUTONOMY:** The AV can safely handle the scene with-  
218 out operator involvement. This includes routine but non-  
219 trivial driving situations such as following vehicles, nav-  
220 igating normal lane traffic, and yielding to pedestrians at  
221 crosswalks.
- 222 • **SUPERVISION:** The AV can continue operating, but en-  
223 vironmental degradation or scene complexity may reduce  
224 the reliability of perception or planning. The remote op-  
225 erator should monitor the situation closely, although im-  
226 mediate intervention is not yet necessary.
- 227 • **TELEOPERATION:** The AV has reached its operational  
228 limits and cannot safely proceed without direct human  
229 control. The remote operator must assume control im-  
230 mediately. Examples include road blockages, police-  
231 directed traffic overriding signals, or large animals ob-  
232 structing the lane.

233 A teleoperation center must oversee  $N$  simultaneous AV  
234 feeds. Consequently, the system must process  $N$  visual  
235 streams under strict latency and compute constraints, pre-

cluding the use of VLMs on every feed. Our objective  
is therefore to train a lightweight VLM policy  $\pi_\theta(y, \mathcal{R} \mid x)$   
that reliably identifies TELEOPERATION cases while  
maintaining a low false-positive rate on AUTONOMY  
scenes.

## 241 4. Methodology

### 242 4.1. System Overview

operates as a modular monitoring layer at the teleoperation  
centre, entirely decoupled from the vehicle’s onboard au-  
tonomy stack. Each autonomous vehicle streams its front-  
camera feed to the centre, where a lightweight student VLM  
(Qwen3-VL-2B-Instruct) processes incoming frames and  
produces a structured triage assessment. This architecture  
ensures that the monitoring system neither competes for on-  
board compute nor introduces a single point of failure into  
the driving pipeline: if the triage VLM fails, the AV con-  
tinues operating under its native autonomy stack while the  
operator retains manual monitoring capability. The student  
model is trained through a two-stage pipeline comprising  
(i) preference data generation via a large teacher VLM, fol-  
lowed by (ii) safety-aware preference alignment, described  
in detail below.

### 258 4.2. Dataset Construction

Because no publicly available datasets explicitly label  
scenes for Level-4 teleoperation triage, we construct a cus-  
tom preference alignment dataset from the nuScenes driving  
corpus.

**Edge-case biased sampling.** We sample approximately  
9,500 camera keyframes from nuScenes [4], with an addi-  
tional 1000 frames reserved for held-out evaluation. Rather  
than sampling uniformly, we apply a priority-biased strat-  
egy that deliberately over-represents semantically challeng-  
ing scenes. Specifically, we assign higher sampling prob-  
ability to frames containing pedestrians near the roadway,  
construction vehicles or traffic cones, complex multi-way  
intersections, and scenes captured under adverse weather or  
lighting conditions (rain, fog, heavy glare). This strategy  
surfaces the long-tail scenarios most likely to cause percep-  
tion or planning failures, precisely the situations where cor-  
rect triage is critical.

**Two-stage annotation pipeline.** Ground-truth reasoning  
traces and triage labels are generated through a two-stage  
process designed to balance annotation throughput with la-  
bel quality:

1. **Machine annotation:** Cosmos Reason 2 8B, a large-  
scale VLM post-trained on driving data with phys-  
ical reasoning capabilities, processes each of the

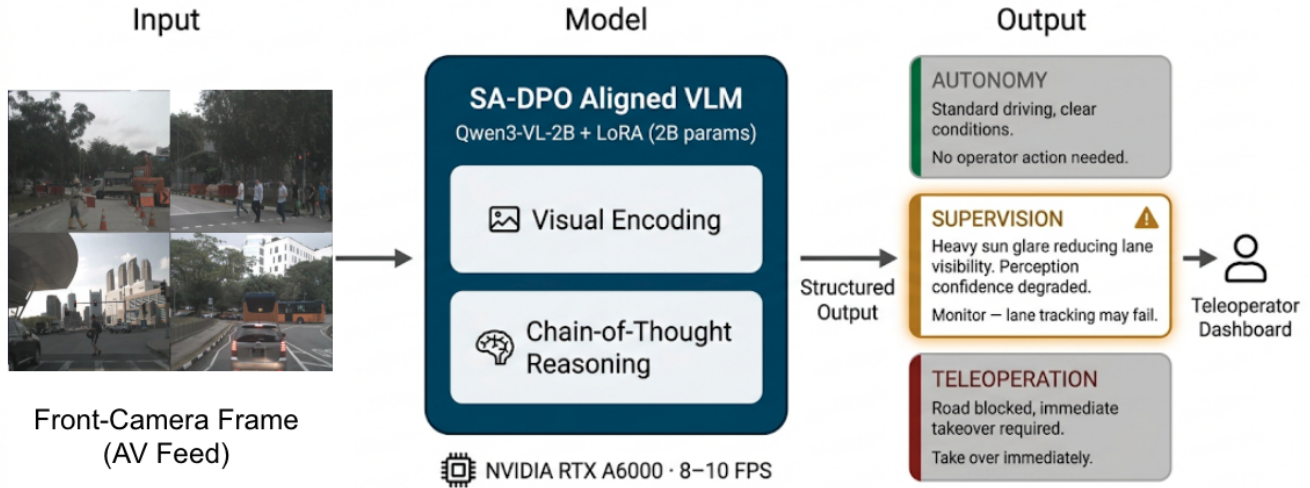


Figure 1. Architecture Diagram

283 9,500 frames. For every frame, the model gen-  
 284 erates a detailed Chain-of-Thought (CoT) reasoning  
 285 trace analysing the driving context identifying rel-  
 286 evant actors, environmental conditions, and poten-  
 287 tial hazards followed by a final triage label  $y \in$   
 288  $\{\text{AUTONOMY, SUPERVISION, TELEOPERATION}\}$ .

289 2. **Human-in-the-loop verification:** Human annotators  
 290 review and, where necessary, correct the machine-  
 291 generated labels. This verification pass is particularly  
 292 important for ambiguous edge cases where the boundary  
 293 between SUPERVISION and TELEOPERATION is context-  
 294 dependent (e.g., a construction zone with partial lane  
 295 blockage).

296 This curated data constitutes the “chosen” preference tier  
 297  $y_w$ . To construct the “rejected” tier  $y_l$ , we run the pre-  
 298 trained, unaligned student model (Qwen3-VL-2B-Instruct)  
 299 on the same 9,500 frames, producing baseline zero-shot re-  
 300 sponses. Together, these form the preference dataset  $\mathcal{D} =$   
 301  $\{(x_i, y_{w,i}, y_{l,i})\}_{i=1}^{9500}$  required for DPO training.

302 **Observed class distribution.** After edge-case biased  
 303 sampling, the resulting dataset exhibits the following class  
 304 frequencies: approximately 49% AUTONOMY, 38% SU-  
 305 PERVISION, and 13% TELEOPERATION. While our biased  
 306 sampling strategy deliberately enriches the representation  
 307 of safety-critical scenes to provide sufficient training signal,  
 308 in uncurated AV logs TELEOPERATION-level events occur  
 309 in fewer than 1–2% of frames. This deployment-time class  
 310 imbalance is precisely what motivates the asymmetric loss  
 311 formulation described in Section 4.4.

### 312 4.3. Teacher-Student Distillation Pipeline

313 Our framework employs reasoning distillation to transfer  
 314 domain-specific driving knowledge from a large, capable

teacher into a compact student model suitable for multi-feed  
 deployment.

317 **Teacher model.** We utilise as the teacher. Unlike  
 318 general-purpose VLMs, has been post-trained on embod-  
 319 ied driving scenarios and exhibits strong physical reason-  
 320 ing capabilities—for instance, understanding that a police  
 321 officer’s hand signal overrides a green traffic light, or that  
 322 heavy solar glare may render lane markings invisible to the  
 323 vehicle’s perception system. These nuanced judgments are  
 324 encoded in the teacher’s CoT reasoning traces, which serve  
 325 as rich supervision signals beyond simple label matching.

326 **Student model and LoRA fine-tuning.** The student  
 327 model, Qwen3-VL-2B-Instruct, is a 2-billion param-  
 328 eter VLM chosen specifically for its favourable compute-  
 329 to-quality trade-off in multi-feed deployment scenarios.  
 330 Rather than fine-tuning all model parameters, we apply  
 331 Low-Rank Adaptation (LoRA) to the attention projection  
 332 matrices, adding fewer than 0.5% additional trainable pa-  
 333 rameters. LoRA enables rapid adaptation with modest GPU  
 334 memory requirements while preserving the student’s pre-  
 335 trained visual and linguistic representations. During train-  
 336 ing, the LoRA weights are optimised via SA-DPO (de-  
 337 scribed below) while the base model weights remain frozen.  
 338 At inference, the LoRA adapters are merged back into the  
 339 base model, incurring zero additional latency.

340 **Preference pair construction.** For each training frame  
 341  $x_i$ , the preference pair  $(y_{w,i}, y_{l,i})$  captures the contrast be-  
 342 tween expert-level reasoning (teacher) and naïve baseline  
 343 reasoning (student). The “chosen” response  $y_{w,i}$  contains  
 344 the teacher’s full CoT trace followed by the correct triage

345 label and a structured action recommendation, while the  
346 “rejected” response  $y_{l,i}$  is the student’s unaligned zero-  
347 shot output on the same frame. This pairing teaches the  
348 student not only *what* the correct triage should be, but  
349 *why*—transferring the semantic reasoning process rather  
350 than merely imitating labels.

#### 351 4.4. Safety-Aware DPO (SA-DPO)

352 Direct Preference Optimization (DPO) aligns language  
353 model outputs with human preferences by optimising the  
354 policy to increase the likelihood of preferred responses rel-  
355 ative to dispreferred ones. The standard DPO loss is:

$$356 \mathcal{L}_{\text{DPO}} = -\log \sigma \left( \beta \left( \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right) \quad (1)$$

357 where  $\pi_{\theta}$  is the policy being trained,  $\pi_{\text{ref}}$  is the frozen refer-  
358 ence policy,  $\beta$  controls the deviation from the reference, and  
359  $\sigma$  is the logistic function. While effective for general align-  
360 ment, standard DPO assigns equal importance to all prefer-  
361 ence pairs regardless of the safety implications of misclas-  
362 sification.

363 **The under-escalation problem.** Although our edge-case  
364 biased sampling enriches TELEOPERATION representation  
365 in the training set, AUTONOMY and SUPERVISION frames  
366 still constitute the vast majority of preference pairs. Be-  
367 cause DPO aggregates gradients across all pairs, the loss  
368 landscape is dominated by these routine examples. Conse-  
369 quently, a model trained with standard DPO learns to sys-  
370 tematically under-escalate—predicting AUTONOMY or SU-  
371 PERVISION even when the scene requires immediate hu-  
372 man takeover because the rare under-escalation errors are  
373 “drowned out” by the frequent correct classifications on  
374 routine scenes. This problem is further amplified at deploy-  
375 ment, where TELEOPERATION-level events occur in fewer  
376 than 1–2% of frames. The behaviour is analogous to the  
377 class-imbalance problem in object detection that motivated  
378 Focal Loss, but manifests in the preference-learning do-  
379 main.

380 **Asymmetric safety weighting.** To address this, we pro-  
381 pose Safety-Aware DPO (SA-DPO), which introduces a  
382 pair-dependent weighting function  $w(y_w, y_l)$ :

$$383 \mathcal{L}_{\text{SA-DPO}} = -w(y_w, y_l) \log \sigma (\beta \cdot \Delta_{\theta}(x, y_w, y_l)) \quad (2)$$

384 where the preference margin is:

$$385 \Delta_{\theta}(x, y_w, y_l) = \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \quad (3)$$

386 The weighting function encodes the asymmetric safety

cost of different misclassification types:

$$w(y_w, y_l) = \begin{cases} \alpha_{\text{high}} & \text{if } y_w = \text{TELEOP.}, y_l = \text{AUTONOMY} \\ \alpha_{\text{med}} & \text{if } y_w = \text{TELEOP.}, y_l = \text{SUPERVISION} \\ 1.0 & \text{otherwise} \end{cases} \quad (4)$$

The most dangerous error—classifying a mandatory  
TELEOPERATION scene as AUTONOMY—receives the high-  
est penalty ( $\alpha_{\text{high}} = 5.0$ ). A less severe but still problematic  
error—downgrading TELEOPERATION to SUPERVISION—  
receives a moderate penalty ( $\alpha_{\text{med}} = 2.0$ ). All other pairs,  
including over-escalation (e.g., AUTONOMY classified as  
TELEOPERATION), are weighted at 1.0, reflecting the op-  
erational reality that false alarms are inconvenient but not  
safety-critical.

**Key advantages over alternative approaches.** Unlike  
data-level balancing strategies (e.g., oversampling TELE-  
OPERATION frames), SA-DPO preserves the natural class  
distribution during training, ensuring the model does not  
develop an artificially inflated prior for rare events. Un-  
like standard class-weighted classification losses that op-  
erate on individual predictions, SA-DPO operates on *pair-  
wise preference comparisons between full reasoning trajec-  
tories*, allowing the model to learn the subtle semantic dif-  
ferences between scenes that require escalation and those  
that do not. For instance, the model learns that “a police  
officer standing on the sidewalk” (AUTONOMY) and “a po-  
lice officer stepping into the roadway with a raised hand”  
(TELEOPERATION) differ by a contextual detail that a sim-  
ple label-weighted classifier would struggle to capture.

#### 4.5. Structured Output for Teleoperator Interface

The model produces a structured natural-language explana-  
tion alongside the triage decision, designed for rapid com-  
prehension by teleoperators monitoring multiple feeds si-  
multaneously. Each output consists of three fields: (i) the  
triage level (AUTONOMY, SUPERVISION, or TELEOPERA-  
TION), (ii) a concise reasoning trace (1–2 sentences) ex-  
plaining the contextual factors driving the decision, and  
(iii) a recommended operator action. This structured for-  
mat enables operators to assess an alert in under two sec-  
onds without interpreting abstract confidence scores or  
bounding-box overlays, significantly reducing cognitive  
load during multi-feed monitoring.

## 5. Evaluation and Discussion

### 5.1. Experimental Setup

We evaluate on a curated benchmark of 1,000 front-camera  
frames sampled from the nuScenes dataset. Ground-truth  
triage labels are generated by and verified through human-  
in-the-loop review, yielding an imbalanced but realistic dis-

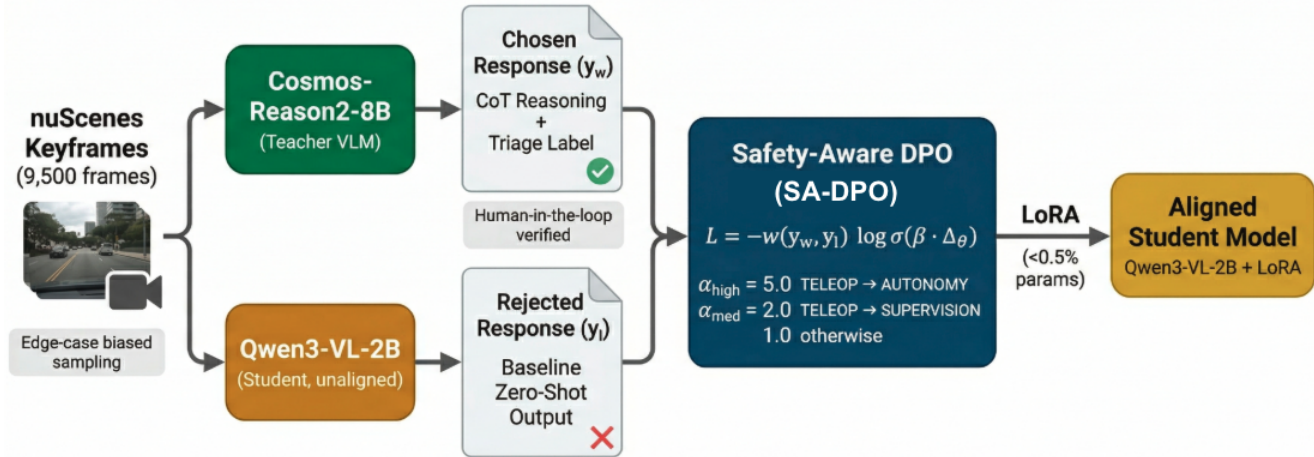


Figure 2. Data Generation and Training Pipeline

Table 1. In-domain nuScenes evaluation ( $N=1,000$ ). Under-Escalation Rate (UE $\downarrow$ ) is lower-is-better.

Method	Acc.	F1	UE $\downarrow$
Baseline (Zero-Shot)	0.475	0.292	71.0%
Standard DPO	0.730	0.624	37.9%
SA-DPO (Ours)	<b>0.921</b>	<b>0.907</b>	<b>2.1%</b>

tribution: 467 AUTONOMY, 388 SUPERVISION, and 145 TELEOPERATION frames. We compare three systems:

1. **Baseline (Zero-Shot)**: The pre-trained without any preference alignment, representing a raw VLM applied directly to the triage task.
2. **Standard DPO**: aligned with the conventional DPO loss using equal weighting for all preference pairs.
3. **SA-DPO ()**: Our proposed model aligned with the safety-aware asymmetric loss that up-weights under-escalation penalties.

All inference is performed using vLLM on a single NVIDIA RTX A6000 (48 GB).

## 5.2. Triage Classification Results

We report Overall Accuracy, Macro-F1, and the safety-critical *Under-Escalation Rate* (UE $\downarrow$ ), defined as the fraction of true TELEOPERATION events misclassified as AUTONOMY.

As shown in Table 1, SA-DPO achieves 92.1% overall accuracy and a Macro-F1 of 0.907, substantially outperforming both the zero-shot baseline (47.5%, 0.292) and standard DPO (73.0%, 0.624). Most critically, SA-DPO reduces the under-escalation rate to **2.1%**—missing only 3 out of 145 true TELEOPERATION events—whereas Standard DPO still misses 37.9% and the unaligned baseline misses

71.0%. The 42.4 percentage-point accuracy gap between SA-DPO and the zero-shot baseline underscores the necessity of preference alignment: without it, the VLM defaults to predicting AUTONOMY for nearly all inputs, rendering it unsafe for deployment. Notably, the jump from Standard DPO to SA-DPO (+19.1% accuracy, +0.283 F1) is driven almost entirely by the asymmetric loss weighting, confirming that safety-aware reweighting is essential for learning from imbalanced triage distributions.

**Per-class recall analysis.** Figure 4 reports recall broken down by triage level. The zero-shot baseline achieves moderate AUTONOMY recall (90.6%) by predominantly predicting the majority class, but catastrophically fails on SUPERVISION (12.4%) and TELEOPERATION (2.8%). Standard DPO improves SUPERVISION recall to 79.1% but still captures only 24.8% of TELEOPERATION events. SA-DPO achieves the highest recall across all three classes—94.4% AUTONOMY, 89.2% SUPERVISION, and 92.4% TELEOPERATION—demonstrating that the asymmetric loss successfully shifts the model’s attention toward safety-critical scenes without sacrificing performance on routine ones.

Figure 3 presents confusion matrices for all three models. The escalation-error trade-off is shown in Figure 5: SA-DPO reduces UE to just 2.1% with only 1.5% over-escalation, a favourable safety profile where the residual false alarms are far less costly than missed interventions.

## 5.3. Multi-Feed Throughput Analysis

A key operational requirement of is its ability to monitor multiple AV feeds concurrently without competing for on-board vehicle compute. We deploy the SA-DPO—aligned via vLLM on a single NVIDIA RTX A6000 GPU (48 GB VRAM). In our benchmarking, the 2B-parameter model

	Baseline (Zero-Shot)			Standard DPO			SA-DPO (Ours)					
True	AUT	429	38	0	AUT	387	64	16	AUT	429	26	12
	SUP	332	48	8	SUP	40	307	41	SUP	19	343	26
	TEL	103	38	4	TEL	55	54	36	TEL	0	12	133
		AUT	SUP	TEL		AUT	SUP	TEL		AUT	SUP	TEL
		Predicted				Predicted				Predicted		

Figure 3. Confusion matrices for the three evaluated models on the 1,000-frame nuScenes benchmark

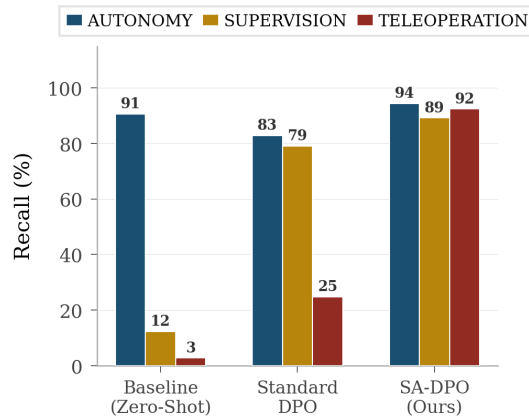


Figure 4. Per-class recall by triage level. SA-DPO achieves balanced, high recall across all three classes

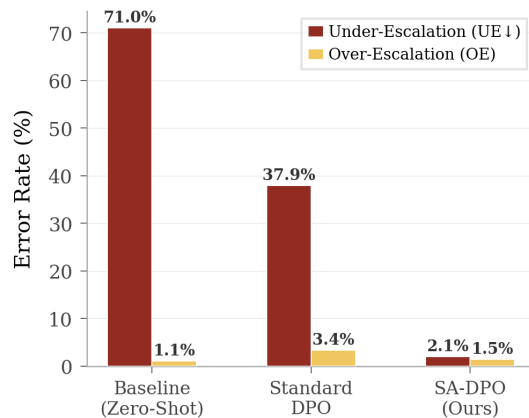


Figure 5. Under-escalation (UE, safety-critical misses) vs. over-escalation (OE, false alarms). SA-DPO reduces dangerous under-escalation to 2.1% while keeping false alarms at 1.5%

processes front-camera frames at **8–10 FPS** per stream, including full structured-output generation. The compact model size also enables deployment on more cost-effective GPUs (e.g., RTX 4090) with only modest throughput reduction, making fleet-scale triage economically viable.

## 6. Conclusion

In this work, we presented a vision-language reasoning framework engineered to optimize the human-robot interaction bottleneck in autonomous vehicle teleoperation. By reformulating teleoperation decision from a heuristic confidence threshold to a structured, semantic reasoning problem, we enable early escalation based on contextual understanding.

Our primary contribution, Safety-Aware DPO (SA-DPO), addresses the critical class imbalance inherent in real-world driving data by applying an asymmetric loss that specifically penalizes missed intervention requests. Our results demonstrate that SA-DPO substantially mitigates under-escalation errors compared to standard alignment techniques, without resorting to artificial data resampling. As AV fleets scale, deploying lightweight, semantically aligned triage VLMs like at the teleoperation center level offers a practical, scalable mechanism for ensuring safe and highly interpretable human-in-the-loop oversight.

## References

- [1] Association for Advancing Automation. Benefits of remote control capabilities for industrial robots. <https://www.automate.org/industry-insights/benefits-remote-control-capabilities-industrial-robots>, 2024. 1
- [2] Tonko EW Bossen, Andreas Møgelmoose, and Ross Greer. Can vision-language models understand and interpret dynamic gestures from pedestrians? pilot datasets and exploration towards instructive nonverbal commands for cooperative autonomous vehicles. In *Proceedings of the Computer*

- 523 Vision and Pattern Recognition Conference, pages 4779–  
524 4788, 2025. 2
- 525 [3] Anthony Brohan, Noah Brown, Justice Carbajal, et al. RT-  
526 2: Vision-language-action models transfer web knowledge  
527 to robotic control. In *Proceedings of The 7th Conference on*  
528 *Robot Learning (CoRL)*, pages 2165–2183, 2023. 2
- 529 [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora,  
530 Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Gi-  
531 ancario Baldan, and Oscar Beijbom. nuscenes: A multi-  
532 modal dataset for autonomous driving. In *Proceedings of*  
533 *the IEEE/CVF conference on computer vision and pattern*  
534 *recognition*, pages 11621–11631, 2020. 3
- 535 [5] J. Y. C. Chen, E. C. Haas, and M. J. Barnes. Human per-  
536 formance issues and user interface design for teleoperated  
537 robots. *IEEE Transactions on Systems, Man, and Cyber-*  
538 *netics, Part C (Applications and Reviews)*, 37(6):1231–1245,  
539 2007. 1, 2
- 540 [6] Cogniteam. Cloud-based teleoperation in robotics. <https://www.cogniteam.com/resources/blogs/cloud-based-teleoperation-in-robotics-2/>, 2024. [Online; accessed 1-Mar-2026]. 2
- 541  
542  
543
- 544 [7] Jacob W Crandall, Michael A Goodrich, Dan R Olsen,  
545 and Curtis W Nielsen. Validating human-robot interaction  
546 schemes in multitasking environments. *IEEE Transactions*  
547 *on Systems, Man, and Cybernetics-Part A: Systems and Hu-*  
548 *mans*, 35(4):438–449, 2005. 1
- 549 [8] Danny Driess, Fei Xia, Mehdi S.M. Sajjadi, et al. Palm-e:  
550 An embodied multimodal language model. *Proceedings of*  
551 *the International Conference on Learning Representations*  
552 *(ICLR)*, 2023. 2
- 553 [9] Duanyu Feng, Bowen Qin, Chen Huang, Zheng Zhang, and  
554 Wenqiang Lei. Towards analyzing and understanding the  
555 limitations of dpo: A theoretical perspective. *arXiv preprint*  
556 *arXiv:2404.04626*, 2024. 3
- 557 [10] Fernride. Teleoperation. <https://www.fernride.com/teleoperation>, 2024. 1
- 558  
559
- 560 [11] M. Gontscharow et al. Scalable remote operation for au-  
561 tonomous vehicles: Integration of cooperative perception  
562 and open source communication. In *2024 IEEE Intelligent*  
563 *Vehicles Symposium (IV)*, pages 43–50, Jeju Island, Korea,  
564 2024. 1, 2
- 565 [12] Michael A. Goodrich and Dan R. Olsen Jr. Seven principles  
566 of efficient human robot interaction. In *Proc. IEEE Interna-*  
567 *tional Conference on Systems, Man and Cybernetics*, pages  
3942–3948, 2003. 1
- 568 [13] Namgyu Ho, Laura Schmid, and Se-Young Yun. Large lan-  
569 guage models are reasoning teachers. In *Proceedings of the*  
570 *61st annual meeting of the association for computational lin-*  
571 *guistics (volume 1: long papers)*, pages 14852–14882, 2023.  
572 2
- 573 [14] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan  
574 Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna,  
575 Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step!  
576 outperforming larger language models with less training data  
577 and smaller model sizes. In *Findings of the Association for*  
578 *Computational Linguistics: ACL 2023*, pages 8003–8017,  
579 2023. 2
- [15] Yifan Jiang, Lei Zhang, Hao Chen, et al. A survey on  
vision–language–action models for autonomous driving. In  
*Proceedings of the IEEE/CVF International Conference on*  
*Computer Vision Workshops (ICCVW)*, 2025. Also available  
as arXiv:2506.24044. 1, 2
- [16] Sidharth Bhanu Kamtam, Qian Lu, Faouzi Bouali,  
Olivier CL Haas, and Stewart Birrell. Network latency in  
teleoperation of connected and autonomous vehicles: A re-  
view of trends, challenges, and mitigation strategies. *Sen-*  
*sors*, 24(12):3957, 2024. 2
- [17] Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren,  
Kai-Wei Chang, and Yejin Choi. Symbolic chain-of-thought  
distillation: Small models can also “think” step-by-step. In  
*Proceedings of the 61st Annual Meeting of the Association*  
*for Computational Linguistics (Volume 1: Long Papers)*,  
pages 2665–2679, 2023. 2
- [18] Qifeng Li, Zhiwei Peng, Lan Feng, and Quanyi Zhang.  
Think2drive: Efficient reinforcement learning by thinking in  
latent world model for quasi-realistic autonomous driving (in  
CARLA-v2), 2024. Also published at ECCV 2024. 2
- [19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and  
Piotr Dollár. Focal loss for dense object detection. In *Pro-*  
*ceedings of the IEEE international conference on computer*  
*vision*, pages 2980–2988, 2017. 3
- [20] Mingjun Lu et al. Teleoperation technologies for enhanc-  
ing connected and autonomous vehicles. *IEEE Transac-*  
*tions on Intelligent Transportation Systems*, 2022. Avail-  
able: <https://weisongshi.org/papers/lu22-teleoperation.pdf>. 2
- [21] Sidi Lu, Ren Zhong, and Weisong Shi. Teleoperation tech-  
nologies for enhancing connected and autonomous vehi-  
cles. In *2022 IEEE 19th International Conference on Mobile*  
*Ad Hoc and Smart Systems (MASS)*, pages 435–443. IEEE,  
2022. 1
- [22] Yingzi Ma, Yulong Cao, Jiachen Sun, Marco Pavone, and  
Chaowei Xiao. Dolphins: Multimodal language model for  
driving. In *European Conference on Computer Vision*, pages  
403–420. Springer, 2024. 2
- [23] A. Matheson, B. Donmez, F. Rehmatullah, P. Jasiobedzki,  
H.-K. Ng, V. Panwar, and M. Li. The effects of predictive  
displays on performance in driving tasks with multi-second  
latency: Aiding tele-operation of lunar rovers. In *Proce-*  
*edings of the Human Factors and Ergonomics Society Annual*  
*Meeting*, pages 21–25, 2013. 2
- [24] Ronja Möller, Antonino Furnari, Sebastiano Battiato, Aki  
Härmä, and Giovanni Maria Farinella. A survey on human-  
aware robot navigation. *Robotics and Autonomous Systems*,  
145:103837, 2021. 1
- [25] Motional. A path forward: Using AI to improve re-  
mote vehicle assistance for AVs. <https://motional.com/news/path-forward-using-ai-improve-remote-vehicle-assistance-avs>, 2024. 2
- [26] Seonghoon Park, Jun Lee, YoungJae Kim, et al. Vlaad: Vi-  
sion and language assistant for autonomous driving. In *Pro-*  
*ceedings of the IEEE/CVF Winter Conference on Applica-*  
*tions of Computer Vision Workshops (WACVW)*, 2024. 2
- [27] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-  
pher D Manning, Stefano Ermon, and Chelsea Finn. Direct

- 638 preference optimization: Your language model is secretly a  
639 reward model. *Advances in neural information processing*  
640 *systems*, 36:53728–53741, 2023. 3
- 641 [28] Daniel J. Rea and Stela H. Seo. Still not solved: A call for re-  
642 newed focus on user-centered teleoperation interfaces. *Frontiers in Robotics and AI*, 9:704225, 2022. 1, 2
- 643 [29] Andreas Schrank, Fabian Walocha, Stefan Brandenburg, and  
644 Michael Oehl. Human-centered design and evaluation of  
645 a workplace for the remote assistance of highly automated  
646 vehicles. *Cognition, Technology & Work*, 26(2):183–206,  
647 2024. 2
- 648 [30] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen,  
649 Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo,  
650 Andreas Geiger, and Hongyang Li. Drivelm: Driving with  
651 graph visual question answering. In *European conference on*  
652 *computer vision*, pages 256–274. Springer, 2024. 2
- 653 [31] Qunying Song, Emelie Engström, and Per Runeson. Industry  
654 practices for challenging autonomous driving systems with  
655 critical scenarios. *ACM Transactions on Software Engineer-*  
656 *ing and Methodology*, 33(4):1–35, 2024. 2
- 657 [32] Felix Tener and Joel Lanir. Driving from a distance: Chal-  
658 lenges and guidelines for autonomous vehicle teleoperation  
659 interfaces. In *Proceedings of the 2022 CHI Conference on*  
660 *Human Factors in Computing Systems (CHI '22)*, pages 1–  
661 13, New York, NY, USA, 2022. Association for Computing  
662 Machinery. 1, 2
- 663 [33] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang,  
664 Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and  
665 Hang Zhao. Drivelm: The convergence of autonomous  
666 driving and large vision-language models. *arXiv preprint*  
667 *arXiv:2402.12289*, 2024. 2
- 668 [34] Adrien Wantiez, Tianming Qiu, Stefan Matthes, and Hao  
669 Shen. Scene understanding for autonomous driving using vi-  
670 sual question answering. In *2023 International Joint Confer-*  
671 *ence on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2023.  
672 1
- 673 [35] Rui Zhao, Qirui Yuan, Jinyu Li, Yuze Fan, Yun Li, and Fei  
674 Gao. DriveLLaVA: Human-level behavior decisions via vi-  
675 sion language model. *Sensors*, 24(13):4113, 2024. 2
- 676