

Remedying the Curse of Autonomous Driving: VLM Driven Training-Free Framework for Efficient Long-Tail Video Detection

Anonymous CVPR submission

Paper ID ***

Abstract

001 *Autonomous driving has made remarkable progress, and*
002 *we have finally witnessed its real-world deployment. These*
003 *advancements have been driven by training models on an*
004 *increasing scale of data. Nowadays, data collection is*
005 *streamlined, with vast amounts—amounting to ~ 100 years*
006 *of driving data—collectible in a single day. However, most*
007 *of this data is routine and does not help the model to gen-*
008 *eralize; on the contrary, it might bias models towards rou-*
009 *tine driving scenarios. This creates a critical bottleneck in*
010 *model generalization, as systems remain vulnerable to long-*
011 *tail scenarios that are rare but safety-critical. Hence, train-*
012 *ing models on such scenarios remains critical for safe and*
013 *scalable deployment.*

014 *Despite extensive research on end-to-end driving mod-*
015 *els, a systematic method for efficiently detecting these long-*
016 *tail events from large-scale dataset is missing. In this*
017 *work, we propose a novel training-free two-stage frame-*
018 *work based on vision language models. In the first stage, a*
019 *small language model is used to process videos for scenario*
020 *summarization. In the second stage, the scenario summa-*
021 *rization from the first stage is processed by a large language*
022 *model to rank the video’s long-tail relevance. Our two-*
023 *stage framework is designed to efficiently process industry-*
024 *scale video datasets and accurately classify the relevance*
025 *of video segments to long-tail events. Experiment results*
026 *show our proposed method surpasses counterpart methods*
027 *by 24% in AUC and provides a $\sim 14\times$ faster inference com-*
028 *pared to counterpart multimodal large language models,*
029 *enabling scalable and targeted data mining for autonomous*
030 *systems, which is critical for generalizing driving models*
031 *for deployment. We hope our study paves the way for fur-*
032 *ther research in this critical field.*

033 1. Introduction

034 In recent years, autonomous driving systems have evolved
035 rapidly, reaching a stage where commercial deployment

has finally become a reality after decades of research and 036
development. Traditionally, these systems have followed 037
a modular design paradigm—comprising separate percep- 038
tion, prediction, and planning components—each trained 039
and optimized independently [2, 26]. However, this dis- 040
joint optimization often leads to compounding errors, mo- 041
tivating a shift toward end-to-end approaches that directly 042
map raw sensory inputs to driving actions [10, 11, 13, 19, 043
22]. End-to-end systems themselves have undergone sig- 044
nificant evolution, transitioning from cascaded, planning- 045
oriented architectures to Vision-Language-Action (VLA) 046
models—large transformer-based neural networks that inte- 047
grate perception, reasoning, and control within a single 048
framework [5, 12, 14, 21, 28]. VLAs have recently gained 049
increasing attention and are emerging as a preferred archi- 050
tecture for real-world deployment [24]. With billions to tril- 051
lions of parameters, these models demonstrate strong gen- 052
eralization capabilities by learning rich, transferable repre- 053
sentations. 054

The success of such generalizable systems stems from 055
large-scale training enabled by fleet-collected driving data. 056
Modern pipelines can accumulate the equivalent of ~ 100 057
years of driving data in a single day [24]. Within 058
these datasets, long-tail scenarios—rare yet safety-critical 059
events—are key to achieving robust and deployable au- 060
tonomous driving. However, these cases constitute $\ll 1\%$ 061
of the data, making their identification and utilization ex- 062
tremely challenging [23]. Several prior works have ex- 063
plored building learnable systems to identify long-tail sce- 064
narios from large-scale driving data [3, 4]. These works 065
leveraged various techniques such as reconstruction, pre- 066
diction and generative models for identification of long- 067
tail videos [3, 8, 15]. However, these traditional methods 068
lack generalization, as they are tailored to specific types of 069
long-tail cases, whereas real-world long-tail distributions 070
are inherently open-ended and diverse. Moreover, these 071
conventional vision-only systems lack reasoning capabili- 072
ties, which further limits their applicability to generalized 073
and complex driving scenarios. An alternative to traditional 074
methods is to leverage powerful multimodal large language 075

076	models (MLLMs). However, these models incur substantial	127
077	latency and memory overhead, resulting in high computa-	128
078	tional costs that hinder deployment at fleet scale.	129
079	To overcome these challenges, we propose long-tail	130
080	video anomaly detection for autonomous driving (LTVAD-	131
081	AD), a systematic training-free framework based on vision	132
082	language models (VLMs) for efficiently detecting long-tail	133
083	videos from massive real-world driving datasets in an open-	134
084	ended manner, enabling scalable discovery beyond pre-	135
085	defined anomaly types. In particular, our framework adopts	136
086	a reasoning-driven two-stage system that integrates scalable	137
087	multimodal understanding with high-level semantic reason-	138
088	ing. In the first stage, a small language model summarizes	139
089	driving videos into concise textual descriptions that capture	140
090	key events and context. In the second stage, a large lan-	141
091	guage model reasons over these summaries to assess and	142
092	rank their relevance to long-tail scenarios. This framework	143
093	bridges low-level perception and high-level abstraction, en-	144
094	abling efficient large-scale analysis and open-ended discov-	145
095	ery of rare, safety-critical driving events.	146
096	To evaluate the effectiveness of our approach, we con-	147
097	duct experiments on a proprietary dataset comprising both	148
098	long-tail and routine driving scenarios, following estab-	149
099	lished practices in the broader video anomaly detection lit-	150
100	erature. The experimental results demonstrate that our pro-	151
101	posed method significantly outperforms counterpart meth-	152
102	ods by 24% in AUC, and provides a $\sim 14\times$ faster inference	153
103	compared to larger video-language models ($\sim 10\times$ more pa-	154
104	rameters) that are impractical for large scale deployment	155
105	due to their latency and memory requirements. Our main	
106	contributions are summarized as follows:	
107	• We propose a systematic training-free framework based	
108	on vision language models for efficiently detecting long-	
109	tail driving scenarios from large-scale real-world datasets	
110	in an open-ended manner.	
111	• We design a reasoning-driven two-stage architecture that	
112	integrates a small language model for scenario under-	
113	standing&summarization, and a large language model for	
114	semantic reasoning and long-tail relevance ranking.	
115	• We provide detailed guidelines for constructing propi-	
116	etary video anomaly detection datasets tailored to au-	
117	tonomous driving applications and future research in this	
118	critical domain.	
119	• Extensive experimental results show that our approach	
120	significantly outperforms comparable baselines in terms	
121	of both long-tail video detection and inference latency,	
122	making it ideal for large scale applications.	
123	2. Related Work	
124	2.1. End-to-End Autonomous Driving	
125	End-to-end autonomous driving methods, which map raw	
126	sensor inputs directly to control outputs, have received re-	
	markable attention in recent years and have rapidly evolved,	127
	paving the way for scalable deployment of autonomous	128
	driving systems. UniAD [11] is one of the seminal works	129
	that inspired this new wave of research in end-to-end au-	130
	tonomous driving. Unlike traditional stand-alone modular	131
	systems, UniAD leverages rasterized representations and	132
	integrates perception, prediction, and planning within a uni-	133
	fied framework, enabling joint optimization across mod-	134
	ules. VAD [13] which uses vectorized scene representa-	135
	tion is another important line of work for end-to-end au-	136
	tonomous driving. Following these two works, a plethora of	137
	end-to-end works have been proposed. Some notable works	138
	include PARA-Drive [22], SparseDrive [19], and Diffusion-	139
	Drive [16]. Despite the notable progress, these vision-only	140
	end-to-end methods face generalization issues as it lacks the	141
	multi-modality and reasoning.	142
	To address the challenges in vision-only E2E systems,	143
	multimodal large language models have recently been in-	144
	vestigated. Initial work in the field distilled features from	145
	large vision language models into vision-only end-to-end	146
	systems to enhance the generalization performance [9, 20].	147
	More recent works have focused on vision language action	148
	models (VLAs) which takes input as videos, language (as	149
	well as other information such as audio or vehicle kinemat-	150
	ics) and directly outputs actions [5, 12, 24, 28]. These mod-	151
	els have achieved state-of-the-art performance and strong	152
	generalization as these large transformer neural networks	153
	have \sim trillion number of parameters and trained on vast	154
	scale of diverse fleet data.	155
	2.2. Video Anomaly Detection	156
	Long-tail scenario detection in autonomous driving is a	157
	video anomaly detection task where any non-routine inter-	158
	esting driving scenario can be considered as anomaly. Early	159
	works in the field leveraged classical CNN or GAN based	160
	methods for this task. Among these works, [8] employs an	161
	autoencoder and use reconstruction error as an indicator for	162
	anomaly score as reconstruction error for non-routine driv-	163
	ing scenarios is anticipated to be high. In another work,	164
	[3] performs future frame prediction, and assign a high	165
	anomaly score if there is a non-predictable relevant object	166
	in a relevant location. However, such works lack reasoning	167
	and limit the long-tail scenarios to a subset, thus are not able	168
	to generalize to open-ended nature of long-tail scenarios.	169
	Given the limitations of conventional methods, more re-	170
	cent works [23] leverages the multimodal LLMs for extract-	171
	ing long-tail scenarios, but does not provide any method-	172
	ological aspects on how MLLMs are leveraged. In this	173
	work, we provide a systematic and detailed approach for	174
	efficiently mining long-tail videos using vision language	175
	models. Moreover, proposed approach is training free, fur-	176
	ther enabling generalization to open-ended nature of long-	177
	tail scenarios.	178

179 **3. Methodology**180 **3.1. Problem Set-up**

181 Video anomaly detection (VAD) systems traditionally aim
182 to *learn* a classifier model by training on a ground truth
183 dataset labeled video-level as normal or anomaly. How-
184 ever, challenge in curating long-tail videos and open-ended
185 nature of long-tail videos renders such methods impractical
186 as they bias towards training scenarios and fail to well-
187 generalize to unseen long-tail scenarios [4].

188 Hence, in this study, we present a *training-free* frame-
189 work leveraging pretrained MLLMs that is not bounded by
190 training distribution, and can leverage advanced reasoning
191 capabilities to generalize to any long-tail scenario. In par-
192 ticular, given a driving scenario video, \mathbf{V} , our goal is to de-
193 sign a training-free framework f to directly classify a given
194 video as normal or anomaly.

195 **3.2. Vanilla LTVAD-AD**

196 In a vanilla setup, multimodal large language models
197 (MLLMs) capable of processing videos [1] can be directly
198 employed for long-tail video classification through care-
199 fully designed prompts. Given a video \mathbf{V} and a scoring
200 prompt P_{score} , the MLLM can be instructed to produce an
201 anomaly score s , reflecting the likelihood that the video de-
202 picts a rare or interesting event (Fig. 1 baseline). To achieve
203 this, we craft the P_{score} (template in Fig. 1, full prompt is
204 provided in Appendix) to enforce the MLLM to output a
205 score $s \in 0.0, 0.1, 0.2, \dots, 1.0$, representing 11 uniformly
206 sampled values in the interval $[0, 1]$, where 0 denotes nor-
207 mal and 1 denotes anomaly, formulated as:

$$208 \quad s = f_{\text{MLLM}}(\mathbf{V}, P_{\text{score}}), \quad (1)$$

209 where $f_{\text{MLLM}}(\cdot)$ denotes the multimodal LLM inference
210 function. The resulting score s is then thresholded to clas-
211 sify the video as either normal or anomaly:

$$212 \quad \hat{y} = \begin{cases} 1, & \text{if } s \geq \tau, \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

213 where τ is a predefined decision threshold.

214 While this setup (with a well-crafted P_{score}) can enable
215 large MLLMs to achieve good performance, processing
216 videos through such models is computationally expensive
217 and impractical for large-scale deployment due to their high
218 inference cost and memory demands. Conversely, small
219 language models (SLMs) are computationally efficient and
220 scalable, but their performance remains limited even with
221 carefully designed prompts, owing to weaker reasoning and
222 representation capabilities. This trade-off between capabil-
223 ity and scalability motivates the need for a more efficient yet
224 reasoning-aware framework for long-tail video understand-
225 ing.

226 **3.3. LTVAD-AD**

227 To address the trade-offs described in 3.2, we propose a sys-
228 tematic training-free two-stage framework based on mul-
229 timodal language models for efficiently detecting long-tail
230 videos in an open-ended manner. As illustrated in Fig. 1c,
231 the core idea is to decouple video understanding and rea-
232 soning across two complementary stages, thereby achieving
233 both scalability and interpretability.

Stage 1: SLMs as Scenario Descriptor. In the first stage,
234 SLMs summarizes raw driving videos into concise textual
235 descriptions that capture key events and context. Long-
236 tail scenarios—such as emergency vehicles, wildlife cross-
237 ings, or other unexpected events—naturally emerge as non-
238 routine information in these summaries. Given the limited
239 capability and context window of SLMs, we split each long
240 video into shorter, non-overlapping segments before pro-
241 cessing to maintain temporal coherence and scalability.
242

Given a video segment \mathbf{V}_i and a description prompt
243 P_{desc} , the SLM generates a textual scenario description
244

$$245 \quad T_i = f_{\text{SLM}}(\mathbf{V}_i, P_{\text{desc}}), \quad (3)$$

246 where $f_{\text{SLM}}(\cdot)$ denotes the SLM model inference function
247 and P_{desc} is simply “Please describe the given video.”.

248 This stage enables efficient large-scale video summariza-
249 tion by offloading heavy visual processing to a lightweight
250 model. Our segmenting strategy further enhances summa-
251 rization quality by improving temporal focus and reducing
252 information dilution across long sequences. Although these
253 SLMs lack advanced reasoning capabilities, the resulting
254 descriptions capture salient scene elements and temporal
255 context, providing a compact yet informative representation
256 for the second-stage reasoning model.

Stage 2: LLM-based Reasoning and Scoring. In the sec-
257 ond stage, a large language model (LLM) performs high-
258 level reasoning over the textual scenario descriptions gener-
259 ated by the SLM. Each description T_i is paired with a rea-
260 soning scoring prompt \hat{P}_{score} (almost identical to P_{score} ,
261 see appendix for details), instructing the model to evaluate
262 whether the segment depicts a rare or anomalous driving
263 event. The LLM produces a relevance score s_i represent-
264 ing the likelihood of the segment belonging to the long-tail
265 distribution:
266

$$267 \quad s_i = f_{\text{LLM}}(T_i, \hat{P}_{\text{score}}), \quad (4)$$

268 where $f_{\text{LLM}}(\cdot)$ denotes the LLM. For each video, we ag-
269 gregate segment-level scores by taking the maximum score
270 across all segments, i.e. $s = \max_i s_i$. The videos are then
271 scored using the same thresholding strategy as in the vanilla
272 setup, where videos with $s \geq \tau$ are identified as anomaly.
273

274 By reasoning directly over textual scenario descriptions
275 that capture key events and contextual details, this stage en-
276 ables efficient large-scale inference without processing raw
277 videos.

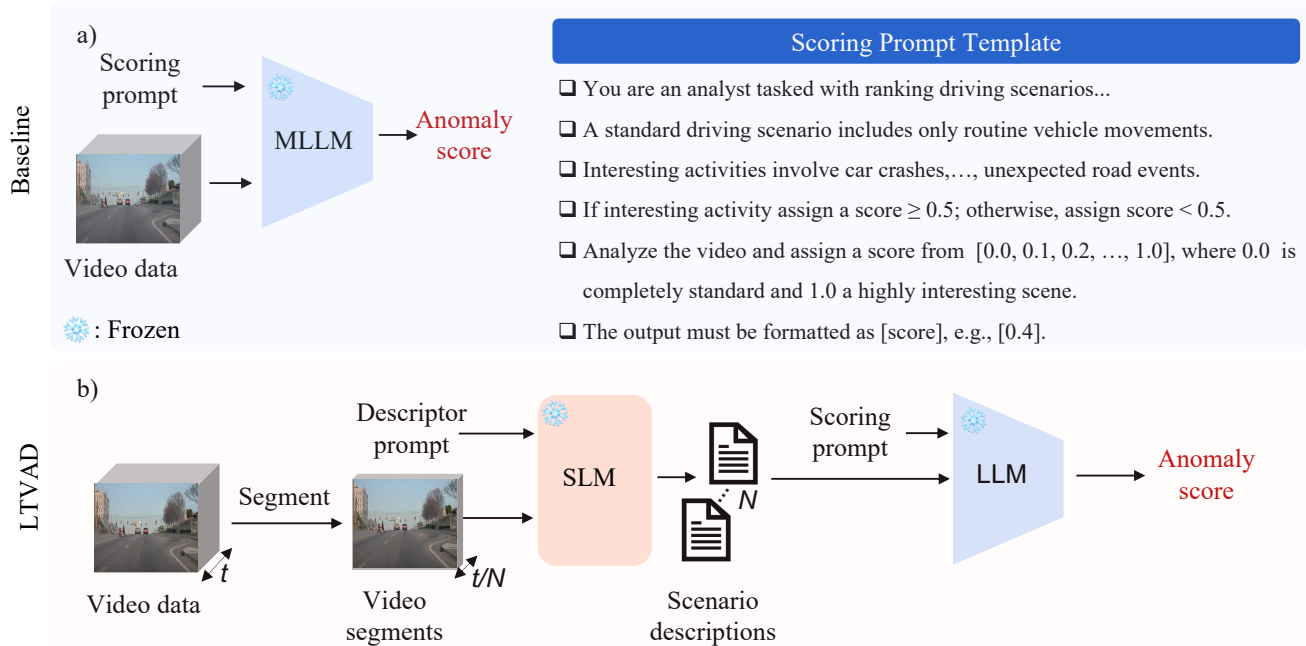


Figure 1. **Overview of proposed framework.** a) Baseline set-up for long-tail video anomaly detection. b) Proposed two stage framework, where first stage temporally segments the video into smaller lengths and perform scenario description using SLMs; second stage employs an LLM to perform anomaly scoring, followed by max pooling over segments for the final anomaly score.

277

4. Experiments and Results

278

4.1. Dataset

279

Given the lack of publicly available datasets for video anomaly detection in autonomous driving, we develop a detailed labeling framework that combines model-based pseudo-labeling with subsequent manual review.

280

281

282

283

284

285

286

287

288

Data Sources and Preprocessing. Our dataset, illustrated in Fig. 2, is collected from real-world driving under diverse conditions—ranging from day to night and across varying weather—using front-view camera recordings. Raw videos are 30 second each, recorded at a frame rate of 30 Hz and a resolution of 1220×1936 pixels.

289

290

291

292

293

294

295

Model-Based Labeling. To determine whether a clip contains an interesting (abnormal) scenario, we adopt a model-based labeling approach. Each 30-second MP4 clip is analyzed using Qwen2.5-VL-72B-Instruct, guided by a structured prompt defining a taxonomy of ten high-level categories of rare or safety-critical driving situations. The model is instructed to:

296

297

298

299

300

301

302

303

1. Output Scenario Detected: [Category] - [Subtype] only when highly confident, or
 2. Return "no" if no scenario is detected or confidence is insufficient.
- This conservative instruction ensures that only high-confidence detections are retained, effectively reducing label noise from uncertain cases.

Category Taxonomy. The prompt defines the following

categories and example subtypes:

304

1. Unpredictable Pedestrian Behavior: jaywalking, pedestrians revealed from occlusion, children running, impaired pedestrians. 305
2. Unusual Vehicle Actions: abrupt stops or lane changes, wrong-way driving, emergency vehicles with active sirens. 306
3. Complex Traffic Interactions: complex roundabouts, non-standard intersections. 307
4. Road Infrastructure Anomalies: active construction zones, unlit or malfunctioning signals, detours, degraded lane markings. 308
5. Obstacles and Debris: fallen cargo, road debris, cones, barriers, or temporary objects. 309
6. Unusual Road Users: tricyclists, e-scooters, agricultural vehicles, non-standard vehicles. 310
7. Animals on the Road: wildlife crossings, domestic pets, livestock on rural roads. 311
8. Adverse Weather Conditions: heavy rain, fog, flooding, snow, or intense sun glare. 312
9. Pickup/Drop-Off Challenges: congested curb zones, long passenger wait times. 313
10. Emergency Situations and Environmental Extremes: active accident scenes, authority-imposed road closures, blocked traffic. 314

Each 30-second clip is processed independently. All clips labeled as "interesting" are then manually reviewed

329

330

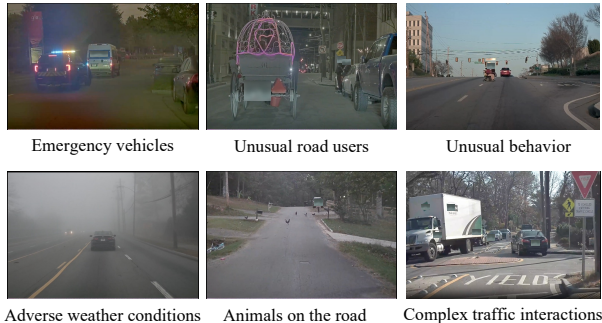


Figure 2. Examples from the category taxonomy.

to confirm correctness and remove false positives. A subset of clips labeled as “no” is also audited to estimate false negatives and refine the model threshold if necessary.

Dataset Composition. Following established practices in video anomaly detection [18, 25], we construct a validation dataset of 300 video clips from a large corpus of aforementioned preprocessed driving videos, tailored for anomaly detection in autonomous driving. We note that training based models requires a much larger annotated dataset, but such large training datasets are not needed for training-free methods as proposed LTVAD [25]. Consistent with [18], our proprietary dataset is balanced, containing 150 normal clips (label = 0) and 150 long-tail edge-case clips (label = 1). The edge-case subset in validation set covers all variety of rare or critical driving scenarios listed in the taxonomy, selected to achieve balanced coverage across categories. This ensures that the dataset reflects a diverse set of long-tail conditions relevant to autonomous driving safety.

4.2. Experimental Set-up and Metrics

Set-up. In this study, we leverage open-source pretrained MLLMs for our experiments. In particular, for stage 1 of our framework, we use Qwen2.5-VL-7B as small MLLM [1], and we employ Qwen2.5-72B-Instruct for stage 2. For baseline comparisons described in Sec. 3.2, we include Qwen2.5-VL-7B and InternVL3-8B as SLM baselines; Qwen2.5-VL-72B, and InternVL3-78B as large MLLM baselines. Given the training-free aspect of our methodology, we did not change any parameter such as fps (default 2) during inference. For scoring threshold, we use $\tau = 0.5$ to strike a balance between recall and precision. All experiments were conducted on four NVIDIA H100 GPUs (80 GB each) using the HuggingFace Transformers library with FlashAttention and BF16 precision.

Metrics. Following established practices in video anomaly detection, we report the Area Under the ROC (receiver operating characteristic) Curve (AUC) as our primary evaluation metric, which plots the true positive rate (TPR) against false positive rate (FPR) [4, 25]. We note that AUC is ag-

nostic to threshold value τ . In addition, we report accuracy, precision, recall, and F1-score for completeness. We further highlight recall (i.e., TPR) and F1 as auxiliary metrics alongside AUC, as they are particularly important for evaluating performance in long-tail anomaly detection scenarios for autonomous driving, where missing rare events can have severe consequences.

4.3. Experimental Results

Main results. Tab. 1 shows the quantitative comparisons between proposed LTVAD-AD with counterpart SLMs and large MLLM models. The experimental results show that proposed framework significantly outperforms counterpart SLM models (Qwen2.5-VL-7B) with a 24% improvement in AUC and further improvements in recall and F-1 score. Furthermore, proposed method able to achieve similar performance (slightly better) in terms of both AUC and Recall compared to Qwen2.5-VL-72B variant without any bells and whistles. More importantly, in terms of computational efficiency, our method achieves a substantial improvement over large MLLMs. In particular, the average processing time for a single 30-sec video clip using a large MLLM (Qwen2.5-VL-72B) is ~ 22 s, whereas our proposed two-stage framework processes the same input in ~ 1.5 s, resulting in a $\sim 14\times$ gain. This efficiency gain arises as the first stage offloads heavy visual processing to a lightweight SLM for video understanding, while the second stage performs text summarization in large batches without incurring significant memory overhead.

Method	Accuracy	Precision	Recall	F1	AUC
InternVL3-8B	0.69	0.83	0.48	0.61	0.78
InternVL3-78B	0.9	0.89	0.91	0.9	0.92
Qwen2.5-VL-7B	0.74	0.89	0.54	0.67	0.75
Qwen2.5-VL-72B	0.88	0.84	0.95	0.89	0.92
LTVAD-AD	0.85	0.79	0.96	0.87	0.93

Table 1. Quantitative comparison of the proposed LTVAD-AD model with counterpart SLM and MLLM methods.

Ablation study. In this section, we investigate the influence key components in our proposed LTVAD-AD framework, namely temporally segmenting videos, including ranking guidance and in-context examples in the scoring prompt.

Tab. 2 shows ablation for various segmenting lengths. In this study, raw videos (30 sec) are temporally segmented into smaller non-overlapping fixed size segments. Temporally segmenting video size by 2 provides improvement over no segmenting. Further segmenting decreases the AUC and F-1 scores.

Tab. 3 shows the comparison between scoring prompts with and without enforced guidance for ranking videos as normal or anomaly. In particular, we ablate the use of “If an interesting activity is present, make sure to assign a score ≥ 0.5 ; otherwise assign a score < 0.5 ” in the prompt. Without

Method	Accuracy	Precision	Recall	F1	AUC
no segmenting	0.85	0.84	0.88	0.86	0.90
segmenting/2	0.85	0.79	0.96	0.87	0.93
segmenting/3	0.76	0.69	0.94	0.79	0.89

Table 2. Ablation study on the impact of video length on long-tail detection for the proposed framework.

413 using this explicit phrase in the prompt, the LLMs tend to
 414 be conservative in ranking videos as anomalies, resulting
 415 in a huge performance drop in recall and F-1 score. Thus,
 416 LLMs needs explicit guidance in correctly ranking videos.

Method	Accuracy	Precision	Recall	F1	AUC
w/o ranking guidance	0.64	0.96	0.29	0.45	0.88
w ranking guidance	0.85	0.84	0.88	0.86	0.90

Table 3. Ablation study on the impact of explicit ranking prompt on long-tail video anomaly detection.

417 Tab. 4 shows the ablation study on the inclusion & exclu-
 418 sion of in-context examples for interesting activities (e.g.,
 419 car crashes). Exclusion of in-context examples does not
 420 lead to a major performance difference highlighting the
 421 LLMs reasoning capability under an explicit scoring guid-
 422 ance.

Method	Accuracy	Precision	Recall	F1	AUC
w/o in-context examples	0.74	0.68	0.91	0.78	0.88
w in-context examples	0.85	0.84	0.88	0.86	0.90

Table 4. Ablation study on the impact of in-context examples in the prompt on long-tail video detection.

423 **Distillation.** SLMs are not specifically trained on au-
 424 tonomous driving data, which can lead to hallucinated sce-
 425 nario descriptions and suboptimal performance in down-
 426 stream anomaly detection. To address this limitation, we
 427 perform supervised fine-tuning of the SLMs (Qwen2.5-
 428 VL-7B) on an autonomous driving dataset to improve do-
 429 main alignment and description quality. For this study, we
 430 adopted a knowledge distillation approach [27] (illustrated
 431 in Fig. 3), where we first employ a stronger multimodal
 432 teacher model (Qwen2.5-VL-72B) to generate scenario de-
 433 scriptions for each training video. Subsequently, these gen-
 434 erated descriptions are treated as pseudo-labels and used to
 435 fine-tune the SLMs. To preserve the generalization ability
 436 of the base model, we adopt parameter-efficient fine-tuning
 437 with LoRA (rank=64, $\alpha=64$), optimizing a standard autore-
 438 gressive language modeling loss [17]. The model is fine-
 439 tuned for 2 epochs on 5K video-description pairs. To illus-
 440 trate the efficiency of the fine-tuning, we rerun the segmen-
 441 ting/3 experiment in Tab. 2 with a fine-tuned SLM back-
 442 bone. Experimental results on Tab. 5 show that, fine-tuning
 443 SLM backbone enhance both accuracy and precision, while
 444 also showing improvements for F1 and AUC scores.

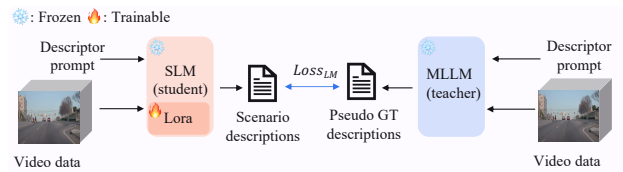


Figure 3. Distillation framework for scenario description generation. Pseudo ground truth (GT) descriptions from a teacher model is used for supervised fine-tuning of student SLM model, with objective of minimizing language modeling loss.

Method	Accuracy	Precision	Recall	F1	AUC
No Fine-tuning	0.76	0.69	0.94	0.79	0.89
Fine-tuning	0.81	0.75	0.93	0.83	0.91

Table 5. Quantitative comparison of SLM fine tuning vs without fine tuning for long-tail video anomaly detection.

5. Limitations

445 The presented framework currently only focuses on front-
 446 view cameras as front-view contain the most vital infor-
 447 mation for driving [7]. Given fleet-data generally collect
 448 multi-view data, additional views might positively help for
 449 better mining long-tail scenarios. However, a concern here
 450 is to find a balance as non-front-view cameras might also in-
 451 troduce large amount of false positives as they do not impact
 452 driving as much of the front-views. Another potential pit-
 453 fall is the lack of explicit motion information extraction and
 454 incorporation into scenario descriptions that can be used for
 455 long-tail video anomaly detection [6]. This can be espe-
 456 cially useful for understanding nuanced scenarios such as
 457 aggressive cut-ins. Lastly, in this study, we adopted a uni-
 458 form sampling for frame selection. However, a lightweight
 459 method for selecting key frames that capture critical mo-
 460 ments can further enhance downstream long-tail detection
 461 performance. We aim to explore these directions as part of
 462 our future work. 463

6. Conclusion

464 In this study, we presented LTVAD-AD, a novel long-
 465 tail video detection framework for autonomous driving. 466
 467 LTVAD-AD consists of two stages, where first stage lever-
 468 ages small language models for scenenario summarization
 469 and second stage employs an LLM to perform anomaly de-
 470 tection on scenario summarization from the first stage. Ex-
 471 tensive experiments on curated driving dataset shows the
 472 efficiency and accuracy of proposed method in detecting
 473 long-tail driving videos over counterpart SLM & MLLM
 474 baselines. The proposed approach is versatile as informa-
 475 tion from first stage is generic and can be used various other
 476 tasks such as embedding and retrieval tasks.

477

References

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

- [1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 3, 5
- [2] Sagar Behere and Martin Törngren. A functional architecture for autonomous driving. In *Proceedings of the first international workshop on automotive software architecture*, pages 3–10, 2015. 1
- [3] Jan-Aike Bolte, Andreas Bar, Daniel Lipinski, and Tim Fingscheidt. Towards corner case detection for autonomous driving. In *2019 IEEE Intelligent vehicles symposium (IV)*, pages 438–445. IEEE, 2019. 1, 2
- [4] Jasmin Breitenstein, Jan-Aike Termöhlen, Daniel Lipinski, and Tim Fingscheidt. Corner cases for visual perception in automated driving: Some guidance on detection approaches. *arXiv preprint arXiv:2102.05897*, 2021. 1, 3, 5
- [5] Haoyu Fu, Diankun Zhang, Zongchuang Zhao, Jianfeng Cui, Dingkan Liang, Chong Zhang, Dingyuan Zhang, Hongwei Xie, Bing Wang, and Xiang Bai. Orion: A holistic end-to-end autonomous driving framework by vision-language instructed action generation. *arXiv preprint arXiv:2503.19755*, 2025. 1, 2
- [6] Yulu Gan, Ligeng Zhu, Dandan Shan, Baifeng Shi, Hongxu Yin, Boris Ivanovic, Song Han, Trevor Darrell, Jitendra Malik, Marco Pavone, et al. Foundationmotion: Auto-labeling and reasoning about spatial movement in videos. *arXiv preprint arXiv:2512.10927*, 2025. 6
- [7] Shenyan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta, Yihang Qiu, Andreas Geiger, Jun Zhang, and Hongyang Li. Vista: A generalizable driving world model with high fidelity and versatile controllability. In *Advances in Neural Information Processing Systems*, 2024. 6
- [8] Mahmudul Hasan, Jonghyun Choi, Jan Neumann, Amit K Roy-Chowdhury, and Larry S Davis. Learning temporal regularity in video sequences. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 733–742, 2016. 1, 2
- [9] Deepti Hegde, Rajeev Yasarla, Hong Cai, Shizhong Han, Apratim Bhattacharyya, Shweta Mahajan, Litian Liu, Risheek Garrepalli, Vishal M Patel, and Fatih Porikli. Distilling multi-modal large language models for autonomous driving. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 27575–27585, 2025. 2
- [10] Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning. In *European Conference on Computer Vision*, pages 533–549. Springer, 2022. 1
- [11] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhao Wang, et al. Planning-oriented autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17853–17862, 2023. 1, 2
- [12] Jyh-Jing Hwang, Runsheng Xu, Hubert Lin, Wei-Chih Hung, Jingwei Ji, Kristy Choi, Di Huang, Tong He, Paul Covington, Benjamin Sapp, et al. Emma: End-to-end multimodal model for autonomous driving. *arXiv preprint arXiv:2410.23262*, 2024. 1, 2
- [13] Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggong Wang. Vad: Vectorized scene representation for efficient autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8350, 2023. 1, 2
- [14] Sicong Jiang, Zilin Huang, Kangan Qian, Ziang Luo, Tianze Zhu, Yang Zhong, Yihong Tang, Menglin Kong, Yunlong Wang, Siwen Jiao, et al. A survey on vision-language-action models for autonomous driving. *arXiv preprint arXiv:2506.24044*, 2025. 1
- [15] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017. 1
- [16] Bencheng Liao, Shaoyu Chen, Haoran Yin, Bo Jiang, Cheng Wang, Sixu Yan, Xinbang Zhang, Xiangyu Li, Ying Zhang, Qian Zhang, et al. Diffusiondrive: Truncated diffusion model for end-to-end autonomous driving. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12037–12047, 2025. 2
- [17] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 6
- [18] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488, 2018. 5
- [19] Wenchao Sun, Xuewu Lin, Yining Shi, Chuang Zhang, Haoran Wu, and Sifa Zheng. Sparsedrive: End-to-end autonomous driving via sparse scene representation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8795–8801. IEEE, 2025. 1, 2
- [20] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang, Zhiyong Zhao, Kun Zhan, Peng Jia, XianPeng Lang, and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. In *Conference on Robot Learning*, pages 4698–4726. PMLR, 2025. 2
- [21] Yan Wang, Wenjie Luo, Junjie Bai, Yulong Cao, Tong Che, Ke Chen, Yuxiao Chen, Jenna Diamond, Yifan Ding, Wenhao Ding, et al. Alpamayo-r1: Bridging reasoning and action prediction for generalizable autonomous driving in the long tail. *arXiv preprint arXiv:2511.00088*, 2025. 1
- [22] Xinshuo Weng, Boris Ivanovic, Yan Wang, Yue Wang, and Marco Pavone. Para-drive: Parallelized architecture for real-time autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15449–15458, 2024. 1, 2
- [23] Runsheng Xu, Hubert Lin, Wonseok Jeon, Hao Feng, Yuliang Zou, Liting Sun, John Gorman, Kate Tolstaya, Sarah Tang, Brandyn White, et al. Wod-e2e: Waymo open dataset for end-to-end driving in challenging long-tail scenarios. *arXiv preprint arXiv:2510.26125*, 2025. 1, 2
- [24] Burhan Yaman, Yunsheng Ma, and Xin Ye. A peek into tesla’s autonomous future: Core tech revealed by vp ashok

- 591 elluswamy at iccv 2025 wdfm-ad. *Medium*, 2025. Accessed:
592 2025-11-06. [1](#), [2](#)
- 593 [25] Luca Zanella, Willi Menapace, Massimiliano Mancini, Yim-
594 ing Wang, and Elisa Ricci. Harnessing large language mod-
595 els for training-free video anomaly detection. In *Proceedings*
596 *of the IEEE/CVF Conference on Computer Vision and Pat-*
597 *tern Recognition*, pages 18527–18536, 2024. [5](#)
- 598 [26] Jingyuan Zhao, Yuyan Wu, Rui Deng, Susu Xu, Jinpeng
599 Gao, and Andrew Burke. A survey of autonomous driving
600 from a deep learning perspective. *ACM Computing Surveys*,
601 57(10):1–60, 2025. [1](#)
- 602 [27] Yingxue Zhou, Tan Zhu, Tao Zeng, Zigeng Wang, and Wei
603 Shen. Learning from llm agents: In-context generative mod-
604 els for text casing in e-commerce ads. In *Proceedings of*
605 *the 2025 Conference on Empirical Methods in Natural Lan-*
606 *guage Processing: Industry Track*, pages 1122–1133, 2025.
607 [6](#)
- 608 [28] Zewei Zhou, Tianhui Cai, Seth Z Zhao, Yun Zhang, Zhiyu
609 Huang, Bolei Zhou, and Jiaqi Ma. Autovla: A vision-
610 language-action model for end-to-end autonomous driving
611 with adaptive reasoning and reinforcement fine-tuning. *arXiv*
612 *preprint arXiv:2506.13757*, 2025. [1](#), [2](#)

613 **A. Appendix**

614 **Prompts for LTVAD-AD.** The complete scoring prompt
615 (P_{score}) used for baselines in our study:

616 You are an analyst tasked with
617 ranking driving scenarios based
618 on the presence of interesting
619 activities.
620 A standard driving scenario
621 includes only routine vehicle
622 movements.
623 Interesting activities
624 involve car crashes, emergency
625 vehicles, vulnerable road users,
626 unusual behavior of vehicles,
627 adverse weather conditions,
628 temporary road hazards, active
629 construction signs, or any other
630 unexpected road events.
631 Analyze the given video and
632 assign a score from the
633 following fixed list: [0.0,
634 0.1, 0.2, 0.3, 0.4, 0.5, 0.6,
635 0.7, 0.8, 0.9, 1.0], where 0.0
636 means the scene is completely
637 standard and 1.0 indicates a
638 highly interesting scene.
639 If any of the interesting
640 activity is present, make sure
641 to assign a high score (≥ 0.5);
642 otherwise, assign a low score
643 (< 0.5).
644 The output must be only one
645 number from the fixed list and
646 formatted as [score], e.g.,
647 [0.4].

648 In our proposed framework, the first stage description
649 prompt (P_{desc}) was the following:

650 Please describe the given video.

651 The second-stage reasoning scoring prompt (\hat{P}_{score}) was
652 identical to P_{score} , with the following modification to incor-
653 porate the scenario summarization from the first stage:

654 ...

655 The video has been summarized as
656 follows: `first-stage scenario`
657 `summarization`.
658 Analyze the video summary,
659 and assign a score from the
660 following fixed list..
661 ...