

Submitted to Bridging Transportation Researchers Conference for presentation in August 2025.

VIDEO-BASED VEHICLE SURVEILLANCE IN THE WILD: LICENSE PLATE, MAKE, AND MODEL RECOGNITION USING VISION-LANGUAGE MODELS

Pouya Parsa

Department of Civil, Environmental, and Geo- Engineering
University of Minnesota
parsa025@umn.edu

Keya Li

Department of Civil, Architectural and Environmental Engineering
The University of Texas at Austin
keya_li@utexas.edu

Kara M. Kockelman, Ph.D., P.E.

Department of Civil, Architectural and Environmental Engineering
The University of Texas at Austin
kkockelm@mail.utexas.edu

Seongjin Choi, Ph.D.

Department of Civil, Environmental, and Geo- Engineering
University of Minnesota
chois@umn.edu

Word Count: 3968 words + 4 table(s) \times 250 = 4968 words

Submission Date: May 23, 2025

ABSTRACT

Automatic license plate recognition (ALPR) and vehicle recognition are critical components in traffic law enforcement, parking management, toll collection, crime-solving and more. Performing ALPR on video captured by hand-held smartphones or dashboard cameras (dashcams) introduces unique challenges compared to fixed traffic cameras, including camera lens motions, suboptimal viewpoints, occlusions, and the absence of prior knowledge of road geometry and location. Traditional ALPR systems rely on specialized camera hardware, hand-crafted OCR pipelines, and tightly coupled sub-modules; so they degrade significantly when faced with small, blurred, or partially occluded plates. Recent advances in large vision–language models (VLMs) are better able to recognize free-form text and semantics directly from arbitrary imagery within a single network.

This paper presents the first comprehensive study applying off-the-shelf VLMs to monocular (single-lens) video captured by both fixed cameras and handheld smartphones, for two parallel tasks: license plate recognition and vehicle make + model recognition. Although often co-deployed, these tasks pose distinct challenges and do not need to be solved jointly. Our end-to-end framework detects candidate vehicles and plates using YOLOv8, ranks frames via a CLIP-based perceptual quality score, composites top frames into rich visual prompts, and queries VLMs through tailored prompt engineering. Experiments on an internal 24-video smartphone benchmark and the public UFPR-ALPR dataset demonstrate that our method achieves 91.67% plate-level top-1 accuracy, nearly tripling the previous 29.7% top-10 baseline. For vehicle make and model recognition, the model attains 70.8% top-1 accuracy compared to 16.9% previously reported. This paper shows how open-source Llama-3.2-Vision (11B) matches GPT-4o in performance with zero API cost, enabling full on-device deployment.

Keywords: automatic license plate recognition, vision–language models, CLIP image quality assessment, prompt engineering, traffic law enforcement, automated enforcement

INTRODUCTION

Nearly 1.2 million people die each year globally in road traffic crashes(1), with over 40,000 fatalities occurring in the United States alone in 2023 (2). Speeding, reckless driving, and hit-and-run collisions account for many of these fatalities, yet the lack of real-time enforcement tools makes preventing them especially difficult. Stationary cameras, while effective, are very expensive to install and maintain—averaging over \$120,000 per location (3). Their fixed placement limits coverage, and many cities, states and nations are reluctant to use them. Moreover, drivers often slow down near enforcement points (and speed up downstream, once out of camera view), a phenomenon known as the “kangaroo effect” (4). A scalable alternative lies in enabling ordinary citizens to assist enforcement by capturing video footage of infractions using smartphones or dashcams. However, leveraging such crowd-sourced data requires robust tools for analyzing low-quality, personal videos, which is a significant technical challenge.

Several nations have developed public reporting platforms to facilitate citizen-driven traffic enforcement. For example, in South Korea, citizens can use the official “Safety e-Report” service¹ to report traffic violations by uploading video evidence directly to law enforcement. In return, some users may receive monetary rewards or recognition, effectively turning smartphones into distributed enforcement tools and encouraging community participation in road safety (5, 6). In

¹<https://www.safetyreport.go.kr/eng/>



FIGURE 1: Comparison between an ideal high-resolution license plate (left plate) and four low-quality plates from real-world footage (right).

New York City, citizens have been helping enforce diesel-truck idling laws; those submitting 3 minutes of video receive 25 percent of any fine obtained from heavy-truck owners, which is close to \$87.50 (7). Such low-cost programs cost-effectively extend enforcement reach while building social accountability among drivers.

Despite their success, these systems still require users to manually enter important vehicle information such as the license plate number, make, and model. In many cases, this information is added after the video upload is complete, or left out entirely when the plate is difficult to read. This extra step increases the burden on users and reduces the number of reports that can be successfully verified. Automating the extraction of this information directly from the video would improve ease of use, minimize input errors, and increase the likelihood that valid reports result in enforceable actions.

To enable automatic traffic law enforcement from citizen-recorded videos, two core recognition tasks are essential: Automated license plate recognition (ALPR) and vehicle make (manufacturer) and model recognition. These tasks serve as the foundation for identifying and tracking offending vehicles across time and locations. ALPR allows for the extraction of unique vehicle identifiers, enabling citation issuance, database cross-referencing, and ownership tracing. Recognizing the make and model of a vehicle provides an additional layer of verification or confirmation, and is especially valuable in cases of partial plate visibility, occlusion, or modified plates. Together, these two tasks form the minimum viable input required for automating enforcement actions from unstructured video evidence.

ALPR, also known as Automatic Number Plate Recognition (ANPR), is a foundational intelligent transportation system (ITS) technology for automated enforcement, enabling identification of vehicles involved in traffic violations, toll evasion, or criminal activities. It can play a critical role in issuing citations, locating stolen vehicles, and monitoring access points (to paid parking lots, tollways, and high-security events). However, traditional ALPR systems depend on high-resolution cameras and controlled environments, making them unreliable when applied to noisy footage from standard smartphones or legacy closed-circuit television (CCTV) camera systems. In such cases, license plates may occupy only a few pixels or appear blurred or occluded, reducing the effectiveness of conventional optical character recognition (OCR) systems. An illustrative comparison is provided in Figure 1, contrasting an ideal high-resolution US license plate with a low-resolution, blurry plate often seen in real-world video footage.

Beyond plate recognition, extracting vehicle make and model further enhances enforcement capabilities by supporting cross-verification, vehicle re-identification, and suspect profiling. But this task typically relies on separate classifiers that are sensitive to viewpoint variations and trained on limited datasets, making them brittle and hard to scale.

Conventional approaches (8–10), for ALPR typically rely on dedicated OCR engines applied to tightly cropped plate regions, while make and model recognition is handled by specialized

CNN-based classifiers trained on curated vehicle datasets (11), (12). These pipelines often assume clean, high-resolution imagery and controlled viewing conditions, which rarely hold in citizen-sourced or dashcam videos. In such scenarios, low resolution, motion blur, obstructions, and unusual angles significantly degrade the performance of traditional classification models. Moreover, maintaining and deploying multiple task-specific models increases system complexity and makes real-time, on-device processing difficult. A practical system therefore needs to be task-agnostic. In other words, it should take an entire video sequence as input and, without relying on separate specialist modules, return the plate text together with make and model.

Vision-language models (VLMs) meet these requirements by training on vast collections of image-caption pairs so that one network can handle both pictures and words. Inside the model, the visual branch converts pixels into internal features that the language branch can immediately read, allowing the system to recognise objects, decipher text, and answer questions without switching between separate modules. In practice, visual understanding means spotting a vehicle and reading a partly blurred plate; reasoning means combining those observations to respond to a prompt such as “What make and model is the car, and what does the plate say?” Because prompts steer the model at run-time, the same network can tackle many tasks with little or no extra training, eliminating the need for standalone OCR or make-model classifiers and making VLMs well suited to real-time traffic-enforcement video.

Specifically, this study evaluates four state-of-the-art VLMs: GPT-4o (13), Llama 3.2-Vision (14), LLaVA (15), and MiniCPM-V (16). With carefully engineered prompts, each model is tasked with inferring vehicle attributes and deciphering license-plate characters, even under low-resolution, blurry, or partially occluded conditions. A primary research question as follows: *Can recent VLMs replace OCR and semantic classifiers in a single unified module, operating directly on crowd-sourced video footage?* This paper answers in the affirmative while providing: (1) an end-to-end pipeline combining object detection, frame-quality ranking, and VLM prompting; (2) the first extensive benchmark of multiple open and proprietary VLMs on plate and make + model recognition, including a cost-accuracy analysis; (3) evidence that an inexpensive 11-billion-parameter VLM matches GPT-4o performance (accuracy) while running locally in real time; and (4) ablation and error analysis guiding future ALPR research.

RELATED RESEARCH

Traditional ALPR systems typically employ a multistage pipeline that involves license plate detection, character segmentation, and character recognition, often relying on handcrafted features and optical character recognition (OCR) techniques (17). These methods, while effective in controlled environments with high-resolution cameras and optimal lighting, tend to degrade in performance when faced with challenges such as motion blur, occlusions, varying viewpoints, and low-quality images, which are prevalent in video captured by handheld devices like smartphones.

Recent advancements in deep learning have significantly improved ALPR robustness. For example, LPRNet (18). LPRNet Paper introduces a lightweight convolutional neural network (CNN) that performs end-to-end license plate recognition, achieving real-time performance with high accuracy on standard datasets. Despite these advances, deep learning-based ALPR systems can still struggle with the variability and noise inherent in unconstrained environments, such as those encountered in smartphone-captured video.

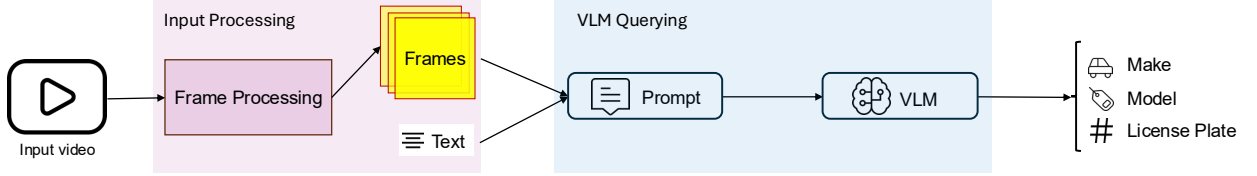


FIGURE 2: Overview of the vision-language pipeline for vehicle attribute extraction. High-quality frames and text prompts are combined to guide the model in predicting vehicle make, model, and license plate information.

Vehicle Make and Model Recognition

Parallel to ALPR, vehicle make and model recognition has gained traction in applications like traffic monitoring, surveillance, and autonomous driving. Deep learning approaches have been particularly effective in this domain. For instance, Deep Learning Vehicle Classification (19) employs CNNs to classify vehicles based on visual features, demonstrating high accuracy on fine-grained vehicle datasets. These systems, however, typically require large annotated datasets and may not generalize well to diverse real-world conditions, such as varying lighting or partial occlusions.

Vision-Language Models (VLMs)

The advent of large vision-language models (VLMs) has introduced a transformative approach to vision tasks by leveraging joint understanding of images and text. VLMs, such as CLIP (20). CLIP Paper, GPT-4o, and LLaVA, are pre-trained on extensive datasets of image-text pairs, enabling them to perform a wide range of tasks, including zero-shot classification and text recognition, without task-specific training. These models excel in understanding complex scenes and can generate textual descriptions from visual inputs, making them well-suited for tasks like ALPR and vehicle recognition.

Recent studies have explored VLMs for optical character recognition (OCR) tasks in dynamic environments. For example, Benchmarking VLMs for OCR evaluates the performance of VLMs like Claude-3, Gemini-1.5, and GPT-4o on OCR in video frames across diverse domains, such as code editors and advertisements. The study highlights VLMs' potential to outperform traditional OCR systems like EasyOCR and RapidOCR in complex scenarios, although challenges like hallucinations and sensitivity to stylized text persist.

Specifically for ALPR, Advancing Vehicle Plate Recognition (21) demonstrates the application of vision-language models (VLMs) to recognize license plates under challenging conditions such as low illumination, motion blur, and tightly packed characters. The authors fine-tune a VLM (PaliGemma) to develop VehiclePaliGemma, achieving 87.6% accuracy on a Malaysian license plate dataset. Their method leverages multitask prompting to identify plates in complex scenes involving multiple vehicles with different colors and models. While effective, their reliance on model fine-tuning introduces an additional barrier for adoption, as it requires access to training infrastructure and expertise. In contrast, our approach employs zero-shot prompting with off-the-shelf VLMs, avoiding the need for retraining and making it easier to adapt and deploy in diverse ALPR scenarios.

METHODOLOGY

Our proposed framework processes an input video to extract vehicle make, model, and license plate information through two main stages: (1) Input Processing and (2) VLM Querying, as illustrated in Figure 2. The Input Processing stage (pink-shaded box) involves selecting the most informative and high-quality frames from the video to enhance recognition accuracy while significantly reducing the input size to the VLM. In the VLM Querying stage (blue-shaded box), these selected frames are paired with carefully engineered textual prompts to form a unified multimodal input. This input is then passed to a VLM, which interprets both the visual and textual information to generate structured outputs such as make, model, and license plate.

Input Processing

The goal of the Input Processing stage is to select frames from the input video that are sharp, well-exposed, informative, and suitable for extracting vehicle information. For example, if a license plate appears in multiple frames, it suffices to use just one frame that best captures the full license plate number. This reduces the input size to the VLM by orders of magnitude compared to using all video frames. Moreover, selecting high-quality frames improves extraction accuracy since poor quality images may lead to errors such as misreading license plates or predicting incorrect vehicle types.

To identify the most informative frames from a video, we use two perceptual image quality metrics: CLIP-IQA (22) and BRISQUE (23). Each scoring method ranks frames based on perceived quality, allowing us to prioritize frames likely to yield accurate recognition results.

CLIP-IQA is a no-reference image quality assessment metric that leverages CLIP embeddings. It computes similarity between an image’s CLIP embedding and a reference quality embedding (e.g., “a high-quality photo” or “a blurry image”) as:

$$\text{CLIP-IQA}(x) = \cos(f_{\text{CLIP}}(x), f_{\text{CLIP}}(t_{\text{ref}})) \quad (1)$$

where $f_{\text{CLIP}}(\cdot)$ denotes the CLIP embedding function and t_{ref} is a textual prompt representing good or poor quality. This score captures perceptual similarity and correlates well with human judgments of image quality (22).

BRISQUE (Blind/Referenceless Image Spatial Quality Evaluator) is a handcrafted feature-based method that uses natural scene statistics (NSS) in the spatial domain to predict perceived quality. It models image distortions by extracting statistical features from locally normalized luminance coefficients and fits them into a support vector regressor trained on human opinion scores. BRISQUE is effective in assessing common distortions such as noise, blur, and compression artifacts without requiring a reference image (23).

Vision-Language Model Querying

To extract structured information from selected video frames, we first engineer textual prompts that specify the attributes or identifiers to be extracted—such as vehicle make, model, or license plate. This prompt design ensures consistency and clarity across queries. These prompts are then paired with representative video frames, previously selected for their quality and relevance, to form a unified multimodal input. This composite input is fed into a VLM, such as GPT-4o (13) or Llama3.2-vision (14), which jointly interprets the visual content and textual instructions. The VLM processes this input to produce structured outputs as well as optional natural language justifications, depending on the prompt specification.

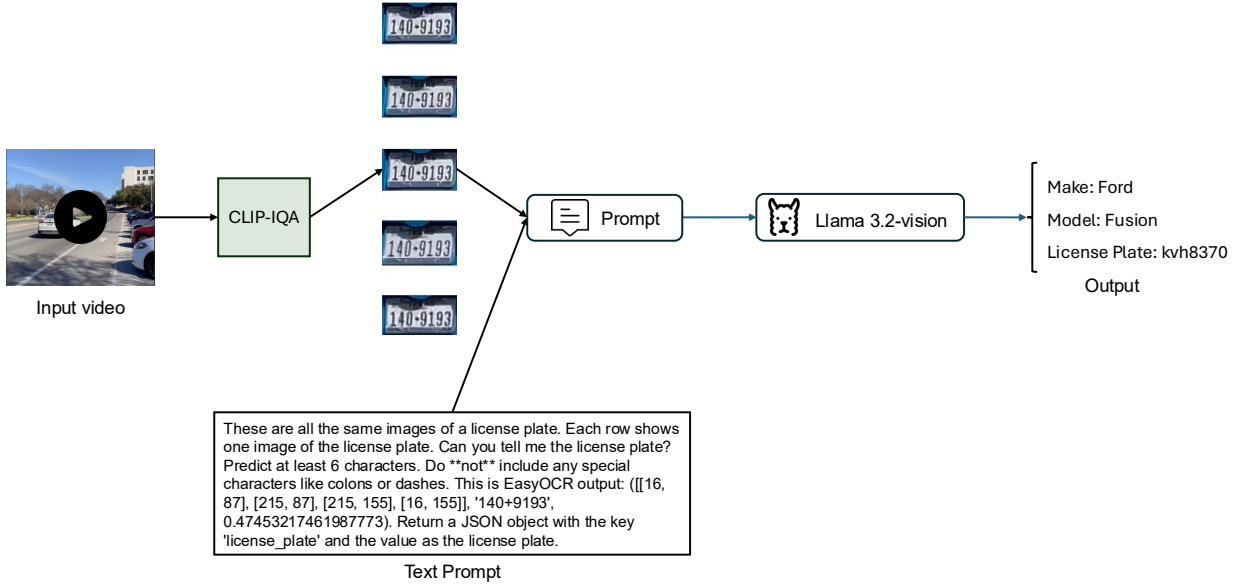


FIGURE 3: The figure illustrates a real-world example of the proposed framework, demonstrating how high-quality frames are selected, composited, and processed by a Vision-Language Model to generate structured predictions.

EXPERIMENTAL SETUP

We implemented our framework on an NVIDIA RTX 6000 Ada GPU using the YOLO model provided by Li et al. (12) which is fine-tuned for two epochs on approximately 2,500 frames from smartphone videos for both vehicle and license plate detection tasks.

Data

We evaluate our approach using the UFPR-ALPR dataset (24), which consists of 4,522 images featuring 150 unique vehicles captured in dynamic camera and vehicle movement scenarios. We utilize the dataset introduced by Li et al. (12), which includes 24 smartphone videos captured at 60 frames per second using an iPhone 12 from a street-side perspective, totaling approximately 12,300 frames. Each video is annotated with ground truth license plate strings and corresponding vehicle make and model labels, verified by two independent raters.

To illustrate the end-to-end pipeline, we provide a working example (figure 3) using the CLIP-IQA metric for frame selection. Given a short video clip of a vehicle in motion, individual frames are first scored by CLIP-IQA, and the top-ranked frames (e.g., those with sharp, unobstructed license plates) are selected and composited into a single image. This composite is then paired with a text prompt (as shown in Figure 3):

This multimodal input is passed to a VLM, specifically Llama 3.2 Vision, which interprets both the composite image and the prompt to extract structured predictions. The model responds with multiple plausible candidates.

Baselines

We explore 24 different configurations by combining methods from three key components of our pipeline: two frame selection techniques (CLIP-IQA and BRISQUE), four VLMs (GPT-4o (13), Llama3.2-vision (14), LLaVA (15), and MiniCPM-V (16)), and three distinct prompting strategies (single call, three options, and three calls).

Frame Selection

To identify the most informative frames from each video, we use two no-reference image quality assessment metrics: CLIP-IQA and BRISQUE. These methods score each frame based on perceived visual quality, allowing us to prioritize those likely to improve recognition accuracy.

Vision-Language Model

For the recognition task, we evaluate four recent and widely available VLMs: GPT-4o, Llama-3.2-Vision, LLaVA, and MiniCPM-V. These models were selected due to their strong performance on multimodal benchmarks and accessibility through open-source release or API endpoints.

Prompting Strategies

We compare three different prompting strategies when querying VLMs:

- **Single Call:** the VLM is queried once to produce a single best guess.
- **Three Options:** the VLM is queried once and prompted to return its top three guesses in a single response.
- **Three Calls:** the VLM is queried three separate times with the same prompt, and a prediction is considered correct if any of the three outputs match the ground truth.

As a non-VLM baseline, we replicate the pipeline proposed by Li et al. (12), which uses a fine-tuned YOLO detector followed by EasyOCR for license plate recognition, with post-processing to select the most frequent prediction.

RESULTS

Table 1 reports zero-shot license-plate recognition accuracy on the UFPR-ALPR benchmark as a function of the frame-quality metric, VLM, and the prompting strategy.

Independent querying is consistently beneficial, but its impact varies by model capacity. For the strongest model, GPT-4o, accuracy rises from 83.1% with a single call to 86.4% with either three options or three calls (using CLIP-IQA), a modest 3-point gain that effectively saturates performance. Mid-tier MiniCPM-V benefits markedly: CLIP-IQA accuracy nearly doubles, climbing from 28.8% (single call) to 54.2% (three calls). The weakest model, Llava, exhibits only limited recovery, never surpassing 22%, and even collapses to 0% under the three-option prompt, suggesting susceptibility to prompt-format changes when the underlying visual grounding is poor.

CLIP-IQA generally produces the highest scores for GPT-4o and MiniCPM-V, whereas BRISQUE matches or slightly exceeds CLIP-IQA for Llama 3.2-Vision under the three-option strategy (84.8% for both metrics). No single metric dominates across models, indicating that either perceptual proxy can serve as an effective pre-filter when paired with a capable VLM.

Table 2 contrasts the influence of the frame-quality metric, the VLM, and the prompting strategy on smartphone data. For the two strongest VLMs, i.e., GPT-4o and Llama 3.2-Vision, zero-shot inference already delivers competitive results (83.3% and 91.7%, respectively, when CLIP-IQA is used), but issuing three independent calls pushes accuracy to 91.7% and 91.7%,

TABLE 1: Accuracy of license plate recognition on UFPR-ALPR dataset across different VLM prediction strategies and image quality measures.

Metric	Model	Single Call	Three Options	Three Calls
CLIP-IQA	GPT-4o	83.05% (49/59)	86.44% (51/59)	86.44% (51/59)
	Llama3.2-vision	66.10% (39/59)	84.75% (50/59)	76.27% (45/59)
	Llava	16.95% (10/59)	00.00% (0/59)	20.34% (12/59)
	MiniCPM-V	28.81% (17/59)	50.85% (30/59)	54.24% (32/59)
BRISQUE	GPT-4o	77.97% (46/59)	84.75% (50/59)	84.75% (50/59)
	Llama3.2-vision	62.71% (37/59)	84.75% (50/59)	69.49% (41/59)
	Llava	22.03% (13/59)	00.00% (0/59)	22.03% (13/59)
	MiniCPM-V	42.37% (25/59)	44.07% (26/59)	50.85% (30/59)
<i>Baseline (Non-VLM): 46% accuracy reported by Li et al. (12)</i>				

TABLE 2: Accuracy of license plate recognition on smartphone dataset across different VLM prediction strategies and image quality measures.

Metric	Model	Single Call	Three Options	Three Calls
CLIP-IQA	GPT-4o	83.33% (20/24)	87.50% (21/24)	91.67% (22/24)
	Llama3.2-vision	91.67% (22/24)	91.67% (22/24)	87.50% (21/24)
	Llava	25.00% (6/24)	25.00% (6/24)	29.17% (7/24)
	MiniCPM-V	54.17% (13/24)	54.17% (13/24)	70.83% (17/24)
BRISQUE	GPT-4o	79.17% (19/24)	87.50% (21/24)	87.50% (19/24)
	Llama3.2-vision	87.50% (21/24)	87.50% (21/24)	91.67% (22/24)
	Llava	33.33% (8/24)	33.33% (8/24)	33.33% (8/24)
	MiniCPM-V	58.33% (14/24)	54.17% (13/24)	79.17% (19/24)
<i>Baseline (Non-VLM): 29.7% top-10 accuracy, reported by Li et al. (12)</i>				

reflecting absolute gains of 8–13 percentage points. The mid-tier MiniCPM-V benefits even more from repetitive calls, climbing from roughly 54–58% with one call to 71–79% after three calls, while the lowest-performing Llava remains largely insensitive to additional queries (≈ 25 –33%). Comparing the two frame-quality metrics, CLIP-IQA yields the best score for GPT-4o, whereas BRISQUE leads for Llama 3.2-Vision; thus, no single metric dominates, but either is sufficient to drive high performance when paired with a strong VLM.

Table 3 evaluates zero-shot VLM performance on make and make + model recognition under single- and three-calls prompting. For the easier *make-only* task, the strongest models (Llama 3.2-Vision and GPT-4o) already achieve high accuracy with a single call (75–79%), and issuing three independent calls yields modest absolute gains of 4–5 percentage points, topping out at

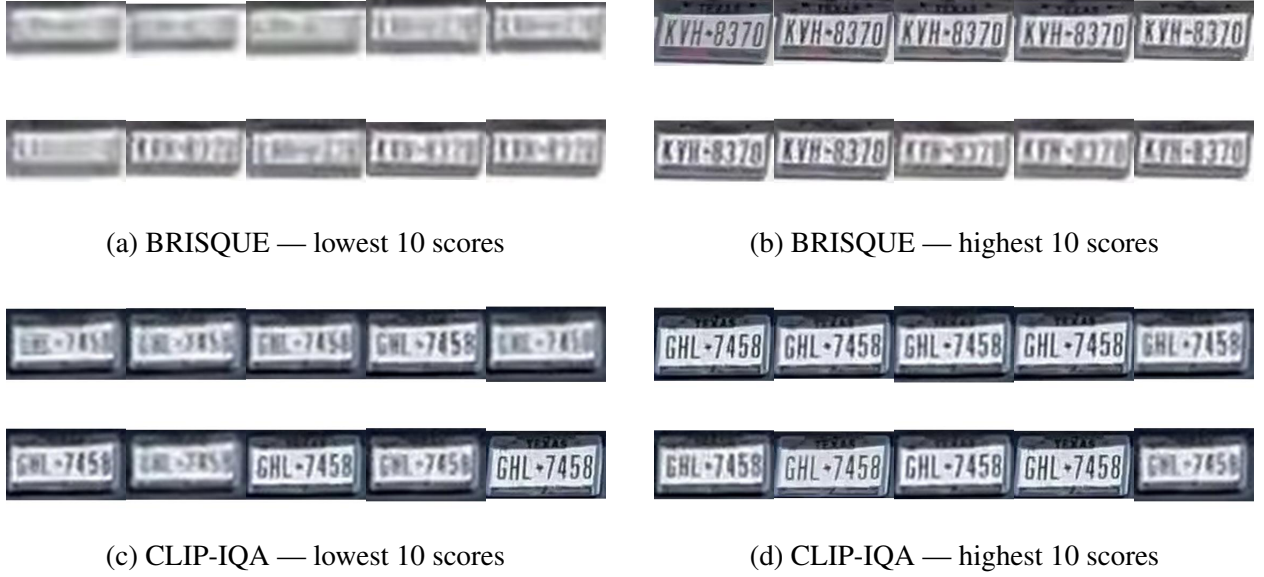


FIGURE 4: Quality extremes for two image-quality assessors. **Top row:** BRISQUE; **Bottom row:** CLIP-IQA. Each panel shows the 10 lowest- or highest-scoring frames (arranged as a 2×5 grid) for the same vehicle.

TABLE 3: Accuracy of make and make + model predictions on the smartphone dataset across different VLM strategies.

	Model	Single Call	Three Calls
Make	GPT-4o	75.00% (18/24)	79.17% (19/24)
	Llama-3.2-Vision	79.17% (19/24)	83.33% (20/24)
	Llava	58.33% (14/24)	58.33% (14/24)
	MiniCPM-V	75.00% (18/24)	79.17% (19/24)
	Baseline Li et al. (12)	48.60%	—
Make + Model	GPT-4o	66.67% (16/24)	70.83% (17/24)
	Llama-3.2-Vision	62.50% (15/24)	70.83% (17/24)
	Llava	16.67% (4/24)	20.83% (5/24)
	MiniCPM-V	20.83% (5/24)	29.17% (7/24)
	Baseline Li et al. (12)	16.89%	—

83.3% for Llama 3.2-Vision. The mid-tier MiniCPM-V follows the same trend, rising from 75.0% to 79.2%. By contrast, Llava remains flat at 58.3%, indicating limited benefit from independent calls when the underlying representation is weak. All four VLMs outperform the traditional CNN baseline reported by (12) (48.6%), underscoring the advantage of prompt-engineered, zero-shot inference for coarse vehicle categorisation.

When the task requires simultaneous make and model recognition, accuracy drops for ev-

ery model, reflecting the added fine-grained visual complexity. Nevertheless, three-call prompting recovers much of this loss: GPT-4o and Llama 3.2-Vision both reach 70.8%, representing 8- and 4-point improvements, respectively, over their single-call results, while MiniCPM-V climbs from 20.8% to 29.2%. Llava again shows only marginal change (16.7%→20.8%). Importantly, the best VLM configurations surpass the non-VLM baseline by more than four-fold (70.8% vs. 16.9%), demonstrating that, even without fine-tuning, contemporary VLMs can deliver state-of-the-art performance on challenging recognition tasks.

Representative examples of both successful and failed predictions are provided in Table 4, illustrating the strengths and limitations of the proposed framework under varying image conditions.

CONCLUSION

This study demonstrates the viability of using vision-language models (VLMs) for robust, end-to-end license plate and vehicle attribute recognition in unconstrained video settings, such as those captured by smartphones. By leveraging a unified pipeline that combines object detection, frame selection through perceptual quality metrics, and prompt-based VLM querying, our method significantly outperforms traditional approaches—achieving up to 91.7% plate recognition accuracy and 70.83% make/model accuracy. Notably, our zero-shot approach eliminates the need for fine-tuning, making it more accessible for deployment on resource-limited devices. The results highlight the potential of open-source VLMs like Llama-3.2-Vision to match proprietary models in accuracy while enabling on-device processing at zero cost, offering a scalable and flexible solution for future intelligent transportation systems.

Limitations

Despite the strong zero-shot performance demonstrated in this study, several limitations remain. First, our pipeline depends critically on the quality and representativeness of the selected frames; extreme motion blur, severe occlusions, or highly unusual viewing angles can still defeat both object detection and VLM prompting, leading to misreads or hallucinations. Second, while off-the-shelf VLMs avoid the need for task-specific training, they can introduce unpredictable errors, especially when interpreting stylized or non-standard plates (e.g., novelty plates, foreign characters) for which the models have little pretraining data. Third, the end-to-end latency and compute requirements, though acceptable on high-end GPUs, may challenge deployment on resource-constrained devices; batching multiple calls to a remote API (e.g., GPT-4o) also incurs variable network latency and potential costs. Finally, our evaluation relies on two datasets (UFPR-ALPR and the 24-video smartphone benchmark) that may not fully capture the global diversity of plate designs, lighting conditions, and camera hardware encountered in real-world citizen enforcement.

Future Directions

To address these challenges, future work could explore several avenues. From a systems perspective, optimizing the pipeline for on-device execution (for example via model quantization, distillation, or pruning) would broaden applicability to smartphones and edge cameras. Expanding our benchmarks to include publicly available dashcam datasets and international plate collections will help quantify generalization. Finally, we aim to equip VLMs with self-validation tools, such as programmatic Google searches or external knowledge APIs, so that when a model’s zero-shot prediction appears implausible (e.g., “Volvo XC60” misread as “Volvo XC90”), it can verify and cor-

rect its output before finalizing the result. Tight integration with enforcement workflows (including automated database lookups, privacy-preserving logging, and human-in-the-loop verification) will also be essential to translate these technical advances into real-world safety and compliance gains.

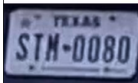
Vehicle	Top-3 plates	Prediction	Ground truth
	  	JVT5339 Toyota RAV4	JVT5339 Toyota RAV4
	  	KVH8370 Ford Fusion	KVH8370 Ford Fusion
	  	PTX1215 Volvo XC90	PTX1215 Volvo XC60
	  	STH0080 Tesla Model 3	STM0080 Tesla Model 3

TABLE 4: Per-vehicle comparison between predicted outputs and ground-truth annotations. The second column presents the top three license plate crops (ranked by quality) provided as input to the model. The first two examples illustrate correct predictions, while the latter two demonstrate failure cases.

REFERENCES

1. World Health Organization, *Road traffic injuries: Key facts*, 2023, accessed: 2025-05-05.
2. National Highway Traffic Safety Administration, *NHTSA Releases 2022 Traffic Deaths, 2023 Early Estimates*, 2024, accessed: 2025-05-05.
3. New York City Independent Budget Office, *Transportation Funds Added for Vision Zero, Traffic Enforcement Cameras*, 2016, accessed: 2025-05-05.
4. Safavi-Naini, S. A. A. et al., Drivers' behavior confronting fixed and point-to-point speed enforcement camera: Agent-based simulation and translation to crash relative risk change. *Scientific Reports*, Vol. 14, 2024, p. Article 52265.
5. Korean National Police Agency, *Smart National Police - Traffic Violation Reporting App*. <https://play.google.com/store/apps/details?id=kr.go.police>, 2023, available on Google Play.
6. Choi, M.-S. and M.-K. Moon, Analysis System for Public Interest Report Video of Traffic Law Violation based on Deep Learning Algorithms. *The Journal of the Korea institute of electronic communication sciences*, Vol. 18, No. 1, 2023, pp. 63–70.
7. Wilson, M., \$87.50 for 3 Minutes: Inside the Hot Market for Videos of Idling Trucks. *The New York Times*, 2022, accessed: 2025-05-09.
8. Salimah, U., V. Maharani, and R. Nursyanti, Automatic License Plate Recognition Using Optical Character Recognition. *IOP Conference Series: Materials Science and Engineering*, Vol. 1115, 2021, p. 012023.
9. Eldharif, E. and T. Fanoush, Automatic License-Plate Recognition Using Optical Character Recognition. *Bulletin of the Natural History Museum Botany*, Vol. 5, 2024.
10. Sugiyono, A. Y., K. Adrio, K. Tanuwijaya, and K. M. Suryaningrum, Extracting Information from Vehicle Registration Plate using OCR Tesseract. *Procedia Computer Science*, Vol. 227, 2023, pp. 932–938, 8th International Conference on Computer Science and Computational Intelligence (ICCSCI 2023).
11. Satar, B. and A. E. Dirik, *Deep Learning Based Vehicle Make-Model Classification*, Springer International Publishing, p. 544–553, 2018.
12. Li, K., J. Malagavalli, L. Goel, T. Wang, and K. Kockelman, Smartphone-based Method for Automated Speed Enforcement. In *2025 Transportation Research Board Annual Meeting*, National Academy of Sciences, 2025.
13. OpenAI, *GPT-4o System Card*, 2024, accessed: 2025-05-05.
14. Meta AI, *LLaMA 3.2: Revolutionizing edge AI and vision with open models*, 2024, accessed: 2025-05-05.
15. Liu, H., C. Li, Q. Wu, and Y. J. Lee, *Visual Instruction Tuning*, 2023, accessed: 2025-05-05.
16. Yao, Y., Z. Liu, M. Sun, and M. Sun, *MiniCPM-V: A GPT-4V Level MLLM on Your Phone*, 2024, accessed: 2025-05-05.
17. Lubna, N. Mufti, and S. A. A. Shah, Automatic Number Plate Recognition: A Detailed Survey of Relevant Algorithms. *Sensors*, Vol. 21, No. 9, 2021.
18. Zherzdev, S. and A. Gruzdev, *LPRNet: License Plate Recognition via Deep Neural Networks*, 2018.
19. Satar, B. and A. E. Dirik, *Deep Learning Based Vehicle Make-Model Classification*, Springer International Publishing, p. 544–553, 2018.

20. Radford, A., J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, *Learning Transferable Visual Models From Natural Language Supervision*, 2021.
21. AlDahoul, N., M. J. T. Tan, R. R. Tera, H. A. Karim, C. H. Lim, M. K. Mishra, and Y. Zaki, *Advancing Vehicle Plate Recognition: Multitasking Visual Language Models with VehiclePaliGemma*, 2024.
22. Wang, J., K. C. K. Chan, and C. C. Loy, *Exploring CLIP for Assessing the Look and Feel of Images*, 2022.
23. Mittal, A., A. K. Moorthy, and A. C. Bovik, No-Reference Image Quality Assessment in the Spatial Domain. *IEEE Transactions on Image Processing*, Vol. 21, No. 12, 2012, pp. 4695–4708.
24. Laroca, R., E. Severo, L. A. Zanlorensi, L. S. Oliveira, G. R. Gonçalves, W. R. Schwartz, and D. Menotti, A Robust Real-Time Automatic License Plate Recognition Based on the YOLO Detector. In *International Joint Conference on Neural Networks (IJCNN)*, 2018, pp. 1–10.