

From Dashcam Videos to Driving Simulations: Stress Testing Automated Vehicles against Rare Events

Yan Miao¹, Georgios Fainekos², Bardh Hoxha², Hideki Okamoto², Danil Prokhorov², Sayan Mitra¹

¹ University of Illinois at Urbana-Champaign, ² Toyota Research Institute North America

Abstract

Testing Automated Driving Systems (ADS) in simulation with realistic driving scenarios is important for verifying their performance. However, converting real-world driving videos into simulation scenarios is a significant challenge due to the complexity of interpreting high-dimensional video data and the time-consuming nature of precise manual scenario reconstruction. In this work, we propose a novel framework that automates the conversion of real-world car crash videos into detailed simulation scenarios for ADS testing. Our approach leverages prompt-engineered Video Language Models (VLM) to transform dashcam footage into SCENIC scripts, which define the environment and driving behaviors in the CARLA simulator, enabling the generation of realistic simulation scenarios. Importantly, rather than solely aiming for one-to-one scenario reconstruction, our framework focuses on capturing the essential driving behaviors from the original video while offering flexibility in parameters such as weather or road conditions to facilitate search-based testing. Additionally, we introduce a similarity metric that helps iteratively refine the generated scenario through feedback by comparing key features of driving behaviors between the real and simulated videos. Our preliminary results demonstrate substantial time efficiency, finishing the real-to-sim conversion in minutes with full automation and no human intervention, while maintaining high fidelity to the original driving events.

1 Introduction

The rapid advancements in Automated Driving Systems (ADS) technology have created an urgent need for robust and realistic testing environments to assure reliability of ADS (Huang et al. 2016). Real-world scenarios, such as crash videos and near-miss events, offer valuable insights into the diverse conditions ADS must navigate, making them an essential source for improving ADS testing. However, replicating these scenarios in real-world settings is both dangerous and impractical. Therefore, converting real-world driving videos into simulation scenarios is a necessary solution, but this process also poses several challenges: the high dimensional video data significantly limits the development of automated methods, while manually reconstructing these scenarios can take experts hours to complete. This time-consuming process underscores the need for a more efficient and automated approach.

In this paper, we propose a novel framework that automates the conversion of real-world driving videos into detailed simulation scenarios. Our approach leverages prompt-engineered Video Language Models (VLMs) to transform dashcam videos into SCENIC scripts, enabling the automatic generation of realistic simulations in CARLA. In addition, rather than aiming solely to create an exact digital copy of the dashcam video, our framework prioritizes capturing the most essential driving behaviors while maintaining flexibility in parameters such as weather and road conditions. This flexibility supports search-based testing, a methodology that systematically varies parameters to identify edge cases and vulnerabilities in ADS. Furthermore, we introduce a similarity metric that iteratively refines the generated simulations by comparing key driving features between the real and simulated videos. The contributions of this work are fourfold: (1) an automated video-to-simulation pipeline that removes the need for manual scenario construction, (2) a similarity metric that bridge the gap between real and simulated scenarios, (3) an iterative feedback loop for scenario refinement using neural network-based feedback from VLM to capture core driving behaviors, and (4) a significant improvement in time efficiency, reducing scenario generation from hours to minutes while maintaining high fidelity to the original events.

2 Related Work

Testing ADS in simulation environments has been the subject of extensive research (Huang et al. 2016). Existing autonomous vehicle simulation platforms such as CARLA (Dosovitskiy et al. 2017) and LGSVL (Rong et al. 2020) allow researchers to generate and manipulate driving scenarios in controlled environments. Efforts such as the SCENIC language (Fremont et al. 2019) enables search-based testing (SBT) so that the generated scenario can be a seed for search based testing with respect to temporal logic requirements (Dreossi et al. 2019; Tuncali et al. 2020), safe driving rules (Hekmatnejad, Hoxha, and Fainekos 2020) and traffic laws (Sun et al. 2022). However, accurately designing these scenarios to closely resemble real-world events remains a challenge.

Recently, efforts have shifted toward automatic real-to-simulation (real-to-sim) conversion approaches that use video data to guide scenario generation. For instance, Bai

et al. (Bai et al. 2024) introduced a system that extracts key information from videos to create driving scenarios in simulation environments. In (Elmaaroufi et al. 2024), the authors automatically generate driving scenarios from police crash reports. Similarly, (Wang et al. 2023) focus on learning realistic human behaviors in real-life scenarios and use learned models to improve simulations. NVIDIA’s STRIVE (Rempe et al. 2022) generates accident-prone driving scenarios by modifying 2D trajectories, but this method is based on controlled scenarios rather than real-world crash videos. Another approach, DEEPCRASHTTEST (Bashetty, Ben Amor, and Fainekos 2020), converts dashcam footage into crash tests by extracting 3D vehicle trajectories but lacks an iterative refinement process to improve simulation accuracy.

While these approaches represent meaningful progress, existing methods are often limited by a reliance on pre-defined trajectories or fail to incorporate iterative feedback to refine the generated scenarios. Furthermore, their primary goal has been to ensure that the generated scenarios are as identical as possible to the original ones (e.g. optimizing on the tracking performance for KITTI tracking dataset (Fritsch, Kuehnl, and Geiger 2013)). In contrast, we want to focus on capturing the core driving behaviors from the original video while introducing flexibility in other parameters, such as weather and road conditions. This flexibility enables search-based testing, allowing for systematic exploration of environmental variations to identify vulnerabilities in ADS and improve their robustness.

To address these gaps, we chose to use VLMs as they enable more flexible and scalable video-to-language translation. Recent works, such as DriveDreamer-2 (Zhao et al. 2024), employ LLMs to generate user-defined driving videos by translating queries into agent trajectories and maps for simulation. Similarly, Text-to-Drive (Nguyen et al. 2024) leverages knowledge-driven language descriptions to synthesize diverse driving behaviors. While these works demonstrate the potential of foundational models in simulation-based driving applications, they either primarily limit to synthetic data and do not leverage real-world driving videos or focus on user-defined queries.

In contrast, our framework integrates real-world driving data, such as dashcam crash videos, to generate simulation scenarios that not only capture critical driving behaviors but also provide flexibility for environmental parameters like road type and weather. Through few-shot prompt engineering, we tailor VLMs to create SCENIC scripts for scenario generation, bridging the gap between exact scenario replication and the need for diversity in simulation environments. This approach ensures realistic and behaviorally accurate test cases while supporting search-based testing for robust ADS validation.

3 Real-to-Sim Scenario Generation Framework

Given a real-life vehicle crash video (e.g., a dash camera recording), our objective is to generate corresponding simulation scenarios that accurately capture the core driving behaviors. In this paper, we define scenario as a driving

event that includes environmental conditions (e.g., weather, road types), vehicle dynamics (e.g., speed, acceleration), and driving behaviors between entities (e.g., overtaking behavior or collisions). This approach enables robust and safe testing of Automated Driving Systems (ADS) by leveraging real-life close-call situations within a controlled simulation environment. Our framework consists of 4 components: (1) conversion of real-world video into SCENIC scripts, (2) generation of simulation videos from SCENIC scripts, (3) similarity analysis between the real and simulated videos, and (4) iterative refinement to ensure the simulated video’s consistency with the original scenario.

3.1 Video-to-Text Generation

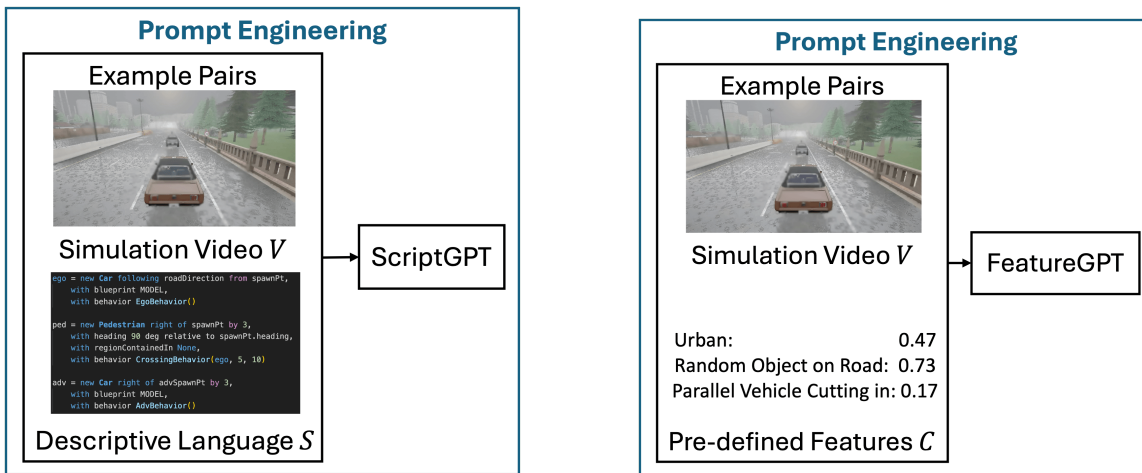
First, we convert the input video into a descriptive script called SCENIC. SCENIC is a probabilistic programming language designed for describing driving scenarios with a high level of precision, enabling detailed specification of environmental factors, road layouts, and vehicle behaviors for use in simulation platforms like CARLA. This conversion is accomplished using a prompt-engineered version of the pre-trained GPT-4o model.

Prompt engineering involves designing and refining input prompts to guide large foundational models, such as GPT-4o (OpenAI 2024), toward producing accurate and desired outputs. In this case, prompt engineering is especially effective for generating detailed scenario descriptions (e.g., SCENIC scripts) from real-world driving videos. During the prompt engineering process, we improve the pre-trained GPT-4o model by providing it with multiple “positive-example” pairs (V_i, S_i) , where V_i is a simulation video generated in CARLA using the corresponding SCENIC script S_i that describes the scenario, as illustrated in Figure 1a. Through this process, the new Video-Language Model, referred to as *ScriptGPT*, learns to map key visual elements, such as weather, road conditions, and vehicle behaviors, into structured and accurate scenario description languages. After sufficient prompt engineering, ScriptGPT is capable of generating a SCENIC script S_{out} for real-world crash video V_{real} .

Although the initial prompt engineering process was conducted using simulation videos paired with their SCENIC scripts, through empirical testing, we found this approach could be extendable to real-world driving videos because the underlying visual and descriptive patterns (e.g., road layouts, traffic behaviors, and environmental factors) are consistent across both domains, allowing the model to effectively generalize its learned capabilities and accurately capture real-world scenarios.

3.2 Text-to-Video Generation

We convert the descriptive language S_{out} to simulation video V_{sim} using SCENIC (Fremont et al. 2018, 2019). SCENIC is a programming language tool designed for specifying driving scenarios through environmental factors, vehicle behaviors, and road conditions. Once we have the SCENIC script S_{out} generated from the real crash video using the prompted-engineered model, we feed it into the SCENIC framework to synthesize a simulation scenario in CARLA, as shown in Figure 2. The script S_{out} serves as



(a) **ScriptGPT**: A Video Language Model derived from GPT-4o, developed through prompt engineering with paired examples of simulation videos and their corresponding descriptive SCENIC scripts.

(b) **FeatureGPT**: A Video Language Model derived from GPT-4o, developed through prompt engineering with paired examples of simulation videos and their corresponding pre-defined features.

Figure 1: Train *ScriptGPT* and *FeatureGPT* using Prompt Engineering

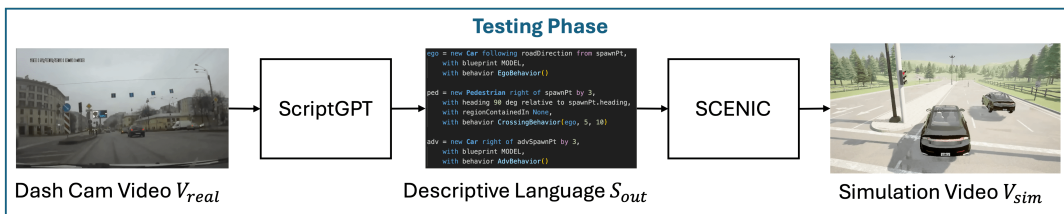


Figure 2: After prompt engineering, the dash cam video is fed into *ScriptGPT*, which synthesizes descriptive language in the SCENIC format. This SCENIC script can then be executed in CARLA to generate a corresponding testing scenario in simulation.

the textual representation of the scene, encoding environmental conditions (e.g., weather, traffic), vehicle dynamics, and road types. The output is a new simulation video V_{sim} , which visually represents the scenario described in S_{out} .

3.3 Similarity Check

Next, we perform a similarity check on the simulated scenario V_{sim} and the original V_{real} , measured on a set of pre-defined features. Ideally, the simulated and original scenario should closely match, allowing us to seamlessly replace the ego vehicle with any ADS (e.g. Baidu’s Apollo planner (Fan et al. 2018) and controller) for testing. However, due to the complexity of real-world scenarios, even with prompt engineering, discrepancies often arise, where the generated simulation may miss key features or introduce extra ones.

To address this, we introduce a similarity metric to ensure that the generated video captures the most important features from the original video. We predefine a set of crucial feature categories, such as the most critical and frequently encountered driving behaviors and environmental conditions, to examine the original crash video V_{real} with the generated simulation V_{sim} . Then, we use another prompt-engineered transformer model, *FeatureGPT* (depicted in Figure 1b), to output a predicted probability (from 0 to 1) for each pre-

defined feature category for a given video. Next, the similarity score vector, $Sim(V_{real}, V_{sim})$, is calculated as a vector of differences between the predicted probabilities across the predefined categories:

$$Sim(V_{real}, V_{sim}) = [C_{real_1} - C_{sim_1}, C_{real_2} - C_{sim_2}, \dots, C_{real_n} - C_{sim_n}],$$

where C_{real_i} and C_{sim_i} represent the predicted probabilities of V_{real} and V_{sim} for feature i .

3.4 Iterative Refinement

Finally, we perform iterative refinement on the *ScriptGPT* with the help of similarity check, as illustrated in Figure 3. If the absolute difference for a given category i is above a predefined threshold τ_i , we refine the SCENIC script S_{out} by feeding the discrepancy for that category back into *ScriptGPT* as an additional prompt (e.g., "there shouldn’t be a leading vehicle overtaking behavior, please improve on that"). This feedback allows *ScriptGPT* to adjust the SCENIC script S_{out} accordingly, generating a new version of the script and producing a new simulation video V_{sim}' .

Although a single SCENIC script can produce multiple simulation videos with consistent core driving behaviors but

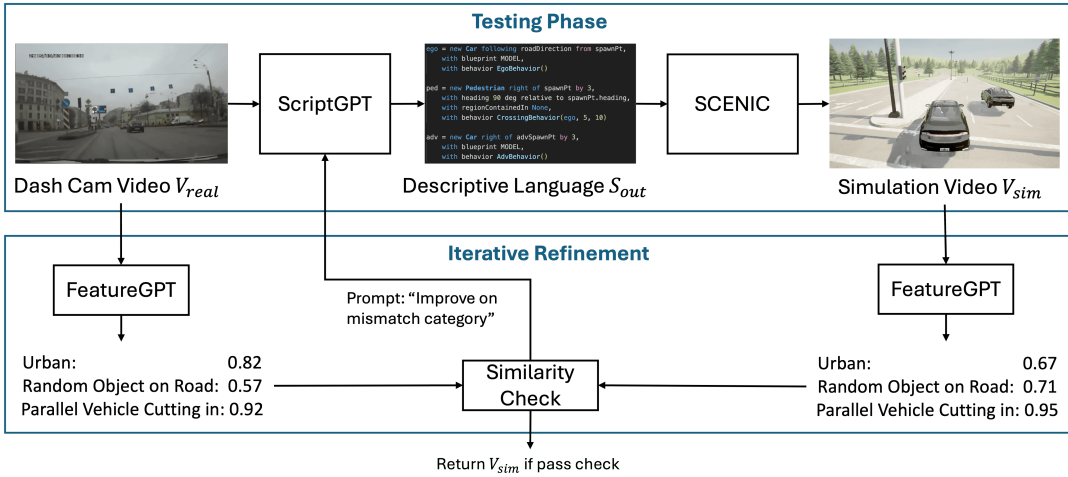


Figure 3: **Iterative Refinement Process**: After obtaining the simulated video V_{sim} from *ScriptGPT* and SCENIC, both the original and simulated videos are fed into *FeatureGPT* to evaluate the probabilities of predefined features. If the difference of any feature between the original and simulated videos exceeds a certain threshold, we iteratively refine *ScriptGPT* by incorporating additional feedback into the SCENIC script, guiding further scenario adjustments until the similarity improves.

variations in other parameters (e.g., road types or vehicle configurations), we utilize only a single video for similarity checks and iterative refinement, because multiple simulation videos has the same driving behaviors and refining on a video is sufficient to ensure that the generated scenario aligns with the essential driving behaviors of the original video. This approach provides the necessary focus for refining critical features while maintaining flexibility for generating diverse scenarios, enabling search-based testing.

The process is repeated iteratively until the difference for each category falls below the predefined threshold, i.e., $\|Sim(V_{real}, V_{sim})_i\| \leq \tau_i$. Once the similarity across all categories pass the threshold check, the final simulation scenario is ready for testing ADS.

4 Experiments & Analysis

4.1 System Setup

Dataset We obtain the collision videos from the Car Crash Dataset (CCD) (Bao, Yu, and Kong 2020). CCD is chosen because it contains real traffic accident videos captured by dashcams mounted on driving vehicles, potentially providing a rich source for developing and testing ADS. Moreover, our framework is also relevant for near misses events.

Simulator We use CARLA (Dosovitskiy et al. 2017), an open-source platform designed to support the development and validation of autonomous driving systems. CARLA is selected for its realistic physics engine and high-fidelity environmental rendering, making it ideal for generating the simulation scenarios needed in our framework.

Description Language SCENIC We use SCENIC as the description language for driving scenarios due to its similarity to Python, which aligns well with GPT-4o’s input data (OpenAI 2024), making prompt engineering doable and more straightforward. Furthermore, SCENIC integrates

seamlessly with the CARLA simulator and it is highly effective at specifying complex driving scenarios, making it well-suited for our application.

Prompt-Engineered *ScriptGPT* To construct the *ScriptGPT* model, we begin by carefully and manually writing 20 SCENIC scripts that cover diverse driving scenarios, such as overtaking, cruising, sudden stops due to obstacles, and turns in varying road and weather conditions. Using the SCENIC library and CARLA, we generate corresponding videos for each scenario and pair them with their respective SCENIC scripts, forming 20 (V_i, S_i) pairs. These pairs are then used as input data for prompt engineering GPT-4o, selected for its ability to learn and generalize from a wide range of examples.

The empirical design choice we made is to design only 20 pairs because through empirical testing, 20 scenarios are sufficient to cover most of the common interesting scenarios. We believe we only need a very small number of prompt engineering due to the powerful generalization ability of pre-trained GPT-4o. Moreover, in practice, since GPT-4o API does not natively accept video formats like .mp4, we preprocess the videos by sampling frames and concatenating them into an n-dimensional array. The same preprocessing technique is applied during testing phases to ensure consistency between training and testing phases, and we also apply the same to the *FeatureGPT*.

Prompt-Engineered *FeatureGPT* We use *FeatureGPT* to enhance our framework’s ability to recognize specific driving behaviors. First, we predefine 10 driving feature categories, as shown in Table 1. Next, we create 20 SCENIC scripts representing scenario videos where these features may or may not appear. Each video is paired with a corresponding 10-dimensional feature vector (e.g., [parallel vehicle overtaking: 0, ..., leading vehicle stopped: 1]), where 0 indicates absence and 1 indicates presence). The input

data is used to prompt-engineer GPT-4o into *FeatureGPT*, which outputs a 10-dimensional probability vector for each video during inference. This allows us to compare and categorize driving behaviors between the original video V_{real} and the generated video V_{sim} . Again we made the empirical design choice of using 10 feature categories and 20 samples to cover most frequently encountered interesting driving behaviors through trial and error.

Iterative Refinement using Similarity Check Once *FeatureGPT* produces feature vectors for both the generated and original videos, we compare them to detect discrepancies. A large discrepancy in any feature indicates that the generated video is either missing a key behavior (negative gap) or introducing an unintended one (positive gap), and then we map the gap into natural language feedback (e.g., "there should be a leading vehicle overtaking behavior, please improve on that") and feed back into *ScriptGPT* to refine the SCENIC script. Gap thresholds τ_i for each feature i , as shown in Table 1, are customized through empirical testing. We assign slightly larger thresholds to weather and road type features compared to other driving behavior features, as these are considered less critical for the testing objectives. This leniency allows for greater variety in environmental parameters, supporting the generation of diverse scenarios while ensuring that core driving behaviors remain accurately captured.

Table 1: Pre-defined feature Category and Threshold

Pre-define Feature	Gap Threshold τ
Sunny / Rainy	0.3
Urban / Highway	0.3
Random Object on Road	0.2
Leading Vehicle Cruising	0.2
Leading Vehicle Stopped	0.2
Parallel Vehicle Cutting in	0.2
Parallel Vehicle Cruising	0.2
Parallel Vehicle Stopped	0.2
Behind Vehicle Overtaking	0.2
Opposite Vehicle Turning	0.2

4.2 Case Study

We present five interesting dashcam videos from the CCD dataset and generate their corresponding simulation scenarios using our framework, as shown in Figure 4, 5, 6, 7, 8. Importantly, each SCENIC script can correspond to multiple simulation scenarios. For example, a highway scenario may result in a 2-lane road or a 3-lane road, and a leading vehicle could be represented as either a truck or a sedan. This variability enables search-based testing by generating diverse scenarios that explore different environmental parameters while preserving the core driving behaviors.

In this section, we showcase one simulation scenario per SCENIC script due to page limitations. However, our framework inherently supports generating multiple testing scenarios from a single script, allowing systematic exploration of

predefined features. As demonstrated in the figures, the generated scenarios faithfully capture the most essential driving behaviors observed in the dashcam videos while introducing controlled variations in other parameters, such as road configuration and vehicle types, to support robust ADS testing.



Figure 4: **Vehicle Cutting In with Pedestrian Crossing Scenario**: in the original dash camera video (top row), the vehicle on the right performs an emergency lane change to the left due to a jaywalking pedestrian in red. In the generated scenario (bottom row) produced by our framework, the vehicle on the right exhibited a similar lane change behavior to the left to avoid a jaywalking pedestrian.



Figure 5: **Opposite Vehicle Invading Lane Scenario**: in the original dash camera video (top row), the vehicle on the opposite lane gradually swifts to ego's lane probably due to loss of focus. In the generated scenario (bottom row) produced by our framework, the vehicle on the opposite lane exhibited a similar lane change behavior to switch to our lane and caused collision.



Figure 6: **Vehicle Spin Scenario**: in the original dash camera video (top row), the vehicle in front of the ego first spins to the left and then collided into the right vehicle. In the generated scenario (bottom row) produced by our framework, the front vehicle exhibited a similar spin and collision behavior



Figure 7: **Animal Crossing Scenario**: in the original dash camera video (top row), an animal attempted to cross the road, prompting the ego vehicle to perform an emergency lane change to the left. In the generated scenario (bottom row) produced by our framework, the ego vehicle exhibited a similar lane change behavior to the left to avoid a jaywalking pedestrian (since CARLA does not have an animal model, the animal is replaced by a pedestrian).

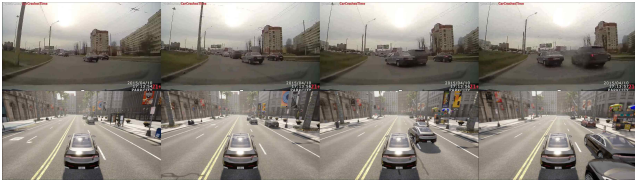


Figure 8: **Vehicle Cutting In with Stopped Object Scenario**: in the original dash camera video (top row), the vehicle on the right perform an emergency brake and tried lane change to the left due to a front parked vehicle. In the generated scenario (bottom row) produced by our framework, the vehicle on the right exhibited a similar lane change behavior to the left to avoid the stopped vehicle and cause a collision.

4.3 Preliminary Qualitative Result

Automated Pipeline After completing the prompt engineering process, our framework is capable of automatically generating the 5 scenarios mentioned in Section 4.2 without any human intervention during the testing phase. We also extended our framework to evaluate 50 randomly selected accidents from the CCD dataset, and found that 32 out of 50 (64%) scenarios generation can be fully automated without any human involvement, while 18 scenarios (36%) encountered errors (i.e. generating a keyword or object that is not supported by SCENIC syntax), thus cannot ran automatically.

Time Efficiency Our framework significantly reduces the time required for real-to-simulation scenario generation. For the 32 videos that can be automatically generated, it takes 1.5 minutes per scenario on average during the testing phase, including iterative refinement, to produce a SCENIC script of approximately 70 lines of code. In contrast, manually coding and debugging a similar real-to-simulation scenario could take an experienced engineer several hours. Even for the 36% scenarios that cannot automatically finish, it only takes minutes for human to fix the syntax error to make the pipeline work again.

Advantage of Iterative Refinement The 5 scenarios in Section 4.2 underwent 1-2 iterations of refinement, resulting in notable improvements in both accuracy and realism.

In the extended evaluation of 50 accidents, iterative refinement happened in 17 scenarios (34%), further showcasing its potential to enhance scenario quality and be generalized to more crash scenarios.

Framework Accuracy While we have not yet conducted a formal quantitative analysis on the framework’s timing efficiency or objectively measure benefits of iterative refinement, the preliminary results provide strong evidence of the concept’s validity. Our framework consistently captures the core driving behaviors from the original videos, indicating its effectiveness in generating accurate and realistic driving scenarios.

5 Limitations & Future Work



Figure 9: Scenarios from the Car Crash Dataset where our current framework cannot handle due to heavy traffic with multiple cars or compromised lighting environmental conditions at night

The current framework faces several challenges. It struggles with scenarios involving poor perception conditions or high complexity, as illustrated in Figure 9. Additionally, the framework’s performance may degrade when testing scenarios (e.g. multi-car complicated driving scenario at night like shown in Figure 9) have not been encountered by *ScriptGPT* and *FeatureGPT* during the prompt engineering process. Therefore, the performance of our framework heavily relies on carefully and manually selecting diverse “positive-examples” for prompt engineering.

For future work, we plan to conduct a human study to quantitatively assess the framework’s time efficiency and accuracy. This will involve timing how long experts take to manually write real-to-simulation conversion scripts and asking them to rate the accuracy of our automated conversions. We also aim to improve the diversity of data samples used in the prompt engineering process to extend our framework’s applicability across the entire Car Crash Dataset. Finally, we envision extending this video-to-video conversion framework to other domains, such as flying tasks or other robotics applications, broadening its use cases and potential impact.

6 Conclusion

In this paper, we have presented a novel framework for automatically converting real-world vehicle crash videos into simulation scenarios using prompt-engineered Video-Language Models. We have deploy multiple techniques including the similarity score vector metric and the iterative refinement process to ensure the generated scenarios

closely align with the original videos. Despite the framework's current limitations, such as reliance on data diversity and the challenges of handling complex or unseen scenarios, through multiple examples, it demonstrates clear potential for improving ADS testing. Future work will focus on expanding data diversity, conducting quantitative human studies, and extending the framework to other robotics domains.

References

- Bai, X.; Luo, Y.; Jiang, L.; Gupta, A.; Kaveti, P.; Singh, H.; and Ostadabbas, S. 2024. Bridging the Domain Gap between Synthetic and Real-World Data for Autonomous Driving. *ACM J. Auton. Transport. Syst.*, 1(2).
- Bao, W.; Yu, Q.; and Kong, Y. 2020. Uncertainty-based Traffic Accident Anticipation with Spatio-Temporal Relational Learning. In *ACM Multimedia Conference*.
- Bashetty, S. K.; Ben Amor, H.; and Fainekos, G. 2020. DeepCrashTest: Turning Dashcam Videos into Virtual Crash Tests for Automated Driving Systems. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 11353–11360.
- Dosovitskiy, A.; Ros, G.; Codevilla, F.; López, A. M.; and Koltun, V. 2017. CARLA: An Open Urban Driving Simulator. *CoRR*, abs/1711.03938.
- Dreossi, T.; Fremont, D. J.; Ghosh, S.; Kim, E.; Ravanbakhsh, H.; Vazquez-Chanlatte, M.; and Seshia, S. A. 2019. Verifai: A toolkit for the formal design and analysis of artificial intelligence-based systems. In *International Conference on Computer Aided Verification*, volume 11561 of *LNCS*, 432–442. Springer.
- Elmaaroufi, K.; Shanker, D.; Cismaru, A.; Vazquez-Chanlatte, M.; Sangiovanni-Vincentelli, A.; Zaharia, M.; and Seshia, S. A. 2024. Generating Probabilistic Scenario Programs from Natural Language. arXiv:2405.03709.
- Fan, H.; Zhu, F.; Liu, C.; Zhang, L.; Zhuang, L.; Li, D.; Zhu, W.; Hu, J.; Li, H.; and Kong, Q. 2018. Baidu Apollo EM Motion Planner. *CoRR*, abs/1807.08048.
- Fremont, D. J.; Dreossi, T.; Ghosh, S.; Yue, X.; Sangiovanni-Vincentelli, A. L.; and Seshia, S. A. 2019. Scenic: A Language for Scenario Specification and Scene Generation. In *Proceedings of the 40th annual ACM SIGPLAN conference on Programming Language Design and Implementation (PLDI)*.
- Fremont, D. J.; Yue, X.; Dreossi, T.; Ghosh, S.; Sangiovanni-Vincentelli, A. L.; and Seshia, S. A. 2018. Scenic: Language-Based Scene Generation. *CoRR*, abs/1809.09310.
- Fritsch, J.; Kuehnl, T.; and Geiger, A. 2013. A New Performance Measure and Evaluation Benchmark for Road Detection Algorithms. In *International Conference on Intelligent Transportation Systems (ITSC)*.
- Hekmatnejad, M.; Hoxha, B.; and Fainekos, G. 2020. Search-based Test-Case Generation by Monitoring Responsibility Safety Rules. In *IEEE Intelligent Transportation Systems Conference (ITSC)*.
- Huang, W.; Wang, K.; Lv, Y.; and Zhu, F. 2016. Autonomous vehicles testing methods review. In *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, 163–168.
- Nguyen, P.; Wang, T.-H.; Hong, Z.-W.; Karaman, S.; and Rus, D. 2024. Text-to-Drive: Diverse Driving Behavior Synthesis via Large Language Models. arXiv:2406.04300.
- OpenAI. 2024. GPT-4 Technical Report. arXiv:2303.08774.
- Rempe, D.; Phillion, J.; Guibas, L. J.; Fidler, S.; and Litany, O. 2022. Generating Useful Accident-Prone Driving Scenarios via a Learned Traffic Prior. arXiv:2112.05077.
- Rong, G.; Shin, B. H.; Tabatabaee, H.; Lu, Q.; Lemke, S.; Možeiko, M.; Boise, E.; Uhm, G.; Gerow, M.; Mehta, S.; Agafonov, E.; Kim, T. H.; Sterner, E.; Ushiroda, K.; Reyes, M.; Zelenkovsky, D.; and Kim, S. 2020. LGSVL Simulator: A High Fidelity Simulator for Autonomous Driving. In *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, 1–6.
- Sun, Y.; Poskitt, C. M.; Sun, J.; Chen, Y.; and Yang, Z. 2022. LawBreaker: An Approach for Specifying Traffic Laws and Fuzzing Autonomous Vehicles. In *Automated Software Engineering (ASE)*.
- Tuncali, C. E.; Fainekos, G.; Prokhorov, D.; Ito, H.; and Kapinski, J. 2020. Requirements-driven Test Generation for Autonomous Vehicles with Machine Learning Components. *IEEE Transactions on Intelligent Vehicles*, 5: 265–280.
- Wang, J.; Yuan, Y.; Luo, Z.; Xie, K.; Lin, D.; Iqbal, U.; Fidler, S.; and Khamis, S. 2023. Learning Human Dynamics in Autonomous Driving Scenarios. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 20739–20749. Los Alamitos, CA, USA: IEEE Computer Society.
- Zhao, G.; Wang, X.; Zhu, Z.; Chen, X.; Huang, G.; Bao, X.; and Wang, X. 2024. DriveDreamer-2: LLM-Enhanced World Models for Diverse Driving Video Generation. arXiv:2403.06845.